

Languages recognised by finite
semigroups, and their generalisations to
objects such as trees and graphs, with an
emphasis on definability in monadic
second-order logic

Mikołaj Bojańczyk

June 16, 2020

The latest version can be downloaded from:

<https://www.mimuw.edu.pl/bojan/2019-2020/algebraic-language-theory-2020>

Contents

| | |
|--|---------------|
| <i>Preface</i> | <i>page v</i> |
| Part I Words | 1 |
| 1 Semigroups, monoids and their structure | 3 |
| 1.1 Recognising languages | 8 |
| 1.2 Green's relations and the structure of finite semigroups | 11 |
| 1.3 The Factorisation Forest Theorem | 19 |
| 2 Logics on finite words, and the corresponding monoids | 29 |
| 2.1 All monoids and monadic second-order logic | 30 |
| 2.2 Aperiodic semigroups and first-order logic | 36 |
| 2.3 Suffix trivial semigroups and temporal logic with F only | 48 |
| 2.4 Infix trivial semigroups and piecewise testable languages | 52 |
| 2.5 Two-variable first-order logic | 58 |
| 3 Infinite words | 64 |
| 3.1 Determinisation of Büchi automata for ω -words | 64 |
| 3.2 Countable words and \circ -semigroups | 74 |
| 3.3 Finite representation of \circ -semigroups | 81 |
| 3.4 From \circ -semigroups to MSO | 91 |
| Part II Monads | 101 |
| 4 Monads | 103 |
| 4.1 Monads and their Eilenberg-Moore algebras | 104 |
| 4.2 A zillion examples | 113 |
| 4.3 Syntactic algebras | 122 |
| 4.4 The Eilenberg Variety Theorem | 134 |

| | | |
|----------|---|-----|
| | Part III Trees and graphs | 143 |
| 5 | Forest algebra | 145 |
| | 5.1 The forest monad | 145 |
| | 5.2 Recognisable languages | 147 |
| | 5.3 Logics for forest algebra | 156 |
| 6 | Hypergraphs of bounded treewidth | 172 |
| | 6.1 Graphs, logic, and treewidth | 172 |
| | 6.2 The hypergraph monad | 179 |
| | 6.3 Definable tree decompositions | 197 |
| | <i>Bibliography</i> | 225 |
| | <i>Author index</i> | 229 |
| | <i>Subject index</i> | 230 |

Preface

These are lecture notes on the algebraic approach to regular languages. The classical algebraic approach is for finite words; it uses semigroups instead of automata. However, the algebraic approach can be extended to structures beyond words, e.g. infinite words, or trees or graphs.

PART ONE

WORDS

1

Semigroups, monoids and their structure

In this chapter, we define semigroups and monoids, and show how they can be used to recognise languages of finite words.

Definition 1.1 (Semigroup). A *semigroup* consists of an underlying set S together with a binary product operation

$$(a, b) \mapsto ab,$$

which is associative in the sense that

$$a(bc) = (ab)c \quad \text{for all } a, b, c \in S.$$

The definition says that the order of evaluation in a semigroup is not important, i.e. that different ways of bracketing a sequence of elements in the semigroup will yield the same result as far as the semigroup product is concerned. For example,

$$((ab)c)(d(e f)) = (((ab)c)d)e)f.$$

Therefore, it makes sense to omit the brackets and write simply

$$abcdef.$$

This means that the product operation in the semigroup can be seen as an operation of type $S^+ \rightarrow S$, i.e. it is defined not just on pairs of semigroup elements, but also on finite nonempty words consisting of semigroup elements.

A *semigroup homomorphism* is a function between semigroups that preserves the structure of semigroups, i.e. a function

$$h : \underbrace{S}_{\text{semigroup}} \rightarrow \underbrace{T}_{\text{semigroup}}$$

which is consistent with the product operation in the sense that

$$h(a \cdot b) = h(a) \cdot h(b),$$

where the semigroup product on the left is in S , and the semigroup product on the right is in T . An equivalent definition of a semigroup homomorphism is requiring the following diagram to commute:

$$\begin{array}{ccc} S^+ & \xrightarrow{h^+} & T^+ \\ \text{product in } S \downarrow & & \downarrow \text{product in } T \\ S & \xrightarrow{h} & T \end{array}$$

In the above, h^+ is the natural lifting of h to words.

A *monoid* is the special case of a semigroup where there is an identity element, denoted by $1 \in S$, which satisfies

$$1a = a1 \quad \text{for all } a \in S.$$

The identity element, if it exists, must be unique. This is because if there are two candidates for the identity, then taking their product reveals the true identity. The product operation in a monoid can be thought of as having type $S^* \rightarrow S$, since with the empty word ε being mapped to 1. A *monoid homomorphism* is a semigroup homomorphism that preserves the identity element. In terms of commuting diagrams, a monoid homomorphism is a function which makes the following diagram commute:

$$\begin{array}{ccc} S^* & \xrightarrow{h^*} & T^* \\ \text{product in } S \downarrow & & \downarrow \text{product in } T \\ S & \xrightarrow{h} & T \end{array}$$

Clearly there is a pattern behind the diagrams. This pattern will be explored in the second part of this book, when talking about monads.

Example 1.2. Here are some examples of monoids and semigroups.

- (1) If Σ is a set, then the set Σ^+ of nonempty words over Σ , equipped with concatenation, is a semigroup, called the *free¹ semigroup over generators Σ* . The *free monoid* is the set Σ^* of possibly empty words.

¹ The reason for this name is the following universality property. The free semigroup is generated by Σ , and it is the biggest semigroup generated by Σ in the following sense. For every semigroup S that is generated by Σ , there exists a (unique) surjective semigroup homomorphism $h : \Sigma^+ \rightarrow S$ which is the identity on the Σ generators.

- (2) Every group is a monoid.
 (3) For every set Q , the set of all functions $Q \rightarrow Q$, equipped with function composition, is a monoid. The monoid identity is the identity function.
 (4) For every set Q , the set of all binary relations on Q is a monoid, when equipped with relational composition

$$a \circ b = \{(p, q) : \text{there is some } r \in Q \text{ such that } (p, r) \in a \text{ and } (r, q) \in b\}.$$

The monoid identity is the identity function. The monoid from the previous item is a sub-monoid of this one, i.e. the inclusion map is a monoid homomorphism.

- (5) Here are all semigroups of size two, up to semigroup isomorphism:

$$\underbrace{(\{0, 1\}, +)}_{\text{addition mod 2}} \quad (\{0, 1\}, \min) \quad \underbrace{(\{0, 1\}, \pi_1)}_{\text{product } ab \text{ is } a} \quad \underbrace{(\{0, 1\}, \pi_2)}_{\text{product } ab \text{ is } b} \quad \underbrace{(\{0, 1\}, (a, b) \mapsto 1)}_{\text{all products are 1}}$$

The first two are monoids.

Compositional functions. Semigroup homomorphisms are closely related with functions that are compositional in the sense defined below. Let S be a semigroup, and let X be a set (without a semigroup structure). A function

$$h : S \rightarrow X$$

is called *compositional* if for every $a, b \in S$, the value $h(a \cdot b)$ is uniquely determined by the values $h(a)$ and $h(b)$. If X has a semigroup structure, then every semigroup homomorphism $S \rightarrow X$ is a compositional function. The following lemma shows that the converse is also true for surjective functions.

Lemma 1.3. *Let S be a semigroup, let X be a set, and let $h : S \rightarrow X$ be a surjective compositional function. Then there exists (a unique) semigroup structure on X which makes h into a semigroup homomorphism.*

Proof Saying that $h(a \cdot b)$ is uniquely determined by $h(a)$ and $h(b)$, as in the definition of compositionality, means that there is a binary operation \circ on X , which is not yet known to be associative, that satisfies

$$h(a \cdot b) = h(a) \circ h(b) \quad \text{for all } a, b \in S. \quad (1.1)$$

The semigroup structure on X uses \circ as the semigroup operation. It remains to prove associativity of \circ . Consider three elements of X , which can be written as $h(a), h(b), h(c)$ thanks to the assumption on surjectivity of h . We have

$$(h(a) \circ h(b)) \circ h(c) \stackrel{(1.1)}{=} (h(ab)) \circ h(c) \stackrel{(1.1)}{=} h(abc).$$

The same reasoning shows that $h(a) \circ (h(b) \circ h(c))$ is equal to $h(abc)$, thus establishing associativity. \square

Commuting diagrams. We finish this section with an alternative description of semigroups which uses commuting diagrams. We include this description, because similar descriptions will be frequently used in this book, e.g. for generalisations of semigroups for infinite words, so we want to start using it as early as possible.

The binary product operation in a semigroup S can be extended to a general operation of type $S^+ \rightarrow S$. The following lemma explains, using commuting diagrams, which operations arise this way.

Lemma 1.4. *An operation $\pi : S^+ \rightarrow S$ arises from some semigroup operation on S if and only if the following two diagrams commute:*

$$\begin{array}{ccc}
 S & & \\
 \text{view a letter as} & \searrow \text{identity} & \\
 \text{a one-letter word} & \downarrow & \\
 S^+ & \xrightarrow{\pi} & S
 \end{array}
 \qquad
 \begin{array}{ccc}
 (S^+)^+ & \xrightarrow{\text{product in free semigroup } S^+} & S^+ \\
 \pi^+ \downarrow & & \downarrow \pi \\
 S^+ & \xrightarrow{\pi} & S
 \end{array}$$

In the above, π^+ stands for the coordinate-wise lifting of π to words of words.

For monoids, the same lemma holds, with $+$ replaced by $*$. There is no need to add an extra diagram for the monoid identity, since the monoid identity can be defined as the image under π of the empty word ε . The axioms $1 \cdot a$ and $a \cdot 1$ then follow from

$$1 \cdot a = \pi(\varepsilon) \cdot \pi(a) = \pi(\varepsilon a) = \pi(a) = a,$$

and a symmetric reasoning for $a \cdot 1$.

Also homomorphisms can be defined using commuting diagrams. A function $h : S \rightarrow T$ is a semigroup homomorphism if and only if the following diagram commutes

$$\begin{array}{ccc}
 S^+ & \xrightarrow{h^+} & T^+ \\
 \text{product in } S \downarrow & & \downarrow \text{product in } T \\
 S & \xrightarrow{h} & T
 \end{array}$$

By replacing $+$ with $*$ we get the definition of a monoid homomorphism.

Exercises

Exercise 1. (1) Show a function between two monoids that is a semigroup homomorphism, but not a monoid homomorphism.

Exercise 2. (1) Show that there are exponentially many semigroups of size n .

Exercise 3. Show that for every semigroup homomorphism $h : \Sigma^+ \rightarrow S$, with S finite, there exists some $N \in \{1, 2, \dots\}$ such that every word of length at least N can be factorised as $w = w_1 w_2 w_3$ where $h(w_2)$ is an idempotent².

Exercise 4. (1) Show that if S is a semigroup, then the same is true for the powerset semigroup, whose elements are possibly empty subsets of S , and where the product is defined coordinate-wise:

$$A \cdot B = \{a \cdot b : a \in A, b \in B\} \quad \text{for } A, B \subseteq S.$$

Exercise 5. (1) Let us view semigroups as a category, where the objects are semigroups and the morphisms are semigroup homomorphisms. What are the product and co-products of this category?

Exercise 6. (2) Let Σ be an alphabet, and let

$$X \subseteq \Sigma^+ \times \Sigma^+$$

be a set of words pairs. Define \sim_X to be least congruence on Σ^+ which contains all pairs from X . This is the same as the symmetric transitive closure of

$$\{(wxy, wxy) : w, y \in \Sigma^+, (x, y) \in X\}.$$

Show that the following problem – which is called the *word problem for semigroups* – is undecidable: given finite Σ, X and $w, v \in \Sigma^+$, decide if $w \sim_X v$.

Exercise 7. (2) Define the *theory of semigroups* to be the set of first-order sentences, which use one ternary relation $x = y \cdot z$, that are true in every semigroup. Show that the theory of semigroups is undecidable, i.e. it is undecidable if a first-order sentence is true in all semigroups.

Exercise 8. (2) Show that the theory of finite semigroups is different from the theory of (all) semigroups, but still undecidable.

² This exercise can be seen as the semigroup version of the pumping lemma.

1.1 Recognising languages

In this book, we are interested in monoids and semigroups as an alternative to finite automata for the purpose of recognising languages. Since languages are usually defined for possibly empty words, we use monoids and not semigroups when recognising languages.

Definition 1.5. Let Σ be a finite alphabet. A language $L \subseteq \Sigma^*$ is *recognised* by a monoid homomorphism

$$h : \Sigma^* \rightarrow M$$

if membership in $w \in L$ is determined uniquely by $h(w)$. In other words, there is a subset $F \subseteq M$ such that

$$w \in L \quad \text{iff} \quad h(w) \in F \quad \text{for every } w \in \Sigma^*.$$

We say that a language is recognised by a monoid if it is recognised by some monoid homomorphism into that monoid. The following theorem shows that, for the purpose of recognising languages, finite monoids and finite automata are equivalent.

Theorem 1.6. *The following conditions are equivalent for every $L \subseteq \Sigma^*$:*

- (1) *L is recognised by a finite nondeterministic automaton;*
- (2) *L is recognised by a finite monoid.*

Proof

2 \Rightarrow 1 From a monoid homomorphism one creates a deterministic automaton, whose states are elements of the monoid, the initial state is the identity, and the transition function is

$$(m, a) \mapsto m \cdot (\text{homomorphic image of } a).$$

After reading an input word, the state of the automaton is its homomorphic image, and therefore the accepting state from the monoid homomorphisms can be used. This automaton computes the monoid product according to the choice of parentheses illustrated in this example:

$$((((ab)c)d)e)f)g.$$

1 \Rightarrow 2 Let Q be the states of the nondeterministic automaton recognising L . Define a function³

$$\delta : \Sigma^* \rightarrow \text{monoid of binary relations on } Q$$

³ This transformation from a nondeterministic (or deterministic) finite automaton to a monoid incurs an exponential blow-up, which is unavoidable in the worst case.

which sends a word w to the binary relation

$$\{(p, q) \in Q^2 : \text{some run over } w \text{ goes from } p \text{ to } q\}.$$

This is a monoid homomorphism. It recognises the language: a word is in the language if and only if its image under the homomorphism contains at least one (initial, accepting) pair.

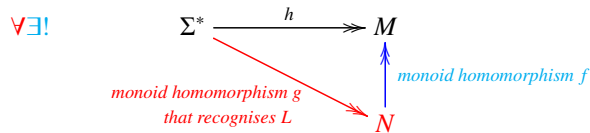
□

The syntactic monoid of a language. Deterministic finite automata have minimisation, i.e. for every language there is a minimal deterministic automaton, which can be found inside every other deterministic automaton that recognises the language. The same is true for monoids, as proved in the following theorem.

Theorem 1.7. *For every language⁴ $L \subseteq \Sigma^*$ there is a surjective monoid homomorphism*

$$h : \Sigma^* \rightarrow M,$$

called the syntactic homomorphism of L , which recognises it and is minimal in the sense explained in the following quantified diagram⁵



Proof The proof is the same as for the Myhill-Nerode theorem about minimal automata, except that the corresponding congruence is two-sided. Define the *syntactic congruence* of L to be the equivalence relation \sim on Σ^* which identifies two words $w, w' \in \Sigma^*$ if

$$uwv \in L \quad \text{iff} \quad uw'v \in L \quad \text{for all } u, v \in \Sigma^*.$$

Define h to be the function that maps a word to its equivalence class under syntactic congruence. It is not hard to see that h is compositional, and therefore by (the monoid version of) Lemma 1.3, one can equip the set of equivalence classes of syntactic congruences with a monoid structure – call M the resulting monoid – which turns h into a monoid homomorphism.

⁴ The language need not be regular, and the alphabet need not be finite.

⁵ Here is how to read the diagram. For every red extension of the black diagram there exists a unique blue extension which makes the diagram commute. Double headed arrows denote surjective homomorphisms, which means that \forall quantifies over surjective homomorphisms, and the same is true for $\exists!$.

It remains to show minimality of h , as expressed by the diagram in the lemma. Let then g be as in the diagram. Because g recognises the language L , we have

$$g(w) = g(w') \quad \text{implies} \quad w \sim w',$$

which, thanks to surjectivity of g , yields some function f from N to M , which makes the diagram commute, i.e. $h = f \circ g$. Furthermore, f must be a monoid homomorphism, because

$$\begin{aligned} f(a_1 \cdot a_2) &= \text{(by surjectivity of } g, \text{ each } a_i \text{ can be presented as } g(w_i) \text{ for some } w_i) \\ f(g(w_1) \cdot g(w_2)) &= \text{(} g \text{ is a monoid homomorphism)} \\ f(g(w_1 w_2)) &= \text{(the diagram commutes)} \\ h(w_1 w_2) &= \text{(} h \text{ is a monoid homomorphism)} \\ h(w_1) \cdot h(w_2) &= \text{(the diagram commutes)} \\ f(g(w_1)) \cdot f(g(w_2)) &= \\ f(a_1) \cdot f(a_2). & \end{aligned}$$

□

Exercise 9. (1) Show that the translation from deterministic finite automata to monoids is exponential in the worst case.

Exercise 10. (1) Show that the translation from (left-to-right) deterministic finite automata to monoids is exponential in the worst case, even if there is a right-to-left deterministic automaton of same size.

Exercise 11. (1) Show that a language $L \subseteq \Sigma^*$ is recognised by a finite commutative monoid if and only if it can be defined by a finite Boolean combination of conditions of the form “letter a appears exactly n times” or “the number of appearances of letter a is congruent to ℓ modulo n ”.

Exercise 12. (1) Prove that surjectivity of g is important in Theorem 1.7.

Exercise 13. (1) Show that for every language, not necessarily regular, its syntactic homomorphism is the function

$$w \in \Sigma^* \quad \mapsto \quad \underbrace{(q \mapsto qw)}_{\substack{\text{state transformation} \\ \text{in the syntactic automaton}}}$$

where the syntactic automaton is the deterministic finite automaton from the Myhill-Nerode theorem.

Exercise 14. (2) Let \mathcal{L} be a class of regular languages with the following closure properties:

- \mathcal{L} is closed under Boolean combinations;
- \mathcal{L} is closed under inverse images of homomorphisms $h : \Sigma^* \rightarrow \Gamma^*$;
- Let $L \subseteq \Sigma^*$ be a language in \mathcal{L} . For every $u, w \in \Sigma^*$, \mathcal{L} contains the inverse image of L under the following operation:

$$v \mapsto uvw.$$

Show that if L belongs to \mathcal{L} , then the same is true for every language recognised by its syntactic monoid.

1.2 Green's relations and the structure of finite semigroups

In this section, we describe some of the structural theory of finite semigroups. This theory is based on Green's relations, which are pre-orders in a semigroup that are based on prefixes, suffixes and infixes.

We begin with idempotents, which are ubiquitous in the analysis of finite semigroups. A semigroup element e is called *idempotent* if it satisfies

$$ee = e.$$

Example 1.8. In a group, there is a unique idempotent element, namely the group identity. There can be several idempotent elements, for example all elements are idempotent in the semigroup

$$(\{1, \dots, n\}, \max).$$

One can think of idempotents as being a relaxed version of identity elements.

Lemma 1.9 (Idempotent Power Lemma). *Let S be a finite semigroup. For every $a \in S$, there is exactly one idempotent in the set*

$$\{a^1, a^2, a^3, \dots\} \subseteq S.$$

Proof Because the semigroup is finite, the sequence a^1, a^2, \dots must contain a repetition, i.e. there must exist $n, k \in \{1, 2, \dots\}$ such that

$$a^n = a^{n+k} = a^{n+2k} = \dots.$$

After multiplying both sides of the above equation by a^{nk-n} we get

$$a^{nk} = a^{nk+k} = a^{nk+2k} = \dots,$$

and therefore $a^{nk} = a^{nk+nk}$ is an idempotent. To prove uniqueness of the idempotent, suppose $n_1, n_2 \in \{1, 2, \dots\}$ are powers such that that a^{n_1} and a^{n_2} are idempotent. Then we have

$$\underbrace{a^{n_1} = (a^{n_1})^{n_2}}_{\substack{\text{because } a^{n_1} \\ \text{is idempotent}}} = a^{n_1 n_2} = \underbrace{(a^{n_1})^{n_2}}_{\substack{\text{because } a^{n_2} \\ \text{is idempotent}}} = a^{n_2}$$

□

Finiteness is crucial for the above lemma, for example the infinite semigroup

$$(\{1, 2, \dots\}, +)$$

contains no idempotents. For $a \in S$, we use the name *idempotent power* for the element a^n , and we use the name *idempotent exponent* for the number n . The idempotent power is unique, but the idempotent exponent is not. It is easy to see that there is always an idempotent exponent which is at most the size of the semigroup, and idempotent exponents are closed under multiplication. Therefore, if a semigroup has n elements, then the factorial $n!$ is an idempotent exponent for every element of the semigroup. This motivates the following notation: we write $a^!$ for the idempotent power of a . The notation usually used in the semigroup literature is a^ω , but we will use ω for infinite words.

The analysis presented in the rest of this chapter will hold in any semigroup which satisfies the conclusion of the Idempotent Power Lemma.

Green's relations

We now give the main definition of this chapter.

Definition 1.10 (Green's relations). Let a, b be elements of a semigroup S . We say that a is a *prefix* of b if there exists a solution x of

$$ax = b.$$

The solution x can be an element of the semigroup, or empty (i.e. $a = b$). Likewise we define the suffix and infix relations, but with the equations

$$\underbrace{xa = b}_{\text{suffix}} \quad \underbrace{xay = b}_{\text{infix}}.$$

In the case of the infix relation, one or both of x and y can be empty.

Figure 1.1 shows a monoid along with the accompanying Green's relations. The prefix, suffix and infix relations are pre-orders, i.e. they are transitive and reflexive⁶. They need not be anti-symmetric, for example in a group every element is an prefix (suffix, infix) of every other element. We say that two elements of a semigroup are in the same *prefix class* if they are prefixes of each other. Likewise we define *suffix classes* and *infix classes*.

Clearly every prefix class is contained in some infix class, because prefixes are special cases of infixes. Therefore, every infix class is partitioned into prefix classes. For the same reasons, every infix class is partitioned into suffix classes. The following lemma describes the structure of these partitions.

Lemma 1.11 (Egg-box lemma). *The following hold in every finite semigroup.*

- (1) *all distinct prefix classes in a given infix class are incomparable:*

a, b are infix equivalent, and a is a prefix of $b \Rightarrow a, b$ are prefix equivalent

- (2) *if a prefix class and a suffix class are contained in the same infix class, then they have nonempty intersection;*
 (3) *all prefix classes in the same infix class have the same size.*

Of course, by symmetry, the lemma remains true after swapping prefixes with suffixes.

Proof

- (1) This item says that distinct prefix classes in the same infix class are incomparable, with respect to the prefix relation. This item of the Egg-box Lemma is the one that will be used most often.

Suppose that a, b are infix equivalent and a is a prefix of b , as witnessed by solutions x, y, z to the equations

$$b = ax \quad a = ybz.$$

⁶ Another description of the prefix pre-order is that a is a prefix of b if

$$aS^1 \supseteq bS^1. \tag{1.2}$$

In the above, S^1 is the monoid obtained from S by adding an identity element, unless it was already there. The sets aS^1, bS^1 are called *right ideals*. Because of the description in terms of inclusion of right ideals, the semigroup literature uses the notation

$$a \geq_{\mathcal{R}} b \stackrel{\text{def}}{=} aS^1 \supseteq bS^1$$

for the prefix relation. Likewise, $a \geq_{\mathcal{L}} b$ is used for the prefix relation, which is defined in terms of left ideals. Also, for some mysterious reason, $a \geq_{\mathcal{J}} b$ is used for the infix relation. We avoid this notation, because it makes longer words smaller.

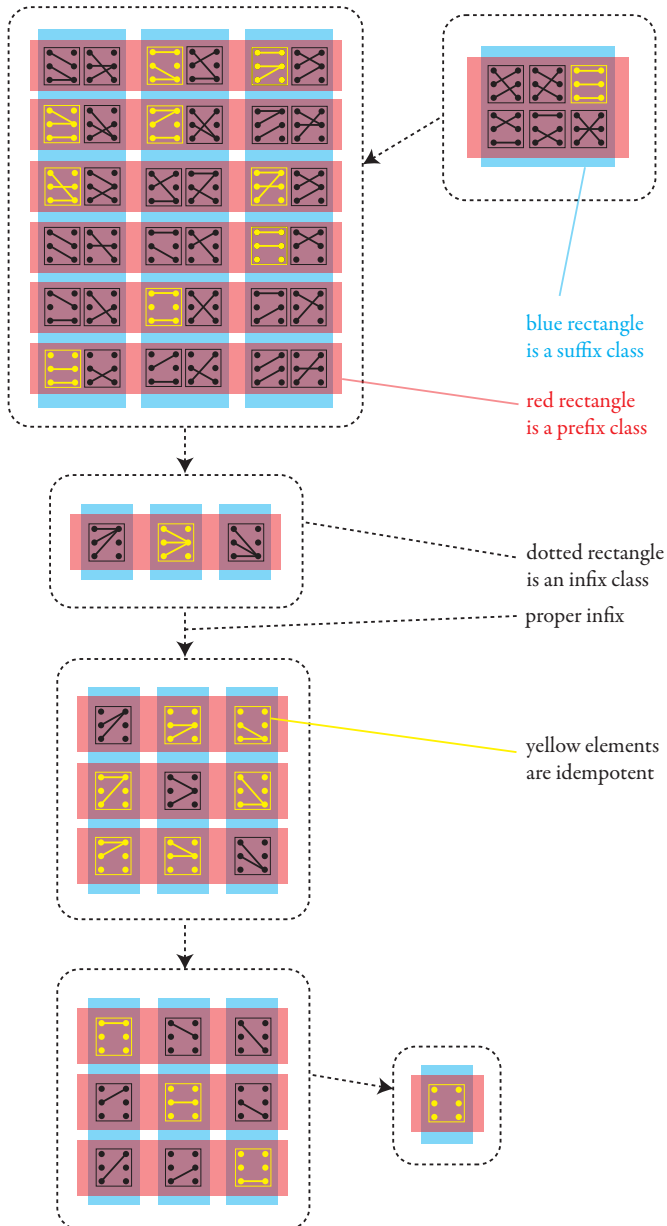
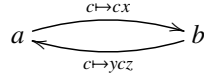


Figure 1.1 The semigroup of partial functions from a three element set to itself, partitioned into prefix, suffix and infix classes. In this particular example, the infix classes are totally ordered, which need not be the case in general.

As usual, each of x, y, z could be empty. This can be illustrated as



Consider the idempotent exponent $! \in \{1, 2, \dots\}$ which arises from Idempotent Power Lemma. We have:

$$\begin{aligned}
 b &= \text{(follow } !+! \text{ times the loop around } a, \text{ then go to } b) \\
 y^{!+!}a(xz)^{!+!}x &= \text{(} y^! \text{ is an idempotent)} \\
 y^!a(xz)^{!+!}x &= \text{(follow } ! \text{ times the loop around } a) \\
 & a(xz)^!x,
 \end{aligned}$$

which establishes that b is a prefix of a , and therefore a, b are in the same prefix class.

- (2) We now show that prefix and suffix classes in the same infix class must intersect. Suppose that a, b are in the same infix class, as witnessed by

$$a = xby.$$

With respect to the infix relation, by is between b and a , and therefore it must be in the same infix class as both of them. We have

$$\begin{array}{c}
 by \text{ is a suffix of } xby = a \\
 \underbrace{\hspace{1.5cm}} \\
 x \quad b \quad y \\
 \underbrace{\hspace{1.5cm}} \\
 b \text{ is a prefix of } by
 \end{array}
 ,$$

and therefore, thanks to the previous item, by is prefix equivalent to b and suffix equivalent to a . This witnesses that the prefix class of b and the suffix class of a have nonempty intersection.

- (3) We now show that all prefix classes in the same infix class have the same size. Take some two prefix classes in the same infix class, given by representatives a, b . We can assume that a, b are in the same suffix class, thanks to the previous item. Let

$$a = xb \quad b = ya$$

be witnesses for the fact that a, b are in the same suffix class. The following claim implies that the two prefix classes under consideration have the same size.

Claim 1.12. *The following maps are mutually inverse bijections*

$$\text{prefix class of } a \xrightleftharpoons[c \rightarrow xc]{c \rightarrow yc} \text{prefix class of } b$$

Proof Suppose that c is in the prefix class of a , as witnessed by a decomposition $c = az$. If we apply sequentially both maps in the statement of the claim to c , then we get

$$xyc = xyaz \stackrel{ya=b}{=} xbz \stackrel{xb=a}{=} az \stackrel{az=c}{=} c.$$

This, and a symmetric argument for the case when c is in the prefix class of b , establishes that the maps in the statement of the claim are mutually inverse. It remains to justify that the images of the maps are as in the statement of the claim, i.e. the image of the top map is the prefix class of b , and the image of the bottom map is the prefix class of a . Because the two maps are mutually inverse, and they prepend elements to their inputs, it follows that each of the maps has its image contained in the infix class of a, b . To show that the image of the top map is in the prefix class of b (a symmetric argument works for the bottom map), we observe that every element of this image is of the form yaz , and therefore it has $b = ya$ as a prefix, but it is still in the same infix class as a, b as we have observed before, and therefore it must be prefix equivalent to b thanks to the item (1) of the lemma. \square

\square

The Egg-box Lemma establishes that each infix class has the structure of a rectangular grid (which apparently is reminiscent of a box of eggs), with the rows being prefix classes and the columns being suffix classes. Let us now look at the eggs in the box: define an \mathcal{H} -class to be a nonempty intersection of some prefix class and some suffix class. By item (2) of the Egg-box Lemma, every pair of prefix and suffix classes in the same infix class lead to some \mathcal{H} -class. The following lemma shows that all \mathcal{H} -classes in the same infix class have the same size.

Lemma 1.13. *If a, b are in the same infix class, then there exist possibly empty x, y such that the following is a bijection*

$$\mathcal{H}\text{-class of } a \xrightarrow{c \mapsto xcy} \mathcal{H}\text{-class of } b$$

Proof Consider first the special case of the lemma, when a and b are in the same suffix class. Take the map from Claim 1.12, which maps bijectively the prefix class of a to the prefix class of b . Since this map preserves suffix classes, it maps bijectively the \mathcal{H} -class of a to the \mathcal{H} -class of b . By a symmetric argument, the lemma is also true when a and b are in the same prefix class.

For the general case, we use item (2) of the Egg-box Lemma, which says that there must be some intermediate element that is in the same prefix class

as a and in the same suffix class as b , and we can apply the previously proved special cases to go from the \mathcal{H} -class of a to the \mathcal{H} -class of the intermediate element, and then to the \mathcal{H} -class of b . \square

The following lemma shows a dichotomy for an \mathcal{H} -class: either it is a group, or the the product of every two elements in that \mathcal{H} -class falls outside the infix class.

Lemma 1.14 (*\mathcal{H} -class Lemma*). *The following conditions are equivalent for every \mathcal{H} -class G in a finite semigroup:*

- (1) G contains an idempotent;
- (2) ab is in the same infix class as a and b , for some $a, b \in G$;
- (3) $ab \in G$ for some $a, b \in G$;
- (4) $ab \in G$ for all $a, b \in G$;
- (5) G is a group (with product inherited from the semigroup)

Proof Implications (5) \Rightarrow (1) \Rightarrow (2) in the lemma are obvious, so we focus on the remaining implications.

(2) \Rightarrow (3) Suppose that ab is in the same infix class as a and b . Since a is a prefix of ab , and the two elements are in the same infix class, item (1) of the Egg-box Lemma implies that ab is in the prefix class of a , which is the same as the prefix class of b . For similar reasons, ab is in the same suffix class as a and b , and therefore $ab \in G$.

(3) \Rightarrow (4) Suppose that there exist $a, b \in G$ with $ab \in G$. We need to show that G contains the product of every elements $c, d \in G$. Since c is prefix equivalent to a there is a decomposition $a = xc$, and for similar reasons there is a decomposition $b = dy$. Therefore, cd is an infix of

$$\underbrace{a}_{xc} \underbrace{b}_{dy} \in G,$$

and therefore it is in the same infix class as G . Since c is a prefix of cd , and both are in the same infix class, the Egg-box Lemma implies that cd is in the prefix class of c . For similar reasons cd is in the suffix class of d . Therefore, $cd \in G$.

(4) \Rightarrow (5) Suppose that G is closed under products, i.e. it is a subsemigroup. We will show that it is a group. By the Idempotent Power Lemma, G contains some idempotent, call it e . We claim that e is an identity element in G , in particular it is unique. Indeed, let $a \in G$. Because a and e are in the same suffix class, it follows that a can be written as xe , and therefore

$$ae = xee = xe = a.$$

For similar reasons, $ea = a$, and therefore e is an identity element in G . The group inverse is defined as follows. Take $! \in \{1, 2, \dots\}$ to be the idempotent exponent which arises from the Idempotent Power Lemma. For every $a \in G$, the power $a^!$ is an idempotent. Since there is only one idempotent in G , we have $a^! = e$. Therefore, $a^{!-1}$ is a group inverse of a .

□

Exercise 15. (1) Show that for every finite monoid, the infix class of the monoid identity is a group.

Exercise 16. (2) Consider a finite semigroup. Show that an infix class contains an idempotent if and only if it is *regular*, which means that there exist a, b in the infix class such that ab is also in the infix class.

Exercise 17. (2) Show that if G_1, G_2 are two \mathcal{H} -classes in the same infix class of a finite semigroup, and they are both groups, then they are isomorphic as groups⁷.

Exercise 18. (2) We say that semigroup is *prefix trivial* if its prefix classes are singletons. Show that a finite semigroup S is prefix trivial if and only if it satisfies the identity

$$(xy)^! = (xy)^!x \quad \text{for all } x, y \in S.$$

Exercise 19. (2) Define the *syntactic semigroup* of a language to be the subset of the syntactic monoid which is the image of the nonempty words under the syntactic homomorphism. The syntactic semigroup may be equal to the syntactic monoid. We say that a language $L \subseteq \Sigma^*$ is definite if it is a finite Boolean combination of languages of the form $w\Sigma^*$, for $w \in \Sigma^*$. Show that a language is definite if and only if its syntactic semigroup S satisfies the identity

$$x^! = x^!y \quad \text{for all } x, y \in S.$$

Exercise 20. (2) Show two regular languages such that one is definite and the other is not, but both have isomorphic syntactic monoids.

⁷ Let us combine Exercises 16 and 17. By Exercises (16) and the \mathcal{H} -class lemma, an infix class is regular if and only if it contains an \mathcal{H} -class which is a group. By Exercise (17), the corresponding group is unique up to isomorphism. This group is called the *Shützenberger group* of the regular infix class.

Exercise 21. (1) Consider semigroups S which satisfy the following property: (*) that there is an infix class $J \subseteq S$ such that every $a \in S$ is an infix of J , or an absorbing zero element. Show that every finite semigroup is sub-semigroup of a product of finite semigroups that satisfy (*).

Exercise 22. (2) Show that every finite semigroup satisfies

$$\forall x_1 \forall x_2 \exists y_1 \exists y_2 \underbrace{z_1 = z_1 z_1 = z_1 z_2 \wedge z_2 = z_2 z_2 = z_2 z_1,}_{\text{where } z_i = x_i y_i}$$

where quantifiers range over elements of the finite semigroup.

Exercise 23. (2) Show that the following problem is decidable:

- **Input.** Two disjoint sets of variables

$$X = \{x_1, \dots, x_n\} \quad Y = \{y_1, \dots, y_m\}$$

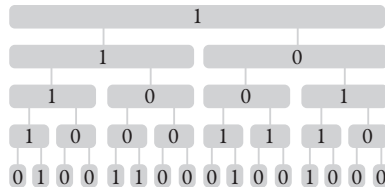
and two words $w, w' \in (X \cup Y)^+$.

- **Question.** Is the following true in all finite semigroups:

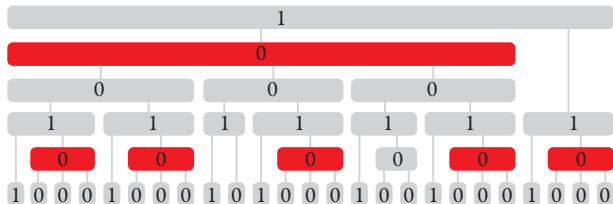
$$\forall x_1 \dots \forall x_n \exists y_1 \dots \exists y_m \underbrace{w = w'}_{\text{same product}}$$

1.3 The Factorisation Forest Theorem

In this section, we show how products in a semigroup can be organised in trees, so that in each node of the tree the products are very simple. The most natural way to do this is to have binary tree, as in the following example, which uses the two semigroup $\{0, 1\}$ with addition modulo 2:



We use the name *factorisation tree* for structures as in the above picture: a *factorisation tree over a semigroup S* is a tree, where nodes are labelled by semigroup elements, such that every node is either a leaf, or is labelled by the



The main result about Simon trees is that their height can be bounded by a constant that depends only on the semigroup, and not the underlying word.

Theorem 1.16 (Factorisation Forest Theorem). *Let S be a finite semigroup. Every word in S^+ admits a Simon tree of height¹⁰ $< 5|S|$.*

The rest of this chapter is devoted to proving the theorem.

Groups. We begin with the special case of groups.

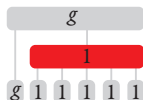
Lemma 1.17. *Let G be a finite group. Every word in G^+ admits a Simon tree of height $< 3|G|$.*

Proof Define the *prefix set* of a word $w \in G^+$ to be the set of group elements that can be obtained by taking a product of some nonempty prefix of w . By induction on the size of the prefix set, we show that every $w \in G^+$ has a Simon tree of height strictly less than 3 times the size of the prefix set. Since the prefix set has maximal size $|G|$, this proves the lemma.

The induction base is when the prefix set is a singleton $\{g\}$. This means that the first letter is g , and every other letter h satisfies $gh = g$. In a group, only the group identity $h = 1$ can satisfy $gh = g$, and therefore h is the group identity. In other words, if the prefix set is $\{g\}$, then the word is of the form

$$g \underbrace{1 \cdots 1}_{\text{a certain number of times}}.$$

Such a word admits a Simon tree as in the following picture:



¹⁰ The first version of this theorem was proved in [36, Theorem 6.1], with a bound of $9|S|$. The optimal bound is $3|S|$, which was shown in [23] Kufleitner, “The Height of Factorization Forests”, 2008, Theorem 1. The proof here is based on Kufleitner, with some optimisations removed.

The height of this tree is 2, which is strictly less than three times the size of the prefix set.

To prove the induction step, we show that every $w \in G^+$ admits a Simon tree, whose height is at most 3 plus the size from the induction assumption. Choose some g in the prefix set of w . Decompose w into factors as

$$w = \underbrace{w_1 w_2 \cdots w_{n-1}}_{\substack{\text{nonempty} \\ \text{factors}}} \underbrace{w_n}_{\substack{\text{could} \\ \text{be empty}}}$$

by cutting along all prefixes with product g . For the same reasons as in the induction base, every factor w_i with $1 < i < n$ has a product which is the group identity.

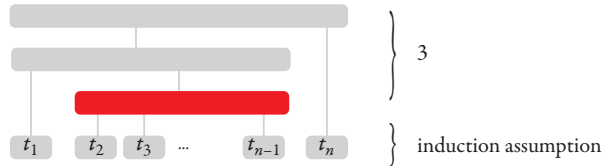
Claim 1.18. *The induction assumption applies to all of w_1, \dots, w_n .*

Proof For the first factor w_1 , the induction assumption applies, because its prefix set omits g . For the remaining blocks, we have a similar situation, namely

$$g \cdot (\text{prefix set of } w_i) \subseteq (\text{prefix set of } w) - \{g\} \quad \text{for } i \in \{2, 3, \dots, n\},$$

where the left side of the inclusion is the image of the prefix set under the operation $x \mapsto gx$. Since this operation is a permutation of the group, it follows that the left side of the inclusion has smaller size than the prefix set of w , and therefore the induction assumption applies. \square

By the above claim, we can apply the induction assumption to compute Simon trees t_1, \dots, t_n for the factors w_1, \dots, w_n . To get a Simon tree for the whole word, we join these trees as follows:

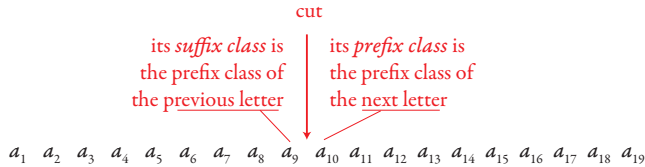


The gray nodes are binary, and the red node is idempotent because every w_i with $1 < i < n$ evaluates to the group identity. \square

Smooth words. In the next step, we prove the theorem for words where all infixes come from the same infix class. We say that a word $w \in S^+$ is *smooth* if all of its nonempty infixes have product in the same infix class. The following lemma constructs Simon trees for smooth words.

Lemma 1.19. *If a word is smooth, and the corresponding infix class is $J \subseteq S$, then it admits a Simon tree of height $< 4|J|$.*

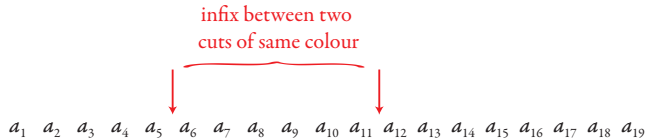
Proof Define a *cut* in a word to be the space between two consecutive letters; in other words this is a decomposition of the word into a nonempty prefix and a nonempty suffix. For a cut, define its *prefix* and *suffix* classes as in the following picture:



For every cut, both the prefix and suffix classes are contained in J , and therefore they have nonempty intersection thanks to item (2) of the Egg-box Lemma. This nonempty intersection is an \mathcal{H} -class, which is defined to be the *colour* of the cut. The following claim gives the crucial property of cuts and their colours.

Claim 1.20. *If two cuts have the same colour H , then the infix between them has product in H .*

Proof Here is a picture of the situation:



The infix begins with a letter from the prefix class containing H . Since the infix is still in the infix class J , by assumption on smoothness, it follows from item (1) of the Egg-box Lemma that the product of the infix is in the prefix class of H . For the same reason, the product of the infix is in the suffix class of H . Therefore, it is in H . \square

Define the *colour set* of a word to be the set of colours of its cuts; this is a subset of the \mathcal{H} -classes in J . Thanks to Lemma 1.12, all \mathcal{H} -classes contained in J have the same size, and therefore it makes sense to talk about the \mathcal{H} -class size in J , without specifying which \mathcal{H} -class is concerned.

Claim 1.21. *Every J -smooth word has a Simon tree of height at most*

$$|\text{colour set of } w| \cdot (3 \cdot \mathcal{H}\text{-class size} + 1).$$

Before proving the claim, we show how it implies the lemma. Since the number of possible colours is the number of \mathcal{H} -classes, the maximal height that can arise from the claim is

$$3 \cdot |J| + (\text{maximal size of colour set}) < 4|J|.$$

It remains to prove the claim.

Proof Induction on the size of the colour set. The induction base is when the colour set is empty. In this case the word has no cuts, and therefore it is a single letter, which is a Simon tree of height zero.

Consider the induction step. Let w be a smooth word. To prove the induction step, we will find a Simon tree whose height is at most the height from the induction assumption, plus

$$3 \cdot (\mathcal{H}\text{-class size}) + 1.$$

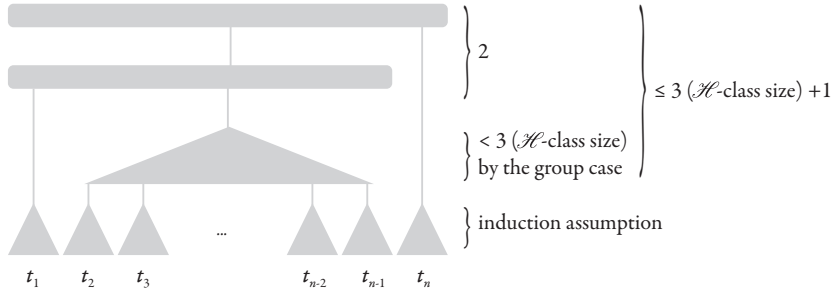
Choose some colour in the colour set of w , which is an \mathcal{H} -class H . Cut the word w along all cuts with colour H , yielding a decomposition

$$w = w_1 \cdots w_n.$$

None of the words w_1, \dots, w_n contain a cut with colour H , so the induction assumption can be applied to yield corresponding Simon trees t_1, \dots, t_n .

If $n \leq 3$, then the Simon trees from the induction assumption can be combined using binary nodes, increasing the height by at most 2, and thus staying within the bounds of the claim.

Suppose now that $n \geq 4$. By Claim 1.20, any infix between two cuts of colour H has product in H . In particular, all w_2, \dots, w_{n-1} have product in H , and the same is true for $w_2 w_3$. It follows that H contains at least one product of two elements from H , and therefore H is a group thanks to item (3) of the \mathcal{H} -class Lemma. Therefore, we can apply the group case from Lemma 1.17 to join the trees t_2, \dots, t_{n-1} . The final Simon tree looks like this:



□

□

General case. We now complete the proof of the Factorisation Forest Theorem. The proof is by induction on the *infix height* of the semigroup, which is defined to be the longest chain that is strictly increasing in the infix ordering. If the infix height is one, then the semigroup is a single infix class, and we can apply Lemma 1.19 since all words in S^+ are smooth. For the induction step, suppose that S has infix height at least two, and let $T \subseteq S$ be the elements which have a proper infix. It is not hard to see that T is a subsemigroup, and its induction parameter is smaller.

Consider a word $w \in S^+$. As in Lemma 1.19, define a cut to be a space between two letters. We say that a cut is *smooth* if the letters preceding and following the cut give a two-letter word that is smooth.

Claim 1.22. *A word in S^+ is smooth if and only if all of its cuts are smooth.*

Proof Clearly if a word is smooth, then all of its cuts must be smooth. We prove the converse implication by induction on the length of the word. Words of length one or two are vacuously smooth. For the induction step, consider a word $w \in S^+$ with all cuts being smooth. Since all cuts are smooth, all letters are in the same infix class. We will show that w is also in this infix class. Decompose the word as $w = vab$ where $a, b \in S$ are the last two letters. By induction assumption, va is smooth. Since the last cut is smooth, a and ab are in the same infix class, and therefore they are in the same prefix class by the Egg-box Lemma. This means that there is some x such that $abx = a$. We have

$$va = vabx = wx$$

which establishes that w is in the same infix class as va , and therefore in the same infix class as all the letters in w . □

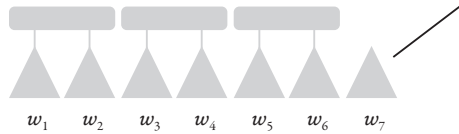
Take a word $w \in S^+$, and cut it along all cuts which are not smooth, yielding a factorisation

$$w = w_1 \cdots w_n.$$

By Claim 1.22, all of the words w_1, \dots, w_n are smooth, and therefore Lemma 1.19 can be applied to construct corresponding Simon trees of height strictly smaller than

$$4 \cdot (\text{maximal size of an infix class in } S - T).$$

Using binary nodes, group these trees into pairs, as in the following picture:



Each pair corresponds to a word with a non-smooth cut, and therefore each pair has product in T . Therefore, we can combine the paired trees into a single tree, using the induction assumption on a smaller semigroup. The resulting height is the height from the induction assumption on T , plus at most

$$1 + 4 \cdot (\text{maximal size of an infix class in } S - T) < 5|S - T|,$$

thus proving the induction step.

Exercises

Exercise 24. (1) Show that for every semigroup homomorphism

$$h : \Sigma^+ \rightarrow S \quad \text{with } S \text{ finite}$$

there is some $k \in \{1, 2, \dots\}$ such that for every $n \in \{3, 4, \dots\}$, every word of length bigger than n^k can be decomposed as

$$w_0 w_1 \cdots w_n w_{n+1}$$

such that all of the words w_1, \dots, w_n are mapped by h to the same idempotent.

Exercise 25. (2) Show optimality for the previous exercise, in the following sense. Show that for every $k \in \{1, 2, \dots\}$ there is some semigroup homomorphism

$$h : \Sigma^+ \rightarrow S \quad \text{with } S \text{ finite}$$

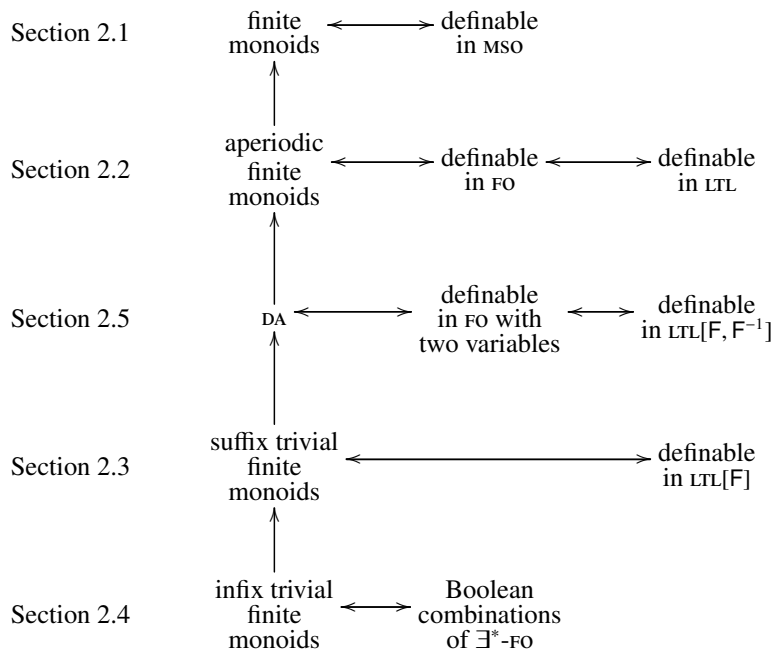
such that for every $n \in \{1, 2, \dots\}$ there is a word of length at least n^k which does not admit a factorisation $w_0 \cdots w_{n+1}$ where all of w_1, \dots, w_n are mapped by h to the same idempotent.

Exercise 26. (2) Let $h : \Sigma^* \rightarrow M$ be a monoid homomorphism. Consider a regular expression over Σ , which does not use Kleene star L^* but only Kleene plus L^+ . Such a regular expression is called h -typed if every subexpression has singleton image under h , and furthermore subexpressions with Kleene plus have idempotent image. Show that every language recognised by h is defined by finite union of h -typed expressions.

2

Logics on finite words, and the corresponding monoids

In this chapter, we show how structural properties of a monoid correspond to the logical power needed to define languages recognised by this monoid. We consider two kinds of logic: monadic second-order logic MSO and its fragments (notably first-order logic FO), as well as linear temporal logic LTL and its fragments. Here is a map of the results from this chapter, with horizontal arrows being equivalences, and the vertical arrows being strict inclusions.



2.1 All monoids and monadic second-order logic

We begin with monadic second-order logic (MSO), which is the logic that captures exactly the class of regular languages.

Logic on words. We assume that the reader is familiar with the basic notions of logic. The following descriptions are meant to fix notation. We use the name *vocabulary* for a set of relation names, each one with associated arity in $\{1, 2, \dots\}$. A *model* over a vocabulary consists of an underlying set (called the *universe* of the model), together with an interpretation of the vocabulary, which maps each relation name to a relation over the universe of corresponding arity. We allow the universe to be empty. For example, a directed graph is the same thing as a model over the vocabulary which contains one binary relation $E(x, y)$.

To express properties of models, we use first-order logic FO and MSO. Formulas of first-order logic over a given vocabulary are constructed as follows:

$$\underbrace{\forall x \quad \exists x}_{\substack{\text{quantification over} \\ \text{elements of the universe}}} \quad
 \underbrace{\varphi \wedge \psi \quad \varphi \vee \psi \quad \neg \varphi}_{\text{Boolean operations}} \quad
 \underbrace{R(x_1, \dots, x_n)}_{\substack{\text{an } n\text{-ary relation name from} \\ \text{the vocabulary applied} \\ \text{to a tuple of variables}}} \quad
 \underbrace{x = y}_{\text{equality}}$$

For the semantics of first-order formulas, we use the notation

$$\mathbb{A}, a_1, \dots, a_n \models \varphi(x_1, \dots, x_n)$$

to say that φ is true in the model \mathbb{A} , assuming that free variable x_i is mapped to $a_i \in \mathbb{A}$. A *sentence* is a formula without free variables.

The logic MSO extends first-order logic by allowing quantification over subsets of the universe (in other words, monadic relations over the universe, hence the name). The syntax of the logic has two kinds of variables: lower case variables x, y, z, \dots describe elements of the universe as in first-order logic, while upper case variables X, Y, Z, \dots describe subsets of the universe. Apart from the syntactic constructions of first-order logic, MSO also allows:

$$\underbrace{\forall X \quad \exists Y}_{\substack{\text{quantification over} \\ \text{subsets of the universe}}} \quad
 \underbrace{x \in X}_{\text{membership}}$$

We do not use more powerful logics (e.g. full second-order logic, which can also quantify over binary relations, ternary relations, etc.).

The following definition associates to each word a corresponding model. With this correspondence, we can use logic to define properties of words.

Definition 2.1 (Languages definable in first-order logic and mso). For a word $w \in \Sigma^*$, define its *ordered model* as follows. The universe is the set of positions in the word. The vocabulary is

$$\underbrace{x \leq y}_{\text{arity 2}} \quad \{ \underbrace{a(x)}_{\text{arity 1}} \}_{a \in \Sigma},$$

where $x \leq y$ is interpreted as the natural order on positions, and with $a(x)$ is interpreted as the set of positions with label a . For a sentence φ of mso over this vocabulary, we define its *language* to be

$$\{w \in \Sigma^* : \text{the ordered model of } w \text{ satisfies } \varphi\}.$$

A language is called *mso definable* if it is of this form. If φ is in first-order logic, then the language is called *first-order definable*.

Example 3. The language $a^*bc^* \subseteq \{a, b, c\}^*$ is first-order definable, as witnessed by the sentence:

$$\underbrace{\exists x}_{\text{there is a position}} \quad \underbrace{b(x)}_{\text{which has label } b} \wedge \underbrace{\forall y \overset{x \leq y \wedge x \neq y}{y < x} \Rightarrow a(y)}_{\text{and every earlier position has label } a} \wedge \underbrace{\forall y y > x \Rightarrow c(y)}_{\text{and every later position has label } b}.$$

□

Example 4. The language $(aa)^*a \subseteq a^*$ of words of odd length is mso definable, as witnessed by the sentence:

$$\underbrace{\exists X}_{\text{there is a set of positions}} \quad \underbrace{\forall x \overset{\forall y y \geq x}{\text{first}(x)} \vee \overset{\forall y y \leq x}{\text{last}(x)} \Rightarrow x \in X}_{\text{which contains the first and last positions,}} \wedge \underbrace{\forall x \forall y \overset{x < y \wedge \forall z z \leq x \vee y \leq z}{x = y + 1} \Rightarrow (x \in X \Leftrightarrow y \notin X)}_{\text{and contains every second position.}}$$

As we will see in Section 2.2, this language is not first-order definable. □

One could imagine other ways of describing a word via a model, e.g. a *successor model* where $x \leq y$ is replaced by a successor relation $x + 1 = y$. The successor relation can be defined in first-order logic in terms of order, but the converse is not true: there are languages that are first-order definable in the order model but not in the successor model, see Exercise 38. For the logic mso, there is no difference between successor and order, since the order can be defined in mso as follows

$$x \leq y \quad \text{iff} \quad \underbrace{\forall X (x \in X \wedge (\forall y \forall z y \in X \wedge y + 1 = z \Rightarrow z \in X))}_{X \text{ contains } x \text{ and is closed under successors}} \Rightarrow y \in X.$$

We now present the seminal Trakhtenbrot-Büchi-Elgot Theorem, which says that mso describes exactly the regular languages.

Theorem 2.2 (Trakhtenbrot-Büchi-Elgot). *A language $L \subseteq \Sigma^*$ is mso definable if and only if it is regular¹.*

This result is seminal for two reasons.

The first reason is that it motivates the search for other correspondences

$$\text{machine model} \quad \sim \quad \text{logic},$$

which can concern either restrictions or generalisations of the regular languages. In the case of restrictions, important examples are first-order logic and its fragments, which will be described later in this chapter. In this book, we do not study the generalisations; we are only interested in regular languages. Nevertheless, it is worth mentioning Fagin’s Theorem, which says that NP describes exactly the languages definable in existential second-order logic².

The second reason is that the Trakhtenbrot-Büchi-Elgot theorem generalises well to structures beyond finite words. For example, there are natural notions of mso definable languages for: infinite words, finite trees, infinite trees, graphs, etc. It therefore makes sense to search for notions of regularity – e.g. based on generalisations of semigroups – which have the same expressive power as mso. This line of research will also be followed in this book.

The rest of Section 2.1 proves the Trakhtenbrot-Büchi-Elgot Theorem.

The easy part is that every regular language is mso definable. Using the same idea as for the parity language in Example 4, the existence of a run of nondeterministic finite automaton can be formalised in mso. If the automaton has n states, then the formula looks like this:

$$\underbrace{\exists X_1 \exists X_2 \cdots \exists X_n}_{\text{existential set quantification}} \quad \underbrace{\text{“the sets } X_1, \dots, X_n \text{ describe an accepting run”}}_{\text{first-order formula}}$$

A corollary is that if we take any mso definable language, turn it into an automaton using the hard implication, and come back to mso using the easy implication, then we get an mso sentence of the form described above.

We now turn to the hard part, which says that every mso definable language

¹ This result was proved, independently, in the following papers:
 [39] Trakhtenbrot, “The synthesis of logical nets whose operators are described in terms of one-place predicate calculus (Russian)”, 1958, Theorems 1 and 2
 [9] Büchi, “Weak second-order arithmetic and finite automata”, 1960, Theorems 1 and 2
 [17] Elgot, “Decision problems of finite automata design and related arithmetics”, 1961, Theorem 5.3

² [18] Fagin, “Generalized first-order spectra and polynomial-time recognizable sets”, 1974, Theorem 6

is regular. This implication is proved in the rest of Section 2.1. For the definition of regularity, we use finite monoids. In other words, we will show that every mso definable language is recognised by a finite monoid. The idea is to construct the finite monoid by induction on formula size. In the induction, we also construct monoids for formulas with free variables, so we begin by dealing with those.

Definition 2.3 (Language of formulas with free variables). For an mso formula

$$\varphi(\underbrace{X_1, \dots, X_n}_{\text{all free variables are set variables}})$$

all free variables are set variables

which uses the vocabulary of the ordered model for words over alphabet Σ , define its *language* to be the set of words w over alphabet $\Sigma \times \{0, 1\}^n$ such that

$$\pi_\Sigma(w) \models \varphi(X_1, \dots, X_n),$$

where π_Σ is the projection of w onto the Σ coordinate, and X_i is the set of positions whose label has value 1 on the i -th bit of the bit vector from $\{0, 1\}^n$.

If the formula φ in the above definition has no free variables, then the above notion of language coincides with Definition 2.1. Therefore, the hard part of the Trakhtenbrot-Büchi-Elgot Theorem will follow immediately from Lemma 2.4 below.

Lemma 2.4. *If $\varphi(X_1, \dots, X_n)$ is a formula of mso where all free variables are set variables, then its language is recognised by a finite monoid.*

Proof Before proving the lemma, we observe that first-order variables can be eliminated from mso. Suppose that we extend mso with the following predicates that express properties of sets

$$\begin{array}{ccc} \underbrace{X \subseteq Y}_{\text{set inclusion}} & \underbrace{X \leq Y}_{\substack{x \leq y \text{ holds for} \\ \text{every } x \in X \text{ and} \\ \text{every } y \in Y}} & \underbrace{X \subseteq a}_{\substack{a(x) \text{ holds for} \\ \text{every } x \in X}}. \end{array} \quad (2.1)$$

The above predicates are second-order predicates in the sense that they express properties of sets; in contrast to the first-order predicates $x \leq y$ and $a(x)$ which express properties of elements. Using the second-order predicates, we can eliminate the first-order variables: instead of quantifying over a position x , we can quantify over a set of positions X , and then say that this set is a singleton:

$$\underbrace{X \neq \emptyset}_{X \text{ is nonempty}} \quad \underbrace{\forall Y Y \subseteq X \Rightarrow Y = \emptyset \vee Y = X}_{\text{and every proper subset of } X \text{ is empty}}.$$

Once elements are represented as singleton sets, the first-order predicates $x \leq y$ and $a(x)$ can be simulated using the second-order predicates from 2.1.

Using the transformation described above, from now on we assume that MSO has only set variables, and it uses the second-order predicates from 2.1. For such formulas, we prove the lemma by induction on formula size.

- *Induction base.* In the induction base, we need to show that for every atomic formula as in (2.1), its language is recognised by a finite monoid. Consider for example the formula $a \subseteq X$. The language of this formula consists of words over alphabet $\Sigma \times 2$ where for every position, the label satisfies:

$$\text{first coordinate is } a \quad \Rightarrow \quad \text{second coordinate is } 1.$$

This language is recognised by the homomorphism into the monoid

$$(\{0, 1\}, \min)$$

which maps letters that satisfy the implication to 1 and other letters to 0. Similar constructions can be done for the remaining predicates, in the case of $X \leq Y$ the monoid is not going to be commutative.

- *Boolean combinations.* For negation, the language of

$$\neg\varphi(X_1, \dots, X_n)$$

is recognised by the same homomorphism as the language of $\varphi(X_1, \dots, X_n)$, only the accepting set needs to be complemented. For conjunction

$$\varphi_1(X_1, \dots, X_n) \wedge \varphi_2(X_1, \dots, X_n),$$

one uses a product homomorphism

$$(h_1, h_2) : (\Sigma \times 2^n)^* \rightarrow M_1 \times M_2 \quad \text{where } h_i : (\Sigma \times 2^n)^* \rightarrow M_i \text{ recognises } \varphi_i,$$

with the accepting set consisting of pairs that are accepting on both coordinates. Disjunction \vee reduces to conjunction and negation using De Morgan's Laws.

- *Set quantification.* By De Morgan's Laws, it is enough to consider existential set quantification

$$\exists X_n \varphi(X_1, \dots, X_n)$$

The language of the quantified formula uses alphabet $\Sigma \times 2^{n-1}$. Let

$$h : (\Sigma \times 2^n)^* \rightarrow M$$

be a homomorphism that recognises the language of the formula $\varphi(X_1, \dots, X_n)$,

which is obtained by induction assumption. To recognise the quantified formula, we will use a powerset construction. Define

$$\pi : (\Sigma \times 2^n)^* \rightarrow (\Sigma \times 2^{n-1})^*$$

to be the letter-to-letter homomorphism which removes the last bit from every input position, and define

$$H : (\Sigma \times 2^{n-1})^* \rightarrow PM \quad H(w) = \{h(v) : \pi(v) = w\}.$$

It is not hard to see that the function H is a homomorphism, with the monoid structure on the powerset monoid defined by

$$A \cdot B = \{a \cdot b : a \in A, b \in B\} \quad \text{for } A, B \subseteq M.$$

The powerset construction clearly preserves finiteness, although at the cost of an exponential blow up. The accepting set consists of those subsets of M which have at least one accepting element.

□

The construction in the above lemma is effective, which means that given a sentence of mso, we can compute in finite time a recognising monoid homomorphism with an accepting set. Therefore, it is decidable if a sentence of mso is true in at least one finite word: check if the image of the monoid homomorphism contains at least one accepting element.

The proof of the “hard” implication in the Trakhtenbrot-Büchi-Elgot Theorem is very generic and will work without substantial changes in other settings, such as infinite words, trees or graphs. The “easy part” will become hard part in some generalisations – e.g. for some kinds of infinite words or for graphs – because these generalisations lack a suitable automaton model.

Exercises

Exercise 27. (3) Define \mathcal{U}_2 to be the monoid with elements $\{a, b, 1\}$ and product

$$xy = \begin{cases} y & \text{if } x = 1 \\ x & \text{otherwise.} \end{cases}$$

Show that every finite monoid can be obtained from \mathcal{U}_2 by applying Cartesian products, quotients (under semigroup congruences), sub-semigroups, and the powerset construction from Exercise 4.

Exercise 28. (2) For an alphabet Σ , consider the model where the universe is Σ^* , and which is equipped with the following relations:

$$\underbrace{x \text{ is a prefix of } y}_{\text{binary relation}} \quad \underbrace{\text{the last letter of } x \text{ is } a \in \Sigma}_{\text{one unary relation for each } a \in \Sigma}$$

Show that a language $L \subseteq \Sigma^*$ is regular if and only if there is a first-order formula $\varphi(x)$ over the above vocabulary such that L is exactly the words that satisfy $\varphi(x)$ in the above structure.

Exercise 29. (1) What happens if the prefix relation in Exercise 28 is replaced by the infix relation?

Exercise 30. (1) Consider the fragment of second-order logic where one can quantify over: elements, unary relations, and binary relations. (This fragment is expressively complete.) Define \equiv_k to be the equivalence on Σ^* which identifies two words if they satisfy the same sentences from the above fragment of second-order logic, up to quantifier rank k . Show that this equivalence relation has finite index, but it is not a semigroup congruence.

2.2 Aperiodic semigroups and first-order logic

We now begin the study of fragments of mso logic. We show that such fragments correspond to structural restrictions on finite monoids. The first – and arguably most important – fragment is first-order logic. This fragment will be described in the Shützenberger-McNaughton-Papert-Kamp Theorem. One part of the theorem says that a language is first-order definable if and only if it is recognised by a finite monoid M which satisfies

$$\underbrace{a^! = a^! a}_{\text{a monoid or semigroup which satisfies this is called aperiodic}},$$

a monoid or semigroup which satisfies this is called *aperiodic*

where $! \in \{1, 2, \dots\}$ is the idempotent exponent from the Idempotent Power Lemma. In other words, in an aperiodic monoid the sequence a, a^2, a^3, \dots is eventually constant, as opposed to having some non-trivial periodic behaviour.

Example 5. Consider the parity language $(aa)^* \subseteq a^*$. We claim that this language is not recognised by any aperiodic monoid, and therefore it is not first-order definable. Of course the same is true for the complement of the language, which was discussed in Example 4.

Suppose that the parity language is recognised by a homomorphism h into some finite monoid M . By Theorem 1.7, there is a surjective homomorphism from the image of h , which is a sub-monoid of M , into the syntactic monoid. In other words, the syntactic monoid is a quotient (i.e. image under a surjective homomorphism) of a sub-monoid of M . Since the syntactic monoid is the two-element group, which is not aperiodic, and since aperiodic monoids are closed under taking quotients and sub-monoids, it follows that M cannot be aperiodic.

The above argument shows that a regular language is first-order definable if and only if its syntactic monoid is aperiodic. Since the syntactic monoid can be computed, and aperiodicity is clearly decidable, it follows that there is an algorithm which decides if a regular language is first-order definable. \square

Apart from first-order logic and aperiodic monoids, the Shützenberger-McNaughton-Papert-Kamp Theorem considers also linear temporal logic and star-free regular expressions, so we begin by defining those.

Linear temporal logic Linear temporal logic (LTL) is an alternative to first-order logic which does not use quantifiers. The logic LTL only makes sense for structures equipped with a linear order; hence the name.

Definition 2.5 (Linear temporal logic). Let Σ be a finite alphabet. Formulas of *linear temporal logic* (LTL) over Σ are defined by the following grammar:

$$\underbrace{a \in \Sigma}_{\substack{\text{the current} \\ \text{position has} \\ \text{label } a}} \quad \varphi \wedge \psi \quad \varphi \vee \psi \quad \neg \varphi \quad \underbrace{\varphi \mathbf{U} \psi}_{\varphi \text{ until } \psi} .$$

The semantics for LTL formulas³ is a ternary relation, denoted by

$$\underbrace{w,}_{\substack{\text{word} \\ \text{in } \Sigma^+}} \underbrace{x}_{\substack{\text{position} \\ \text{in } w}} \models \underbrace{\varphi}_{\text{LTL formula}},$$

which is defined as follows. A formula $a \in \Sigma$ is true in positions with label a . The semantics of Boolean combinations are defined as usual. For formulas of the form $\varphi \mathbf{U} \psi$, the semantics⁴ are

$$w, x \models \varphi \mathbf{U} \psi \stackrel{\text{def}}{=} \underbrace{\exists y \ x < y}_{\substack{\text{there is some} \\ \text{position strictly} \\ \text{after } x}} \wedge \underbrace{w, y \models \psi}_{\text{which satisfies } \psi} \wedge \underbrace{\forall z \ x < z < y \Rightarrow w, z \models \varphi}_{\text{and such that all intermediate positions satisfy } \varphi} .$$

We say that an LTL formula is true in a word, without specifying a position, if

³ The semantics make sense for any linear order, possibly infinite, with positions coloured by Σ .

⁴ We use a variant of the until operator which is sometimes called *strict until*.

the formula is true in the first position of that word; this only makes sense for nonempty words. A language $L \subseteq \Sigma^*$ is called *LTL definable* if there is an LTL formula φ that defines the language on nonempty words:

$$w \in L \quad \text{iff} \quad w \models \varphi \quad \text{for every } w \in \Sigma^+.$$

For example, the formula aUb defines the language $\Sigma a^* b \Sigma^*$.

Example 6. To get a better feeling for LTL, we discuss some extra operators that can be defined using until. We write \perp for any vacuously false formula, e.g. $a \wedge \neg a$, likewise \top denotes any vacuously true formula. Here are some commonly used extra operators:

$$\underbrace{X\varphi}_{\text{the next position satisfies } \varphi} \stackrel{\text{def}}{=} \perp U \varphi, \quad \underbrace{F\varphi}_{\text{some strictly later position satisfies } \varphi} \stackrel{\text{def}}{=} \top U \varphi, \quad \underbrace{\varphi U^* \psi}_{\text{non-strict until}} \stackrel{\text{def}}{=} \psi \vee (\varphi U \psi).$$

Similarly, we define a non-strict version of the operator F , with $F^*\varphi = \varphi \vee F\varphi$. For example, the formula

$$F^*(a \wedge \underbrace{\neg F\top}_{\text{last position}})$$

says that the last position in the word has label a . \square

Almost by definition, every LTL definable language is also first-order definable. Indeed, by unfolding the definition, one sees that for every LTL formula there is a first-order formula $\varphi(x)$ that is true in the same positions.

Star-free languages. Apart from first-order logic, aperiodic monoids, and LTL, another equivalent formalism is going to be star-free expressions. As the name implies, star-free expressions cannot use Kleene star. However, in exchange they are allowed to use complementation (without star and complementation one could only define finite languages). For an alphabet Σ , the star-free expressions are those that can be defined using the following operations on languages:

$$\underbrace{a \in \Sigma}_{\text{the language that contains only the word } a} \quad \underbrace{\emptyset}_{\text{empty language}} \quad \underbrace{LK}_{\text{concatenation}} \quad \underbrace{L + K}_{\text{union}} \quad \underbrace{\bar{L}}_{\text{complementation with respect to } \Sigma^*}$$

Note that the alphabet needs to be specified to give meaning to the complementation operation. A language is called *star-free* if it can be defined by a star-free expression.

Example 7. Assume that the alphabet is $\{a, b\}$. The expression $\bar{\emptyset}$ describes the

full language $\{a, b\}^*$. Therefore

$$\bar{0} \cdot a \cdot \bar{0}$$

describes all words with at least one a . Taking the complement of the above expression, we get a star-free expression for the language b^* . \square

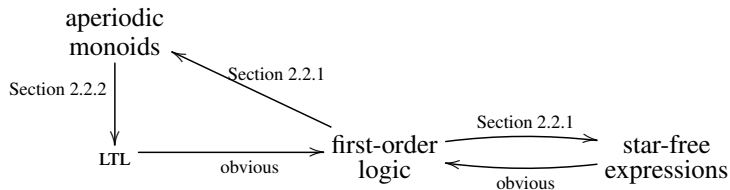
Like for LTL formulas, almost by definition every star-free expression describes a first-order definable language. This is because to every star-free expression one can associate a first-order formula $\varphi(x, y)$ which selects a pair of positions $x \leq y$ if and only if the corresponding infix (including x and y) belongs to the language described by the expression.

Equivalence of the models. The Shützenberger-McNaughton-Papert-Kamp Theorem says that all of the models discussed so far in this section are equivalent.

Theorem 2.6 (Shützenberger-McNaughton-Papert-Kamp). *The following are equivalent⁵ for every $L \subseteq \Sigma^*$:*

- (1) recognised by a finite aperiodic monoid;
- (2) star-free;
- (3) first-order definable;
- (4) LTL definable.

The rest of Section 2.2 is devoted to proving the theorem, according to the following plan:



⁵ This theorem combines three equivalences.

The equivalence aperiodic monoids and star-free expressions was shown in [33] Schützenberger, “On finite monoids having only trivial subgroups”, 1965, p. 190

The equivalence of star-free expressions and first-order logic was shown in [26] McNaughton and Papert, *Counter-free automata*, 1971, Theorem 10.5

The equivalence of first-order logic and LTL, not just for finite words, was shown in [22] Kamp, “Tense Logic and the Theory of Linear Order”, 1968

2.2.1 From first-order logic to aperiodic monoids and star-free expressions

In this section, we prove two inclusions: first-order logic is contained in both aperiodic monoids and star-free expressions.

Ehrenfeucht-Fraïssé games. In the proof, we use Ehrenfeucht-Fraïssé games, which are described as follows. An Ehrenfeucht-Fraïssé game is played by two players, called Spoiler and Duplicator. A configuration of the game is a pair of words (one red and one blue), each one with with a n -tuple of distinguished word positions

$$\underbrace{w, x_1, \dots, x_n}_{\text{in } \Sigma^* \quad \text{positions in } w} \quad \underbrace{w, x_1, \dots, x_n}_{\text{in } \Sigma^* \quad \text{positions in } w}$$

For such a configuration and $k \in \{0, 1, \dots\}$, the k -round game is played as follows. If there is a quantifier-free formula that distinguishes the two sides (red and blue), then Spoiler wins immediately and the game is stopped. Otherwise, the game continues as follows. If $k = 0$, then Duplicator wins. If $k > 0$ then Spoiler chooses one of the colours, and a distinguished position x_{n+1} in the word of the chosen colour. Duplicator responds with a matching distinguished position in the word of the other colour, and the game continues with $k - 1$ rounds from the configuration with $n + 1$ distinguished positions, which is obtained by adding the new distinguished positions. This completes the definition of Ehrenfeucht-Fraïssé games.

The point of Ehrenfeucht-Fraïssé games is that they characterise the expressive power of first-order logic, as stated in Theorem 2.7 below. The number of number of rounds in the games corresponds to quantifier rank of a formula, which is the nesting depth of quantifiers, as illustrated in the following example

$$\underbrace{\forall x (a(x) \Rightarrow \underbrace{(\exists y y < x \wedge b(y))}_{\text{quantifier rank 1}} \wedge \underbrace{(\exists y y > x \wedge b(y))}_{\text{quantifier rank 1}})}_{\text{quantifier rank 2}}$$

The correspondence of logic and games is given in the following theorem:

Theorem 2.7. *For every configuration of the game and $k \in \{0, 1, \dots\}$, Duplicator has a winning strategy in the game if and only if the two sides of the configuration satisfy the same formulas of first-order logic with quantifier rank at most k .*

Proof Straightforward induction on k . □

For $k \in \{0, 1, 2, \dots\}$ and $w, w' \in \Sigma^*$, we write

$$w \equiv_k w'$$

if the two words satisfy the same sentences of first-order logic with quantifier rank at most k , or equivalently, Duplicator has a winning strategy in the k -round game over the two words (with no distinguished positions). The following lemma characterises equivalence classes of \equiv_{k+1} in terms of equivalence classes of \equiv_k by using only Boolean combinations and concatenation.

Lemma 2.8. *Let $w, w' \in \Sigma^*$ and $k \in \{0, 1, \dots\}$. Then $w \equiv_{k+1} w'$ if and only if*

$$w \in LaK \Leftrightarrow w' \in LaK$$

holds for every $a \in \Sigma$ and every $L, K \subseteq \Sigma^$ which are equivalence classes of \equiv_k .*

Proof For the left-to-right implication, we observe that LaK can be defined by a first-order sentence of quantifier rank $k + 1$, which existentially quantifies over some position x with label a and then checks (using quantifier rank k) that the part before x belongs to L and the part after x belongs to K . Therefore, if w and w' satisfy the same sentences of quantifier rank $k + 1$, they must belong to the same languages of the form LaK .

Consider now the right-to-left implication. Here it will be useful to consider variant of the Ehrenfeucht-Fraïssé game, call it the *local game*. Consider a configuration of the Ehrenfeucht-Fraïssé game of the form

$$w, x_1, \dots, x_n \quad w', x_1, \dots, x_n.$$

where the red distinguished positions are listed in increasing order $x_1 < \dots < x_n$, and the same is true for the blue positions. Let us partition the positions x of w into the following $2n + 1$ sets, some of which may be empty:

$$\underbrace{X_0}_{x < x_1} \quad \{x_1\} \quad \underbrace{X_1}_{x_1 < x < x_2} \quad \{x_2\} \quad \cdots \quad \underbrace{X_{n-1}}_{x_{n-1} < x < x_n} \quad \{x_n\} \quad \underbrace{X_{n+1}}_{x_1 < x} \quad (2.2)$$

. Similarly, we partition the positions in the blue word w' . We say that a strategy of player Spoiler is local if there is some $i \in \{0, \dots, n\}$ such that all positions chosen by Spoiler in the strategy belong to X_i (for positions in the red word w) or X_i (for positions in the blue word w').

Claim 2.9. *If Spoiler has a winning strategy, then he also has a local one.*

Proof There is no benefit for Spoiler in using two different blocks of the partition described in (2.2). \square

A corollary of this claim is the if $n = 1$, the Spoiler has a winning strategy in the $(k + 1)$ -round game for the configuration

$$w, x_1 \quad w, x_1$$

if and only if: (1) the distinguished positions have different labels; or (2) Spoiler has a winning strategy in the k -round game for the parts strictly before the distinguished position; or (3) Spoiler has a winning strategy for the parts strictly after the distinguished position. This gives the right-to-left implication in the lemma. \square

We now use the lemma above to prove the inclusion of first-order logic in both star-free expressions and aperiodic monoids.

From first-order logic to star-free. It is enough to show that every equivalence class of \equiv_k is star-free. This is proved by induction on k . For the induction base, there is only one equivalence class, namely all words, which is clearly star-free. Consider now the induction step. Consider an equivalence class M of \equiv_{k+1} . Let be X the set of triples

$$\underbrace{L}_{\substack{\text{equivalence} \\ \text{class of } \equiv_k}} \quad \underbrace{a}_{\substack{\text{letter in } \Sigma}} \quad \underbrace{K}_{\substack{\text{equivalence} \\ \text{class of } \equiv_k}}.$$

By Lemma 2.8, if L, a, K are as above, then M is either contained in LaK or disjoint with LaK . Therefore, the equivalence class M is equal to the following Boolean combination of concatenations

$$\bigcap_{\substack{(L,a,K) \in X \\ M \subseteq LaK}} LaK \quad \cap \quad \bigcap_{\substack{(L,a,K) \in X \\ LaK \cap M = \emptyset}} \overline{LaK}.$$

This is a star-free expression, if we assume that L and K are described by star-free expressions from the induction assumption. Since every first-order definable language is a finite union of equivalence classes of \equiv_k for some k , the result follows.

From first-order logic to aperiodic monoids. An corollary of Lemma 2.8 is the following compositionality property for first-order logic on words.

Corollary 2.10. *For every alphabet Σ and $k \in \{0, 1, \dots\}$, the equivalence relation \equiv_k on Σ^* is a monoid congruence with finitely many equivalence classes.*

Proof Induction on k . To see that there are finitely many equivalence classes, we use Lemma 2.8, which says that an equivalence class of \equiv_{k+1} can be viewed as a set of triples (equivalence class of \equiv_k , letter from Σ , equivalence

class of \equiv_k), and there are finitely many possible sets of such triples. We now show that \equiv_{k+1} is a monoid congruence, i.e.

$$w \equiv_{k+1} w \text{ and } v \equiv_{k+1} v \quad \text{implies} \quad wv \equiv_{k+1} wv.$$

By Lemma 2.8, to prove the conclusion of the above implication, it is enough to show that wv and wv belong to the same languages of the form LaK as in the lemma. This follows immediately from the assumption of the implication, and the induction assumption of the lemma which says \equiv_k is a monoid congruence. (In the proof we also need the observation that \equiv_{k+1} refines \equiv_k , which follows from the definition of \equiv_k .) \square

By the above corollary, the function h_k which maps a word to its equivalence class under \equiv_k is a monoid homomorphism, into a finite monoid. This homomorphism recognises every language that is defined by a first-order sentence of quantifier rank at most k , by definition of \equiv_k . Therefore, every first-order definable language is recognised by h_k for some k . It remains to show that the monoid uses by such a homomorphism is aperiodic. To prove this, we use Lemma 2.8 and induction on k to show that

$$w^{2^k-1} \equiv_k w^{2^k} \quad \text{for every } w \in \Sigma^* \text{ and } k \in \{1, 2, \dots\}.$$

2.2.2 From aperiodic monoids to LTL

The last, and most important, step in the proof is constructing an LTL formula based on an aperiodic monoid⁶. In this part of the proof, semigroups will be more convenient than monoids. We will use LTL to define colourings, which are like languages but with possibly more than two values: a function from Σ^+ to a finite set of colours is called LTL definable if for every colour, the words sent that colour are an LTL definable language. For example, a semigroup homomorphism into a finite semigroup is a colouring.

Lemma 2.11. *Let S be a semigroup, and let $\Sigma \subseteq S$. The colouring*

$$w \in \Sigma^+ \quad \mapsto \quad \text{product of } w$$

is LTL definable.

By applying the lemma to the special case of S being a monoid, and substituting each monoid element for the letters that get mapped to it in the recognising homomorphism, we immediately get the implication from finite aperiodic monoids to LTL.

⁶ The proof in this section is based on [40] Wilke, ‘‘Classifying Discrete Temporal Properties’’, 1999, Section 2

It remains to prove the lemma. The proof is by induction on two parameters: the size of the semigroup S , and the size of the subset Σ . These parameters are ordered lexicographically, with the size of S being more important. Without loss of generality, we assume that Σ generates S , i.e. every element of S is the product of some word in Σ^+ .

The induction base is treated in the following claim.

Claim 2.12. *If either S or Σ has size one, then Lemma 2.11 holds.*

Proof If the semigroup has one element, there is nothing to do, since colourings with one possible colour are clearly LTL definable. Consider the case when the Σ contains only one element $a \in S$. By aperiodicity, the sequence

$$a, a^2, a^3, \dots$$

is eventually constant, because all powers bigger than the threshold $!$ give the same result. The product is therefore easily seen to be an LTL definable colouring. \square

We are left with the induction step. For $c \in S$, consider the function

$$a \in S \mapsto ca \in S.$$

Claim 2.13. *If $a \mapsto ca$ is a permutation of S , then it is the identity.*

Proof Suppose that $a \mapsto ca$ is a permutation of S , call it π . By aperiodicity,

$$\pi^! \circ \pi = \pi^!.$$

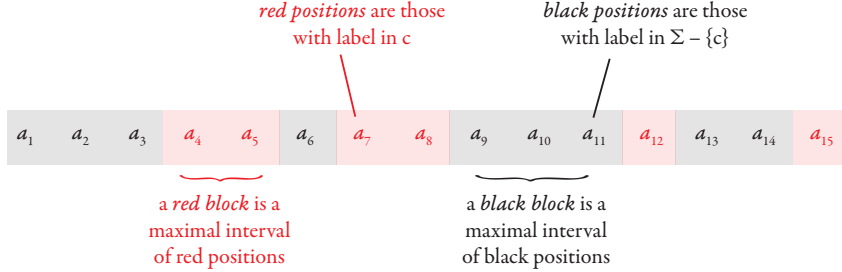
Since permutations form a group, it follows that π is the identity. \square

If the function $a \mapsto ca$ is the identity for every $c \in \Sigma$, then the product of a word is the same as its last letter; and such a colouring is clearly LTL definable. We are left with the case when there is some $c \in \Sigma$ such that $a \mapsto ca$ is not the identity. Fix this c for the rest of the proof. Define T to be the image of the function $a \mapsto ca$, this is a proper subset of S by assumption on c .

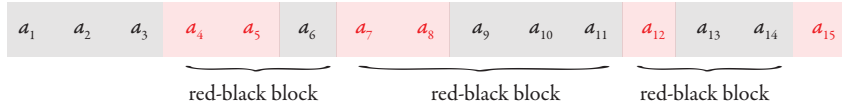
Claim 2.14. *T is a sub-semigroup of S .*

Proof If two semigroup elements have prefix c , then the same is true for their product. \square

In the rest of the proof, we use the following terminology for a word $w \in \Sigma^+$:



We first describe the proof strategy. For each black block, its product can be computed in LTL using the induction assumption on a smaller set of generators. The same is true for red blocks. Define a *red-black block* to be any union of a red block plus the following (non-empty) black block; as illustrated below:



For every red-black block, its product is in T because it begins with c and has at least two letters. Furthermore, the product can be computed in LTL, by using the products of the red and black blocks inside it. Using the induction assumption on a smaller semigroup, we compute the product of the union of all red-black blocks. Finally, the product of the entire word is obtained by taking into account the blocks that are not part of any red-black block.

The rest of this section is devoted to formalising the above proof sketch. In the formalisation, it will be convenient to reason with word-to-word functions. We say that a function a function of type $\Sigma^* \rightarrow \Gamma^*$ is an LTL *transduction* if it has the form

$$a_1 \cdots a_n \in \Sigma^* \quad \mapsto \quad f(a_1 \cdots a_n)f(a_2 \cdots a_n) \cdots f(a_n)$$

for some LTL definable colouring $f : \Sigma^+ \rightarrow \Gamma + \varepsilon$. By substituting formulas, one easily shows the following composition properties:

$$\begin{aligned}
 (\text{LTL colourings}) \circ (\text{LTL transductions}) &\subseteq \text{LTL colourings} \\
 (\text{LTL transductions}) \circ (\text{LTL transductions}) &\subseteq \text{LTL transductions.}
 \end{aligned}$$

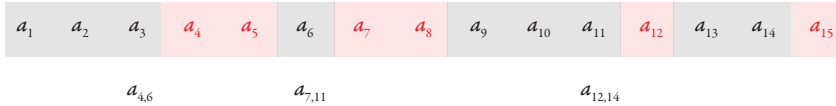
We use LTL transductions to decorate an input word $w \in \Sigma^+$ with extra information that will serve towards computing its product.

- (1) For each position that precedes a block (i.e. the next position begins a new block), write in that position the value of the next block. For the remaining positions, do not write anything. Use two disjoint copies of S to distinguish the values of the red and black blocks. Here is a picture:



In the above picture, $a_{i,j}$ denotes the product of the infix $\{i, \dots, j\}$. The function described in this step is an LTL transduction, thanks to the induction assumption on smaller alphabets⁷.

- (2) Take the output of the function in the previous step, and for each red letter (the product of a red block), multiply it with the next letter (which is the product of a black block). As a result, we get the values of all red-black blocks which do not begin in the first position. Here is a picture:



The function in this step is clearly an LTL transduction.

By induction assumption on a smaller semigroup, the product operation $T^+ \rightarrow T$ is an LTL colouring. By composing the functions described above with the semigroup product in T , we see that

$$w \in \Sigma^+ \mapsto \text{value of the union of red-black blocks}$$

is an LTL colouring. The values of the (at most two) blocks that do not participate in above union can also be computed using LTL colourings, and therefore the product of the entire word can be computed.

Exercises

⁷ To make this formal, we need a simple closure property of LTL that is described in Exercise 36.

Exercise 31. (2) Show that for every sentence of first-order logic, there is a sentence that is equivalent on finite words, and which uses at most three variables (but these variables can be repeatedly quantified).

Exercise 32. (2) Show that the following are equivalent for a finite semigroup:

- (1) aperiodic;
- (2) \mathcal{H} -trivial, which means that all \mathcal{H} -classes are singletons;
- (3) no sub-semigroup is a non-trivial group.

Exercise 33. (1) Consider the successor model of a word $w \in \Sigma^*$, which is defined like the ordered model, except that instead of $x < y$ we have $x + 1 = y$. Given an example of a regular language that is first-order definable using the ordered model, but not using the successor model.

Exercise 34. (1) Show two languages which have the same syntactic monoid, and such that only one of them is first-order definable in the successor model. In particular, one of the closure properties from Exercise 14 must fail for this logic.

Exercise 35. (3) Let Σ be a finite alphabet and let \vdash, \dashv be fresh symbols. For $k, \ell \in \{0, 1, \dots\}$, we say that $w, w' \in \Sigma^*$ are (k, ℓ) -locally equivalent if

$$\vdash w \dashv \text{ has at least } i \text{ occurrences of infix } v \quad \text{iff} \quad \vdash w' \dashv \text{ has at least } i \text{ occurrences of infix } v$$

holds for every $i \in \{0, \dots, k\}$ and every $v \in \Sigma^*$ of length at most ℓ . Show that $L \subseteq \Sigma^*$ is first-order definable in the successor model if and only if it is a union of equivalence classes of (k, ℓ) -local equivalence, for some k, ℓ .

Exercise 36. (1) Let $\Gamma \subseteq \Sigma$ and let $L \subseteq \Gamma^*$. If L is definable in LTL , then the same is true for

$$\{w \in \Sigma^* : L \text{ contains the maximal prefix of } w \text{ which uses only letters from } \Sigma\}.$$

Exercise 37. (2) Consider $\text{LTL}[X]$, i.e. the fragment of LTL where the only operator is X . Show that this fragment is equal to the definite languages from Exercise 19.

Exercise 38. (1) Show that if a language is first-order definable in the successor model, then the syntactic semigroup satisfies the following equality

$$eafbecf = ecfbeaf \quad \text{for all } \underbrace{e, f, a, b, c}_{\text{idempotents}}$$

Exercise 39. (3) Show that the identity in Exercise 38, together with aperiodicity, is equivalent to first-order definability in the successor model.

Exercise 40. (3) Consider the following extension of LTL with group operators. Suppose that G is a finite group, and let

$$\{\varphi_g\}_{g \in G},$$

be a family of already defined formulas such that every position in an input word is selected by exactly one formula φ_g . Then we can create a new formula, which is true in a word of length n if

$$1 = g_1 \cdots g_n,$$

where $g_i \in G$ is the unique group element whose corresponding formula selects position i . Show that this logic defines all regular languages.

2.3 Suffix trivial semigroups and temporal logic with F only

In the previous section, we showed that first-order logic corresponds to the monoids without groups, which is the same thing as monoids with trivial \mathcal{H} -classes (Exercise 32). What about monoids with trivial suffix classes, prefix classes, or infix classes? Trivial infix classes will be described in Section 2.4. In this section, we give a logical characterisation of trivial suffix classes (of course, a symmetric statement holds for trivial prefix classes).

In the characterisation, we use the fragment of LTL where until is replaced by the following operators

$$\underbrace{\top \text{U} \varphi}_{F\varphi} \quad \underbrace{\neg \text{F} \neg \varphi}_{G\varphi} \quad \underbrace{\varphi \vee \text{F} \varphi}_{F^* \varphi} \quad \underbrace{\neg \text{F}^* \neg \varphi}_{G^* \varphi}.$$

Since all of the above operators can be defined in terms of F, we write $\text{LTL}[F]$ for the resulting logic.

Theorem 2.15. *The following conditions are equivalent for $L \subseteq \Sigma^*$:*

- (1) Recognised by a finite suffix trivial monoid.
- (2) Defined by a finite union of regular expressions of the form

$$\underbrace{\Sigma_0^* a_1 \Sigma_1^* a_2 \cdots a_n \Sigma_n^*}_{\text{we call such an expression suffix unambiguous}} \quad \text{where } a_i \in \Sigma - \Sigma_i \text{ for } i \in \{1, \dots, n\}.$$

In the above, some of the sets $\Sigma_i \subseteq \Sigma$ might be empty, in which case $\Sigma_i^* = \{\varepsilon\}$.

- (3) Defined by a Boolean combination of $\text{LTL}[F]$ formulas of the form $F^* \varphi$.

To see why the formulas in item (3) need to be guarded by F^* , consider the $\text{LTL}[F]$ formula a which defines the language “words beginning with a ”. This language is not recognised by any finite suffix trivial monoid.

Proof

- (1) \Rightarrow (2) We will show that for every finite suffix trivial monoid M , and every $F \subseteq M$, the language

$$\{w \in M^* : F \text{ contains the product of } w\}$$

is defined by a finite union of suffix unambiguous expressions. It will follow that for every monoid homomorphism into M , the recognised language is defined by a similar expression, with monoid elements substituted by the letters that map to them (such a substitution preserves suffix unambiguity).

It is of course enough to consider the case when F contains only one element, call it $a \in M$. The proof is by induction on the position of a in the suffix ordering.

The induction base is when a is the monoid identity. By Exercise 15, the suffix class of the identity is a group, and a group must be trivial in a suffix trivial monoid. It follows that a word has product a if and only if it belongs to a^* , which is a suffix unambiguous expression.

We now prove the induction step. Consider a word with product a . This word must be nonempty, since otherwise its product would be the identity. Let i be the maximal position in the word such that the suffix starting in i also has product a . By suffix triviality, every position $< i$ is labelled by a letter in

$$\Sigma_0 = \{b \in M : ba = a\}.$$

Let b be the product of the suffix that starts after i , not including i , and let c be the label of position i . By choice of i , b is a proper suffix of a and $a = cb$.

Summing up, words with product a are defined by the expression

$$\bigcup_{\substack{b, c \in M \\ b \text{ is a proper suffix of } a \\ a = cb}} \Sigma_0^* c \cdot (\text{words with product } b),$$

Apply the induction assumption to b , yielding a finite union of suffix unambiguous expressions, and distribute the finite union across concatenation. It remains to justify that the resulting expressions are also suffix unambiguous. This is because none of the expressions that define words with product b can begin with Σ_1^* with $c \in \Sigma_1$, since otherwise we would contradict the assumption that $cb = a \neq b$.

(2) \Rightarrow (3) Since the formulas from item (3) are closed under union, it is enough to show that every suffix unambiguous expression

$$\Sigma_0^* a_1 \Sigma_1^* a_2 \cdots a_n \Sigma_n^*$$

can be defined by a formula as in (3). For $i \in \{0, \dots, n\}$, define L_i to be the suffix of the above expression that begins with Σ_i^* . By induction on i , starting with n and progressing down to 0, we show that L_i can be defined by a formula φ_i as in item (3). In the induction base, we use the formula

$$\varphi_n = \mathbf{G}^* \underbrace{\bigvee_{a \in \Sigma_n} a}_{\text{all positions have label in } \Sigma_n}.$$

For the induction step, we first define the language $a_i L_i$, using a formula of $\text{LTL}[F]$ (which is not in the shape from item (3)):

$$\psi_i = a_i \wedge (\mathbf{F}\varphi_i) \wedge \mathbf{G} \bigvee_{j>i} \varphi_j.$$

Because the expression is suffix unambiguous, the formula ψ_i selects at most one position in a given input word; this property will be used below. The language L_{i-1} is then defined by

$$\varphi_{i-1} = \mathbf{F}^* \psi_i \wedge \underbrace{\mathbf{G}^* ((\mathbf{F}\psi_i) \Rightarrow \bigwedge_{a \in \Sigma_0} a)}_{\substack{\text{if a position is to the left} \\ \text{of the unique position} \\ \text{satisfying } \psi_i, \text{ then} \\ \text{it has label in } \Sigma_0.}}$$

(3) \Rightarrow (1) Define the rank of a formula in $\text{LTL}[F]$ to be the nesting depth of the operator F . For $k \in \{0, 1, \dots\}$, define \approx_k to be the equivalence relation on Σ^+ which identifies two words if they satisfy the same formulas of rank at most k . The key observation is the following pumping lemma.

Claim 2.16. For every $k \in \{0, 1, 2, \dots\}$ we have

$$w(xy)^i u \approx_k wy(xy)^j u \quad \text{for every } w \in \Sigma^+, x, y, u \in \Sigma^* \text{ and } i, j \geq k.$$

Proof Induction on k . For $k = 0$, we observe that the equivalence class under \approx_0 depends only on the first letter, and the two words on both sides in the claim have the same letter because w is nonempty.

Consider now the induction step, when going from k to $k + 1$. By unravelling the definition of \approx_{k+1} , we need to show that if $i, j \geq k + 1$, then for every nonempty proper suffix of the word on one side of equivalence

$$w(xy)^i u \approx_{k+1} wy(xy)^j u$$

there is a nonempty proper suffix on the other side of the equivalence, such that the two suffixes are equivalent under \approx_k . Suppose first that v is a nonempty suffix of the left side. If v is a suffix of $(xy)^{k+1}u$, then the same v is a suffix of the right side. Otherwise, we can use the induction assumption. Consider now a nonempty proper suffix v of the right side. Here we argue in the same way as previously, except that there is one extra case, when

$$v = z(xy)^i u \quad \text{for some } z \text{ that is a suffix of } y.$$

In this case, the \approx_k -equivalent suffix on the left side is $z(xy)^k u$. □

By unravelling the definition of the syntactic monoid, in terms of two-sided congruences, we infer from the above claim that for every rank k formula φ of $\text{LTL}[F]$, the syntactic monoid M of $F^*\varphi$ satisfies

$$(xy)^! = y(xy)^! \quad \text{for all } x, y \in M. \tag{2.3}$$

The same is also true for syntactic monoids of Boolean combinations of such formulas. To finish the proof, we observe that property (2.3) is true in a finite monoid if and only if it is suffix trivial. Indeed, if a monoid is suffix trivial, then $(xy)^!$ and $y(xy)^!$ must be in the same suffix class, and hence equal. Conversely, if a, b are in the same suffix class, then there must be some x, y such that $b = xa$ and $a = yb$; it follows that

$$a = y(xy)^! b \stackrel{(2.3)}{=} (xy)^! b = b.$$

□

Exercises

Exercise 41. (1) Let Σ be an alphabet and let $c \notin \Sigma$ be a fresh letter. Show that $L \subseteq \Sigma^+$ satisfies the conditions of Theorem 2.15 if and only if cL is definable in $\text{LTL}[F]$.

2.4 Infix trivial semigroups and piecewise testable languages

In this section, we describe the languages recognised by finite monoids that are infix trivial. For languages recognised by finite infix trivial monoids, a prominent role will be played embeddings (also known as the Higman ordering).

Definition 2.17 (Embedding). We say that $w \in \Sigma^*$ *embeds* in $v \in \Sigma^*$, denoted by $w \hookrightarrow v$, if there is an injective function from positions in w to positions in v , which preserves the order on positions and the labels.

In other words, w embeds in v if and only if w can be obtained from v by removing zero or more positions. For example “ape” embeds into “example”. It is easy to see that embedding is an ordering: it is reflexive, transitive and anti-symmetric (although it will cease to be anti-symmetric for infinite words). We say that a language $L \subseteq \Sigma^*$ is *upward closed* if

$$v \hookrightarrow w \wedge v \in L \Rightarrow w \in L.$$

Symmetrically, we define downward closed languages. The main result about embedding is that it is a well-quasi order, as explained in the following lemma.

Lemma 2.18 (Higman’s Lemma). *For every upward closed $L \subseteq \Sigma^*$ there is a finite subset $U \subseteq L$ such that*

$$L = \underbrace{\{w \in \Sigma^* : v \hookrightarrow w \text{ for some } v \in U\}}_{\text{we call this the upward closure of } U}$$

Here is a logical corollary of Higman’s lemma.

Theorem 2.19. *A language is upward closed if and only if it can be defined in the ordered model by an \exists^* -sentence, i.e. a sentence of the form*

$$\underbrace{\exists x_1 \exists x_2 \cdots \exists x_n}_{\text{only existential quantifiers}} \underbrace{\varphi(x_1, \dots, x_n)}_{\text{quantifier-free}}.$$

Proof Clearly every \exists^* -sentence defines an upward closed language. Higman’s Lemma gives the converse implication, because the upward closure of every finite set is definable by an \exists^* -sentence. \square

Embeddings will also play an important role in the characterisation of languages recognised by monoids that are infix trivial. Before stating the characterisation, we introduce one more definition, namely zigzags⁸. For languages $L, K \subseteq \Sigma^*$, define a *zigzag between L and K* to be a sequence

$$\underbrace{w_1}_{\in L} \hookrightarrow \underbrace{w_2}_{\in K} \hookrightarrow \underbrace{w_3}_{\in L} \hookrightarrow \underbrace{w_4}_{\in K} \hookrightarrow \underbrace{w_5}_{\in L} \hookrightarrow \underbrace{w_6}_{\in K} \hookrightarrow \dots$$

In other words, this is a sequence that is growing with respect to embeddings, and such that odd-numbered elements are in L and even-numbered elements are in K . The zigzag does not need to be strictly growing, but it will be if L and K are disjoint.

We are now ready for the characterisation of infix trivial monoids.

Theorem 2.20. *The following conditions are equivalent⁹ for every $L \subseteq \Sigma^*$:*

- (1) *recognised by a finite monoid that is infix trivial;*
- (2) *is a finite Boolean combination of upward closed languages;*
- (3) *there is no infinite zigzag between L and its complement.*

We use the name *piecewise testable* for languages as in item (2) of the above theorem. Equivalence of items (2) and (3) is a corollary of the following lemma, when applied to $K = \Sigma^* - L$.

Lemma 2.21 (Zigzag Lemma). *Let $L, K \subseteq \Sigma^*$. The following are equivalent:*

- (1) *there are zigzags between L and K of every finite length;*
- (2) *there is an infinite zigzag between L and K ;*
- (3) *there is no piecewise testable language $M \subseteq \Sigma^*$ such that*

$$\underbrace{L \subseteq M \quad M \cap K = \emptyset}_{\text{we say that } M \text{ separates } L \text{ and } K}$$

Proof

- (1) \Rightarrow (2) Assume that zigzags between L and K can have arbitrarily long finite lengths. Define a directed acyclic graph G as follows. Vertices are words in L , and there is an edge $w \rightarrow v$ if

$$w \hookrightarrow u \hookrightarrow v \quad \text{for some } u \in K.$$

⁸ Zigzags and the Zigzag Lemma are inspired by
 [14] Czerwinski et al., “A Characterization for Decidable Separability by Piecewise Testable Languages”, 2017

⁹ Equivalence of items (1) and (2) was first proved in
 [37] Simon, “Piecewise testable events”, 1975

For a vertex $v \in L$ of this graph, define its *potential*

$$\alpha(v) \in \{0, 1, \dots, \omega\}$$

to be the least upper bound on the lengths of paths in the graph that start in v . This can be either a finite number, or ω if the paths have unbounded length.

We first show that some vertex must have potential ω . By assumption on arbitrarily long zigzags, potentials have arbitrarily high values. By definition of the graph, α is monotone with respect to (the opposite of the) embedding, in the following sense:

$$v \hookrightarrow v' \quad \text{implies} \quad \alpha(v) \geq \alpha(v') \quad \text{for every } v, v' \in L.$$

By Higman's Lemma, the language L , like any set of words, has finitely many minimal elements with respect to embedding. By monotonicity, one of these minimal words must therefore have potential ω .

For the same reason as above, if a word has potential ω , then one of its successors (words reachable in one step in the graph) must also have potential ω ; this is because there are finitely many successors that are minimal with respect to embedding. This way, we can construct an infinite path in the graph which only sees potential ω , as in the König Lemma.

- (2) \Rightarrow (3) Suppose that there is a zigzag between L and K of infinite length. Every upward closed set selects either no elements of the zigzag, or all but finitely many elements of the zigzag. It follows that every finite Boolean combination of upward closed sets must contain, or be disjoint with, two consecutive elements of the zigzag. Therefore, such a Boolean combination cannot separate L from K .
- (3) \Rightarrow (1) We prove the contra-positive: if zigzags between L and K have bounded length, then L and K can be separated by a piecewise testable language. For $w \in L$ define its *potential* to be the maximal length of a zigzag between L and K that starts in w ; likewise we define the potential for $w \in K$, but using zigzags between K and L . Define $L_i \subseteq L$ to be the words in L with potential exactly $i \in \{1, 2, \dots\}$, likewise define $K_i \subseteq K$. Our assumption is that the potential is bounded, and therefore L is a finite union of the languages L_i , likewise for K . By induction on $i \in \{0, 1, \dots\}$, we will show that

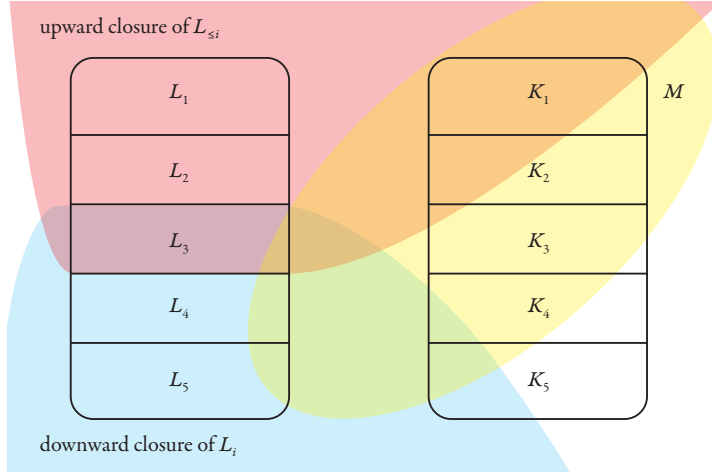
$$\underbrace{L_1 \cup \dots \cup L_i}_{L_{\leq i}} \quad \text{and} \quad \underbrace{K_1 \cup \dots \cup K_i}_{L_{\leq i}}$$

can be separated by a piecewise testable language, call it M_i . In the induction base, both languages are empty, and can therefore be separated by the empty language, which is clearly piecewise testable. Consider the induction step,

where we find the separator M_i . We will use the following sets

$$\underbrace{L_{\leq i} \uparrow}_{\text{upward closure of } L_{\leq i}} \qquad \underbrace{L_i \downarrow}_{\text{downward closure of } L_i} \qquad \underbrace{M}_{\text{a separator of } K_{< i} \text{ and } L_{< i} \text{ from the induction assumption}}$$

All of these sets are piecewise testable: the first one is upward closed, the second one is the complement of an upward closed set, and the third one is obtained from the induction assumption. These sets are depicted in the following picture, with $i = 3$:



The set $L_{\leq i} \uparrow$ contains $L_{\leq i}$ by definition. It is also disjoint with K_i , because otherwise there would be some words

$$\underbrace{w}_{L_{\leq i}} \quad \hookrightarrow \quad \underbrace{v}_{K_i},$$

and therefore w would need to have potential $i + 1$. For similar reasons, the set $L_i \downarrow$ is disjoint with K_i and $L_{\leq i-1}$. Putting these facts together, we see that

$$M_i = L_{\leq i} \uparrow - (M - L_i \downarrow)$$

separates $L_{\leq i}$ from $K_{\leq i}$.

□

The Zigzag Lemma proves that equivalence of the conditions about infinite zigzags and piecewise testability in Theorem 2.20. To finish the proof of the Theorem, we show that the syntactic monoid of L is finite and infix trivial (which is the same as saying that some recognising monoid is finite and infix trivial) if and only if there is no infinite zigzag between L and its complement.

Suppose first that the syntactic monoid of L is not finite or infix trivial. If the syntactic monoid is not finite, then the language cannot be piecewise testable, since piecewise testable languages are necessarily regular. Assume therefore that the syntactic monoid is finite but not infix trivial. This means that the syntactic monoid is either not prefix trivial, or not suffix trivial. By symmetry, we only consider the case where the syntactic monoid is not suffix trivial. This means that there exist a, b in the syntactic monoid such that

$$(ab)^! \neq b(ab)^!$$

By unravelling the definition of the syntactic monoid, the above disequality can be easily used to create an infinite zigzag between L and its complement.

It remains to show that if the syntactic monoid of L is finite and infix trivial, then there is no zigzag between L and its complement. Let M be the syntactic monoid. For $a, b \in M$, define a zigzag between a and b to be a zigzag, over alphabet M , between the words with product a and the words with product b . If M recognises L , then a zigzag between L and its complement can be used, by extraction, to obtain a zigzag between $a \neq b \in M$. The following lemma shows that this cannot happen, thus completing the proof of Theorem 2.20.

Lemma 2.22. *Let M be finite and infix trivial, and let $a, b \in M$. If there is an infinite zigzag between a and b , then $a = b$.*

Proof The proof is by induction on the infix ordering lifted to pairs:

$$(x, y) \leq (a, b) \stackrel{\text{def}}{=} x \text{ is an infix of } a \text{ and } y \text{ is an infix of } b.$$

The induction base is proved the same way as the induction step. Suppose that we have proved the lemma for all pairs $(x, y) < (a, b)$.

Claim 2.23. *If there is an infinite zigzag between a and b , then there exists $n \in \{0, 1, \dots\}$ and monoid elements $\{a_i, b_i, c_i\}_i$ such that*

$$\begin{aligned} a &= && c_1 && c_2 && \cdots && c_n \\ b &= &b_0 & c_1 & b_1 & c_2 & b_2 & \cdots & b_{n-1} & c_n & b_n \\ a &= &a_0 & c_1 & a_1 & c_2 & a_2 & \cdots & a_{n-1} & c_n & a_n \end{aligned}$$

and for every $i \in \{0, \dots, n\}$ there is an infinite zigzag between a_i and b_i .

Proof Consider an infinite zigzag between a and b of the form

$$w_1 \hookrightarrow w_2 \hookrightarrow \cdots$$

Let the letters in w_1 be $c_1, \dots, c_n \in M$. For $j \geq 2$, define an *important*

position in w_j to be any position that arises by starting in some position of w_1 , and then following the embeddings

$$w_1 \hookrightarrow w_2 \hookrightarrow \cdots \hookrightarrow w_j.$$

By distinguishing the important positions in w_j , we get a factorisation

$$w_j = \underbrace{w_{j,0}}_{M^*} c_1 \underbrace{w_{j,1}}_{M^*} c_2 \cdots c_{n-1} \underbrace{w_{j,n-1}}_{M^*} c_n \underbrace{w_{j,n}}_{M^*}.$$

By definition of important positions, for every $i \in \{0, \dots, n\}$ the following sequence is growing with respect to embedding:

$$w_{2,i} \hookrightarrow w_{3,i} \hookrightarrow \cdots.$$

By extracting a subsequence, we can assume that for every $i \in \{0, 1, \dots, n\}$, the above chain is a zigzag between b_i and a_i , for some $b_i, a_i \in M$. This proves the conclusion of the claim. \square

Claim 2.24. *If there is an infinite zigzag between a and b , then there exist $c, c' \in M$ such that $a = cc'$ and $cb = b = bc'$.*

Proof Apply Claim 2.23, yielding monoid elements with

$$\begin{array}{rcccccccc} a & = & & c_1 & & c_2 & & \cdots & & c_n \\ b & = & b_0 & c_1 & b_1 & c_2 & b_2 & \cdots & b_{n-1} & c_n & b_n \\ a & = & a_0 & c_1 & a_1 & c_2 & a_2 & \cdots & a_{n-1} & c_n & a_n \end{array}$$

For every $i \in \{0, \dots, n\}$, we can see that $(b_i, a_i) \leq (b, a)$. If the inclusion is strict, then the induction assumption of the lemma yields $b_i = a_i$. Otherwise, the inclusion is not strict, and therefore

$$(a_i, b_i) = (a, b).$$

If the inclusion is strict for all i , then the third and second rows in the conclusion of Claim 2.23 are equal, thus proving $a = b$, and we are done. Otherwise, there is some $i \in \{0, \dots, n\}$ such that $(b_i, a_i) = (b, a)$. By infix triviality, every interval in the second row that contains i will have product b . It follows that

$$\underbrace{c_j b = b}_{\text{for all } j \leq i} \quad \underbrace{bc_j = b}_{\text{for all } j > i}$$

It is now easy to see that the conclusion of the claim holds if we define c to be the prefix of $a = c_1 \cdots c_n$ that ends with c_i , and define c' to be the remaining suffix. \square

Apply the above claim, and a symmetric one with the roles of a and b swapped, yielding elements c, c', d, d' such that

$$a = cc' \quad cb = b = bc' \quad b = dd' \quad da = a = d'. \quad (2.4)$$

We can now prove the conclusion of the lemma:

$$a \overset{(2.4)}{=} (dc)^!(c'd')! \underset{\text{infix triviality}}{=} (cd)!(d'c')! \overset{(2.4)}{=} b.$$

□

Exercises

Exercise 42. (2) Prove Higman's Lemma.

Exercise 43. (2) Give a polynomial time algorithm, which inputs two non-deterministic automata, and decides if their languages can be separated by a piecewise testable language.

Exercise 44. (1) Consider ω -words, i.e. infinite words of the form

$$a_1 a_2 \cdots \quad \text{where } a_1, a_2, \dots \in \Sigma.$$

Embedding naturally extends to ω -words (in fact, any labelled orders). Show that the embedding on ω -words is also a well-quasi order, i.e. every upward closed set is the upward closure of finitely many elements.

2.5 Two-variable first-order logic

We finish this chapter with one more monoid characterisation of a fragment of first-order logic. A corollary of the equivalence of first-order logic and LTL (or of the equivalence of first-order logic and star-free expressions) is that, over finite words, first-order logic is equivalent to its three variable fragment. What about one or two variables?

It is an easy exercise to show that first-order logic with one variable defines

exactly the languages that are recognised by monoids that are aperiodic and commutative:

$$a^! = a^{!+1} \quad ab = ba \quad \text{for all } a, b.$$

The more interesting case is first-order logic with two variables, which we denote by FO^2 . This logic is characterised in the following theorem.

Theorem 2.25. *For a language $L \subseteq \Sigma^*$, the following are equivalent:*

- (1) *Definable in two variable first-order logic;*
- (2) *Recognised by a finite monoid M with the following property¹⁰: M is a aperiodic, and if an infix class $J \subseteq M$ contains an idempotent, then J is a sub-semigroup of M .*

We use the name DA for the monoids (more generally, finite semigroups) that satisfy the property in item (2). In the exercises, we add several other equivalent conditions for the above theorem, including the temporal logic $\text{LTL}[F, F^{-1}]$ and the following fragment of first-order logic:

$$(\text{definable by a } \exists^* \forall^* \text{-sentence}) \quad \cap \quad (\text{definable by a } \forall^* \exists^* \text{-sentence}).$$

The rest of Section 2.5 is devoted to proving the theorem. We begin with an equational description of DA . (A stronger equational description is given in Exercise 45.)

Lemma 2.26. *A finite semigroup S is in DA if and only if it satisfies:*

$$\underbrace{w^! = w^! v w^!}_{\text{the equality means that}} \quad \text{for all } w, v \in M^* \text{ with } v \leftrightarrow w.$$

the two words have the same product

Proof We first prove that the identity implies that S is DA . Let e be an idempotent. We need to show that if a, b are infix equivalent to e , then the same is true for ab . Because a, b are infixes of e , and e is an idempotent, one can find word w in S^+ which has product e , and where both a and b appear. In particular, $ab \leftrightarrow w$. By the identity in the lemma, we know that $e = eabe$, and therefore ab is an infix of e .

We now show that if S is in DA , then the identity is satisfied. Let $v \leftrightarrow w$ be as in the identity. Let e be the product of $w^!$, and let J be the infix class of e . This infix class is a semigroup, by definition of DA . For every letter a that appears in the word w , there is a suffix of $w^! w^!$ which begins with a and has product in

¹⁰ This property appears in [34] Schützenberger, “Sur Le Produit De Concatenation Non Ambigu”, 1976 where it is used to characterise certain unambiguous regular expressions, see Exercise 49.

J . Let $a' \in J$ be the product of this suffix. Since J is a semigroup, it follows that $ea' \in J$ and therefore also $ea \in J$. Since $ea \in J$ holds for every letter that appears in w , it follows that $eve \in J$. This means that eve is in the \mathcal{H} -class of e , and therefore $e = eve$ by aperiodicity (which is part of the definition of DA), thus establishing the identity. \square

We now prove the theorem.

To prove the implication (1) \Rightarrow (2), we show that for every language definable in FO^2 , its syntactic monoid belongs to DA . By Lemma 2.26 and unravelling the definition of the syntactic monoid, it is enough to show that for every $w_1, w_2, v, w \in \Sigma^*$ and $n \in \{0, 1, \dots\}$, if $v \hookrightarrow w$ then the words

$$w_1 w^n w_2 \quad w_1 w^n v w^n w_2$$

satisfy the same FO^2 sentences of quantifier rank at most n . This is shown using a simple Ehrenfeucht-Fraïssé argument.

For the implication (2) \Rightarrow (1), we use the following lemma.

Lemma 2.27. *Let M be a monoid in DA , and let $a_1, a_2 \in M$. Then*

$$w \in M^* \quad \mapsto \quad \underbrace{a_1 \cdot (\text{product of } w) \cdot a_2}_{\in M}$$

is a colouring definable in FO^2 , i.e. every inverse image is definable in FO^2 .

If we apply the above lemma to a_1 and a_2 being the monoid identity, we conclude that the product operation is definable in FO^2 . This implies that every language recognised by the monoid is definable in FO^2 , thus proving the implication (1) \Leftarrow (2) in the theorem.

Proof Induction on the following parameters, ordered lexicographically:

- (1) size of M ;
- (2) number of elements that properly extend a_1 in the prefix ordering;
- (3) number of elements that properly extend a_2 in the suffix ordering.

The induction base is when M has one element, in which case the colouring in the lemma is constant, and therefore definable in FO^2 .

Let us also prove another variant of the induction base, namely when induction parameters (2) and (3) are zero, which means that a_1 is maximal in the prefix ordering and a_2 is maximal in the suffix ordering. It follows that, for

$$\mathcal{H}\text{-class of } a_1 a a_2 = \mathcal{H}\text{-class of } a_1 b a_2 \quad \text{for all } a, b \in M.$$

Since DA implies aperiodicity, which implies \mathcal{H} -triviality, the colouring in the statement of the lemma is constant, and therefore definable in FO^2 .

It remains to prove the induction step. Because of the two kinds of induction base that were considered above, we can assume that one of the parameters (2) or (3) is nonzero. By symmetry, assume that a_1 is not maximal in the prefix ordering.

Claim 2.28. *For every $a \in M$, the following is a sub-monoid of M :*

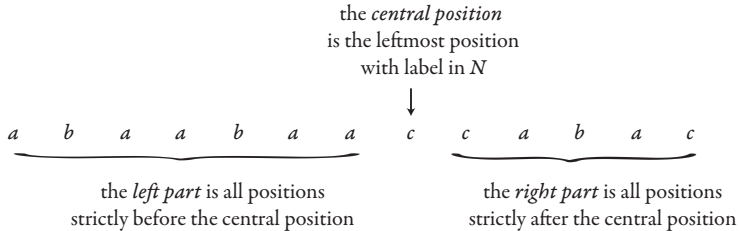
$$\underbrace{\{b \in M : ab \text{ is prefix equivalent to } a\}}_{\text{we call this set the prefix stabiliser of } a}$$

Proof The prefix stabiliser clearly contains the monoid identity. It remains to show that it is closed under composition. Let b, c be in the prefix stabiliser of a . Using the definition of the prefix stabiliser, it is easy to construct a word $w \in M^*$, such that $bc \hookrightarrow w$ and $aw = a$. By Lemma 2.26, it follows that

$$a = aw = aw^1 = aw^1bcw^1 = abcw^1,$$

which establishes that bc is in the prefix stabiliser of a . □

Let $N \subseteq M$ be the prefix stabiliser of a_1 ; our assumption says that N is a proper subset of M , and by the above claim it is also a sub-monoid. We decompose a word $w \in M^*$ into three parts, as explained in the following picture:



There is an FO^2 formula which selects the central position. Since all labels in the left part are from N , we can use the induction assumption on a smaller monoid to prove that the colouring

$$w \quad \mapsto \quad \text{product of left part}$$

is definable in FO^2 . (When using the induction assumption, we restrict all quantifiers of the formulas from the induction assumption so that they quantify over positions that are to the left of the central position.) Let c be the product of the prefix up to and including the central position; as we have shown above, this product can be computed in FO^2 . By definition of the central position, we know

that a_1 is a proper prefix of a_1c , and therefore we can use the induction assumption to prove that

$$w \mapsto a_1c \cdot (\text{product of right part}) \cdot a_2$$

is a colouring definable in FO^2 . The conclusion of the lemma follows. \square

Exercises

Exercise 45. (2) Show that a semigroup belongs to DA if and only if it satisfies the identity

$$(ab)^! = (ab)^!a(ab)^! \quad \text{for all } a, b.$$

Exercise 46. (2) Show that FO^2 has the same expressive power as $\text{LTL}[F, F^{-1}]$, which is the extension of $\text{LTL}[F]$ with the following past operator:

$$w, x \models F^{-1}\varphi \stackrel{\text{def}}{=} \exists y y < x \wedge w, y \models \varphi.$$

Exercise 47. (3) Define the *syntactic ordering* on the syntactic monoid, which depends on the accepting set F , as follows:

$$a \leq b \stackrel{\text{def}}{=} \forall x, y \in M \ xay \in F \Rightarrow xby \in F.$$

Show that a language can be defined by a first-order sentence of the form

$$\underbrace{\exists x_1 \cdots \exists x_n \forall y_1 \cdots \forall y_m}_{\text{such a formula is called an } \exists^* \forall^* \text{-sentence}} \overbrace{\varphi(x_1, \dots, x_n, y_1, \dots, y_m)}^{\text{quantifier-free}}$$

if and only if

$$w^! \leq w^!vw^! \quad \text{for all } \underbrace{v \hookrightarrow w}_{\text{Higman ordering}}$$

Hint¹¹: use Exercise 26.

¹¹ An effective characterisation of $\exists^* \forall^*$ -sentence was first given in [2] Arfi, “Polynomial Operations on Rational Languages”, 1987, Theorem 3. The proof was simplified in [27] Pin and Weil, “Polynomial closure and unambiguous product”, 1997, Theorem 5.8. The solution which uses Exercise 26 is based on [27]. Characterisations of fragments of first-order logic such as $\exists^* \forall^*$ are widely studied, see [28] Place and Zeitoun, “Going Higher in First-Order Quantifier Alternation Hierarchies on Words”, 2019

Exercise 48. (3) Show that L is definable in FO^2 if and only if both L and its complement can be defined using $\exists^*\forall^*$ -sentences.

Exercise 49. (2) We say that a regular expression

$$\Sigma_0^* a_1 \Sigma_1^* \cdots \Sigma_{n-1}^* a_n \Sigma_n^*$$

is *unambiguous* if every word w admits at most one factorisation

$$w = w_0 a_1 w_1 \cdots w_{n-1} a_n w_n \quad \text{where } w_i \in \Sigma_i^* \text{ for all } i \in \{1, \dots, n\}.$$

Show that a language is a finite disjoint union of unambiguous expressions if and only if its syntactic monoid of L is in DA^{12} .

¹² This exercise is based on

[34] Schützenberger, “Sur Le Produit De Concatenation Non Ambigu”, 1976

3

Infinite words

In this chapter, we study infinite words.

In Section 3.1, we begin with the classical model of infinite words, namely ω -words. In ω -word, the positions are ordered like the natural numbers. We show how the structure of finite semigroups, which was developed in Section 1.2, can be applied to prove McNaughton's Theorem about determinisation of ω -automata.

In Section 3.2, we move to more general infinite words, where the positions can be any countable linear order, e.g. the rational numbers. For this kind of infinite words, we define a suitable generalisation of semigroups, and show that it has the same expressive power as monadic second-order logic.

3.1 Determinisation of Büchi automata for ω -words

An ω -word is a function from the natural numbers to some alphabet Σ . We write Σ^ω for the set of all ω -words over alphabet Σ . To recognise properties of ω -words, we use Büchi automata.

Definition 3.1 (Büchi automata). The syntax of a *nondeterministic Büchi automaton* is the same as the syntax of a nondeterministic finite automaton for finite words, namely it consists of:

$$\underbrace{Q}_{\text{states}} \quad \underbrace{\Sigma}_{\text{input alphabet}} \quad \underbrace{I, F \subseteq Q}_{\text{initial and final states}} \quad \underbrace{\delta \subseteq Q \times \Sigma \times Q}_{\text{transition relation}}.$$

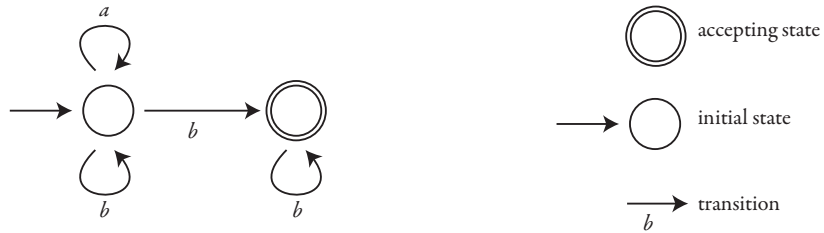
The difference, with respect to nondeterministic automata on finite words, is in the semantics: a word in Σ^ω is accepted by the automaton if there exists a run which begins in an initial state, and which satisfies the *Büchi condition*:

some accepting state appears infinitely often in the run.

A *deterministic Büchi automaton* is the special case when there is one initial state, and the transition relation is a function from $Q \times \Sigma$ to Q .

The following example shows that deterministic Büchi automata are weaker than than nondeterministic ones.

Example 8. Consider the language of ω -words over alphabet $\{a, b\}$ where letter a appears finitely often. This language is recognised by a nondeterministic Büchi automaton as in the following picture:



The idea is that the automaton nondeterministically guesses some position which will not be followed by any a letters; this guess corresponds to the horizontal transition with label b in the picture.

This language is not recognised by any deterministic Büchi automaton. Toward a contradiction, imagine a hypothetical deterministic Büchi automaton which recognises the language. Run this automaton on b^ω . Since a appears finitely often in this ω -word, the corresponding run (unique by determinism) must use an accepting state in some finite prefix. Extend that finite prefix by appending ab^ω . Again, the word must be accepted, so an accepting state must be eventually visited after the first a . By repeating this argument, we get a word which has infinitely many a 's and where the (unique) run of the deterministic automaton sees accepting states infinitely often; a contradiction. \square

The above shows that languages recognised by deterministic Büchi automata are not closed under Boolean combinations. This turns out to be the only limitation of the model, as shown in the following theorem.

Theorem 3.2. *If a language of ω -words is recognised by a nondeterministic Büchi automaton, then it is a Boolean combination of languages recognised by deterministic Büchi automata¹.*

¹ A Boolean combination of deterministic Büchi automata is the same thing as what is known as a *deterministic Muller automaton*. Therefore, the theorem is the same McNaughton's Theorem,

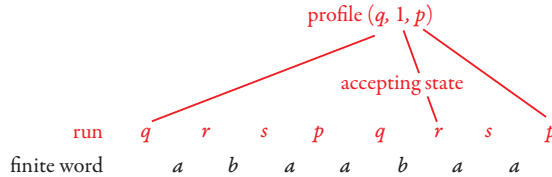
The converse implication in the above theorem is also true, which is left as an exercise for the reader (Exercise 51). A language is called ω -regular if it is recognised by a nondeterministic Büchi automaton, or equivalently, by a Boolean combination of deterministic Büchi automata. The ω -regular languages are closed under Boolean combination thanks to the deterministic characterisation. The original application of Büchi automata was Büchi's proof² that they recognise exactly the same languages of ω -words as monadic second-order logic; this application is a simple corollary of Theorem 3.2, see Exercise 55.

There are several combinatorial proofs for the determinisation result in Theorem 3.2. In this section, we present an algebraic proof, which leverages the structural theory of finite semigroups developed in Section 1.2.

Let \mathcal{A} be a nondeterministic Büchi automaton, with states Q and input alphabet Σ . For an ω -word, define its ω -type to be the set of states from which the word is accepted. We also define the type for finite words, but here we need to store a bit more information. For a run of the automaton over a finite word, define the *profile* of the run to be the triple (q, i, p) where q is the source state of the run, p is the target state of the run, and

$$i = \begin{cases} 0 & \text{if the run does not use any accepting state} \\ 1 & \text{if the run uses some accepting state.} \end{cases}$$

Here is a picture of a run with its profile:



Define the *type* of a finite word $w \in \Sigma^+$ to be the set of profiles of runs over this word. It is not hard to see that the function

$$w \in \Sigma^+ \quad \mapsto \quad \text{type of } w \in \underbrace{\mathbf{P}(Q \times \{0, 1\} \times Q)}_S$$

[25] McNaughton, "Testing and generating infinite sequences by a finite automaton", 1966, p. 524

which says that nondeterministic Büchi automata can be determinised into deterministic Muller automata.

² [8] Büchi, "On a decision method in restricted second order arithmetic", 1962

is a semigroup homomorphism, with a naturally defined semigroup structure on S . The following lemma shows that types and ω -types are compatible.

Lemma 3.3. *If $w_i \in \Sigma^+$ and $v_i \in \Sigma^+$ have the same type for every $i \in \{1, 2, \dots\}$, then $w_1 w_2 \dots \in \Sigma^\omega$ and $v_1 v_2 \dots \in \Sigma^\omega$ have the same ω -type.*

Proof By substituting parts of an accepting run, while preserving the Büchi condition. \square

Thanks to the above lemma, it makes sense to talk about the ω -type of a word $w \in S^\omega$ built out of types. In particular, it makes sense to say whether or not a word $w \in S^\omega$ is accepted by \mathcal{A} , since this information is stored in the type. A special case of this notation is ae^ω , where $a, e \in S$, which is the ω -type of the ω -word that begins with letter a and has all other letters equal to e . The importance of this special case is explained by the following lemma about factorisations of ω -words³

Lemma 3.4. *For every $w \in S^\omega$ there exist $a, e \in S$, such that e is an idempotent, $ae = e$, and there is a factorisation*

$$w = \underbrace{\quad}_{w_0}^{\text{type } a} \underbrace{\quad}_{w_1}^{\text{type } e} \underbrace{\quad}_{w_2}^{\text{type } e} \underbrace{\quad}_{w_3}^{\text{type } e} \dots$$

Proof Define a *cut* in w to be the space between two positions. Consider an undirected edge-labelled graph, defined as follows. Vertices are cuts. For every two distinct cuts, there is an undirected edge, labelled by the type of the finite word that connects the two cuts. By Ramsey's Theorem A, see Exercise 50, there exists a type $a \in S$ and an infinite set X of vertices, such every two distinct vertices from X are connected by an edge with label e . Define the decomposition from the lemma to be the result of cutting w along all cuts from X . By assumption on X , every word w_i with $i > 0$ has type e . Idempotence of e follows from

$$\underbrace{\underbrace{\quad}_{w_i}^e \quad \underbrace{\quad}_{w_{i+1}}^e}_e.$$

Finally, we can assure that $ae = a$ by joining the first two groups. \square

A corollary of Lemmas 3.3 and 3.4 is that $w \in L$ if and only if

(*) there is a factorisation as in Lemma 3.4 such that $ae^\omega \in L$.

³ This lemma was first observed by Büchi in [8, Lemma 1] where it was used to prove that nondeterministic Büchi automata are closed under complementation, without passing through a deterministic model.

So far, we are doing the same argument as in Büchi's original complementation proof from [8]. In his proof, Büchi observed that variant of (*) with $ae^\omega \notin L$, which characterises the complement of L , can be expressed by a nondeterministic Büchi automaton, and therefore nondeterministic Büchi automata are closed under complementation.

Now, we diverge from Büchi's proof, since we are interested in determinisation and not complementation. Here, semigroups will be helpful. Since it is not clear how to express condition (*) using a deterministic Büchi automaton, we will reformulate it. This is done in the following lemma. In the lemma, we say that a pair $(a, b) \in S^2$ appears infinitely often in an ω -word $w \in \Sigma^\omega$ if for every $n \in \{1, 2, \dots\}$ one can find a factorisation

$$w = \underbrace{x}_{\text{type } a} \underbrace{y}_{\text{type } b} z$$

such that x has length at least n .

Lemma 3.5. *An ω -word $w \in \Sigma^\omega$ is accepted by \mathcal{A} if and only if*

(**) *there exist $a, e \in S$, with e idempotent, $ae = e$, and $ae^\omega \in L$, such that all of the following conditions are satisfied:*

- (1) *(a, e) appears infinitely often; and*
- (2) *if (b, c) appears infinitely often, then c is an infix of e .*

Proof The top-down implication, which says that every word accepted by \mathcal{A} must satisfy (**), is an immediate consequence of Lemma 3.4. We are left with the bottom-up implication. Suppose that w satisfies (**), as witnessed by $a, e \in S$. By condition (1), there is a decomposition

$$w = w_1 v_1 w_2 v_2 w_3 v_3 \cdots$$

such that for every $i \in \{1, 2, \dots\}$ the word v_i has type e and the word $w_1 v_1 \cdots w_i v_i$ has type a . Let a_i be the type of w_i . The ω -type of w is equal to

$$a_1 e a_2 e a_3 e \cdots,$$

which by Lemma 3.3 and idempotence of e is equal to the ω -type of

$$a_1 e \underbrace{ea_2 e}_{\text{call this } g_2} \underbrace{ea_3 e}_{\text{call this } g_3} \cdots,$$

By condition (2), there is some n such that ea_i is an infix of e for all $i \geq n$. Therefore, g_i is an infix of e for $i \geq n$. Since g_i begins and ends with e , it follows that g is in the \mathcal{H} -class of e for all $i > n$. Since this \mathcal{H} -class, call it G ,

contains the idempotent e , it must be a group by the \mathcal{H} -class lemma. We now have

$$\begin{aligned}
\omega\text{-type of } w &= (a_1 e a_2 \cdots a_{n-1} e_{n-1} = a) \\
a g_n g_{n+1} g_{n+2} \cdots &= (\text{by Lemma 3.4, for some } g, f \in G) \\
a g f^\omega &= (\text{because } e \text{ is the unique idempotent in } G) \\
a g e^\omega &= (\text{some power of } g \text{ is the idempotent } e) \\
a g^\omega &= (\text{for the same reason}) \\
a e^\omega &
\end{aligned}$$

and therefore w must belong to L . □

To finish the determinisation construction in Theorem 3.2, it remains to show that condition (***) from the above lemma is a finite Boolean combination of languages recognised by deterministic Büchi automata. This will follow from the following lemma.

Lemma 3.6. *For every $a, e \in S$ the property “ (a, e) appears infinitely often” is recognised by a deterministic Büchi automaton.*

Proof Let $L \subseteq \Sigma^*$ be the set of words which can be decomposed as

$$w = \underbrace{u}_{\text{type } b} \underbrace{v}_{\text{type } e} \quad \text{for some } b \in S \text{ such that } aeb = a.$$

This is easily seen to be a regular language, and hence it is recognised by some finite deterministic automaton \mathcal{D} . The deterministic Büchi automaton \mathcal{B} recognising the property in the statement of the lemma is defined as follows. Its space is the disjoint union of the set of types S and the states of \mathcal{D} . The initial state is the type in S of the empty word. The automaton \mathcal{B} begins to read input letters, keeping in its state the type of the prefix read so far in its state, until the prefix has type ae . Then it switches to the initial state q_0 of the automaton \mathcal{D} . For states of \mathcal{D} , the state update function of \mathcal{B} is as follows:

$$\delta_{\mathcal{B}}(q, \sigma) \mapsto \begin{cases} \delta_{\mathcal{A}}(q, \sigma) & \text{if } q \text{ is not accepting} \\ \delta_{\mathcal{A}}(q_0, \sigma) & \text{otherwise.} \end{cases}$$

The Büchi accepting states of \mathcal{B} are the same as in \mathcal{D} . □

This completes the proof of Theorem 3.2.

Semigroups for ω -words. There is an implicit algebraic structure in the proof of Theorem 3.2, which is formalised in the following definition.

Definition 3.7. An ω -semigroup consists of:

- two sets S_+ and S_ω , called the *finite sort* and the *ω -sort*, respectively.
- a finite product $\pi_+ : (S_+)^+ \rightarrow S_+$, associative as in the sense of semigroups;
- an ω -product $\pi_\omega : (S_+)^{\omega} \rightarrow S_\omega$, associative in the following sense:

$$\pi_\omega(w_1 w_2 \cdots) = \pi_\omega(\pi_+(w_1) \pi_+(w_2) \cdots) \quad \text{for every } w_1, w_2, \dots \in S_+.$$

An example of an ω -semigroup is the automaton types that were used in the proof of Theorem 3.2. Another example is the *free ω -semigroup over a set Σ* , where the finite sort is Σ^+ , the ω -sort is Σ^ω , and the two products are defined in the natural way. The same proof as in Theorem 3.2 shows that a language is ω -regular if and only if it is recognised by a homomorphism into an ω -semigroup which is finite (on both sorts). This is discussed in more detail in some of the exercises at the end of this section.

The associativity axiom on the ω -product can be represented using a commuting diagram, in the same spirit as for Lemma 1.4:

$$\begin{array}{ccc} ((S_+)^+)^{\omega} & \xrightarrow{\omega\text{-product in free } \omega\text{-semigroup over } S_+} & (S_+)^{\omega} \\ (\pi_+)^{\omega} \downarrow & & \downarrow \pi_\omega \\ (S_+)^{\omega} & \xrightarrow{\pi_\omega} & S_\omega \end{array}$$

In the above diagram, $(\pi_+)^{\omega}$ denotes the coordinate-wise lifting of π_+ to ω -words of finite words.

Exercises

Exercise 50. (2) Prove the following result, called *Ramsey's Theorem A*⁴. Consider an infinite undirected graph, where every two distinct vertices are a connected by an edge that is labelled by one of finitely many colours. Then the graph contains an infinite monochromatic clique, which means that there exists a colour e and an infinite set X of vertices, such that every two distinct vertices from X are connected by an edge with colour e .

Exercise 51. (1) Prove the converse implication in Theorem 3.2.

⁴ [29] Ramsey, "On a problem of formal logic", 1929, Theorem A

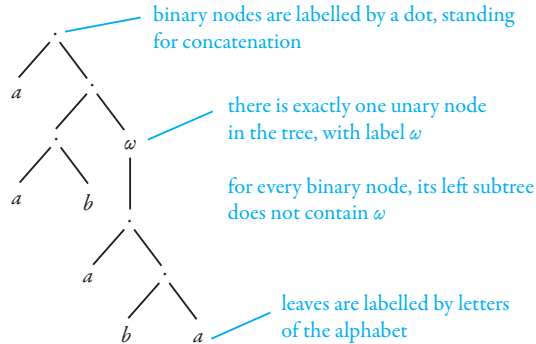
Exercise 52. (1) We say that an ω -word is *ultimately periodic* if it has the form wu^ω , for some finite words $w, u \in \Sigma^\omega$. Show that every nonempty ω -regular language contains an ultimately periodic ω -word.

Exercise 53. (1) Show that two ω -regular languages are equal if and only if they contain the same ultimately periodic ω -words.

Exercise 54. (1) Show that an ω -word w is ultimately periodic if and only if $\{w\}$ is an ω -regular language.

Exercise 55. (2) To an ω -word we associate an ordered model, in the same way as for finite words. Show that a language is mso definable (using the ordered model) if and only if it is ω -regular.

Exercise 56. (2) Define an ω -term to be any tree as in the following picture:



Every ω -term represents some ultimately periodic ω -word, but several ω -terms might represent the same ultimately periodic ω -word. Show that two ω -terms represent the same ultimately periodic ω -word if and only if one can be transformed into the other using the equations:

$$(xy)z = x(yz) \quad (xy)^\omega = x(yx)^\omega \quad \underbrace{(x^n)^\omega = x^\omega}_{\text{for every } n \in \{1, 2, \dots\}}$$

where x, y, z stand for ω -terms.

Exercise 57. (1) Let $L \subseteq \Sigma^\omega$. Consider the following equivalence relations on Σ^+ .

- Right equivalence is defined by

$$w \sim w' \stackrel{\text{def}}{=} wv \in L \Leftrightarrow w'v \in L \text{ for every } v \in \Sigma^\omega.$$

- Two-sided congruence is defined by

$$w \sim w' \stackrel{\text{def}}{=} uvw \in L \Leftrightarrow uw'v \in L \text{ for every } u \in \Sigma^*, v \in \Sigma^\omega.$$

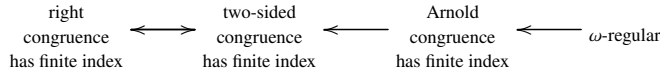
- Arnold congruence is defined by

$$w \sim w' \stackrel{\text{def}}{=} \bigwedge \begin{cases} u(wv)^\omega \in L \Leftrightarrow u(w'v)^\omega \in L & \text{for every } u, v \in \Sigma^*. \\ uwv \in L \Leftrightarrow uw'v \in L & \text{for every } u \in \Sigma^*, v \in \Sigma^\omega. \end{cases}$$

Show that the latter two, but not necessarily the first one, are semigroup congruences, i.e. they satisfy

$$\bigwedge_{i \in \{1,2\}} w_i \sim w'_i \quad \text{implies} \quad w_1 w_2 \sim w'_1 w'_2.$$

Exercise 58. (2) Consider the equivalence relations defined in Exercise 57. Prove that the arrows in the following diagram are true implications, and provide counter-examples the missing arrows:



Exercise 59. (2) Define the *Arnold semigroup* of a language $L \subseteq \Sigma^\omega$ to be the quotient of Σ^+ under Arnold congruence. Let $L \subseteq \Sigma^\omega$ be a ω -regular. Show that L is definable in first-order logic if and only if its Arnold semigroup is aperiodic.

Exercise 60. (2) The temporal logic $\text{LTL}[F]$ can also be used to define languages of ω -words. Let $L \subseteq \Sigma^\omega$ be a ω -regular. Show that L is definable in LTL if and only if its Arnold semigroup is suffix-trivial.

Exercise 61. (1) Show an ω -regular language where the Arnold semigroup is infix trivial, but which cannot be defined by a Boolean combination of \exists^* -sentences.

Exercise 62. (1) Define a *safety automaton* to be an automaton on ω -words with the following acceptance condition: all states in the run are accepting.

Show that deterministic and nondeterministic safety automata recognise the same languages.

Exercise 63. (1) Show that an ω -regular language of ω -words is recognised by a safety automaton (deterministic or nondeterministic, does not matter by Exercise 62) if and only if

$$uw^!v \in L \Leftrightarrow u(w^!)^\omega \in L \quad \text{for every } u, w \in \Sigma^+ \text{ and } v \in \Sigma^\omega,$$

where $! \in \{1, 2, \dots\}$ is the exponent obtained from the Idempotent Power Lemma as applied to the Arnold semigroup of L .

Exercise 64. (2) For a finite alphabet Σ , we can view Σ^ω as metric space, where the distance between two different ω -words is defined to be

$$\frac{1}{2^{(\text{length of longest common prefix})}}$$

This is indeed a distance, i.e. it satisfies the triangle inequality. Let $L \subseteq \Sigma^\omega$ be ω -regular. Show that L is recognised by a safety automaton if and only if it is a closed set with respect to this distance.

Exercise 65. (2) Find a condition on the Arnold semigroup of an ω -regular language which characterises the clopen languages (i.e. languages which are both closed and open with respect to the distance from Exercise 64)

Exercise 66. (1) We use the topology from Exercise 64. Define a G_δ set to be any countable intersection of open sets. Show that every ω -regular language is a finite Boolean combination of G_δ sets.

Exercise 67. (2) Let $L \subseteq \Sigma^\omega$ be an ω -regular language, and define $!$ as in Exercise 62. Show that L is recognised by a deterministic Büchi automaton if and only if:

$$u(wv^!)^!v^\omega \in L \Rightarrow u(wv^!)^\omega \in L \quad \text{for every } u, w, v \in \Sigma^+.$$

Exercise 68. (2) Let $L \subseteq \Sigma^\omega$. Define an ω -congruence to be any equivalence relation \sim on Σ^+ which is a semigroup congruence and which satisfies

$$\bigwedge_{i \in \{1, 2, \dots\}} w_i \sim w'_i \quad \text{implies} \quad w_1 w_2 \cdots \in L \Leftrightarrow w'_1 w'_2 \cdots \in L. \quad (3.1)$$

Show that a language is ω -regular if and only if it has an ω -congruence of finite index.

Exercise 69. (2) Define *semi- ω -congruence* for a language $L \subseteq \Sigma^\omega$ to be an equivalence relation on finite words which satisfies (3.1), but which is not necessarily a semigroup congruence. Show that if there is a semi- ω -congruence of finite index, then there is an ω -congruence of finite index.

Exercise 70. (2) We say that \sim is the *syntactic ω -congruence* of $L \subseteq \Sigma^\omega$ if it is an ω -congruence, and every other ω -congruence for L refines \sim . Show that if a language is ω -regular, then it has a syntactic ω -congruence, which is equal to the Arnold congruence.

Exercise 71. (2) Show a language of ω -words which does not have a syntactic ω -congruence.

3.2 Countable words and \circ -semigroups

In this section, we move to \circ -words. These are words where the set positions is a countable linear order. The positions could be some finite linear order, as in finite words, or the natural numbers, as in ω -words, but the positions could also be dense, as in the rational numbers. One advantage of \circ -words, as compared to ω -words, is that they can be concatenated, which is useful when defining the corresponding generalisation of semigroups.

For finite words, as well as for ω -words, the approach via semigroups can be seen as an alternative to existing automata models. This is no longer the case for \circ -words. There is no known corresponding automaton model, and therefore \circ -semigroups are the only known model of recognisability.

Definition 3.8 (\circ -words). A Σ -labelled linear order consists of a set X of *positions*, equipped with a total order and a labelling of type $X \rightarrow \Sigma$. Two such objects are considered *isomorphic* if there is a bijection between their positions, which preserves the order and labelling. Define a \circ -word over Σ to be

any isomorphism class of countable⁵ Σ -labelled linear orders. We write Σ° for the set of \circ -words⁶.

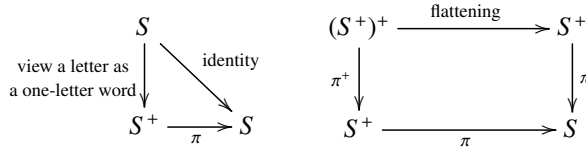
Every finite word is a \circ -word, likewise for every ω -word. Another example is labelled countable ordinals, e.g. any \circ -word where the positions are $\omega + \omega$. Below is a more fancy example, which uses a dense set of positions.

Example 3.9 (Shuffles). A classical exercise on linear orders is that the rational numbers are the unique – up to isomorphism – countable linear order which is dense and has no endpoints (i.e. neither a least nor greatest element). This is proved by constructing, using the back-and-forth method, an isomorphism between any two such orders. The same argument shows that for every countable Σ there is a \circ -word over Σ which has no endpoints, and which satisfies

$$\bigwedge_{a \in \Sigma} \underbrace{\forall x \forall y \exists z \quad x < z < y \wedge a(z)}_{\text{label } a \text{ is dense}}.$$

We use the name *shuffle of Σ* for the above \circ -word. Shuffles will play an important role in semigroups for \circ -words.

We now define the generalisation of semigroups for \circ -words. We use the approach to associativity via commuting diagrams that was described in Lemma 1.4. Recall from that lemma that a semigroup product on a set S could be defined as any operation π which makes the following diagram commute:



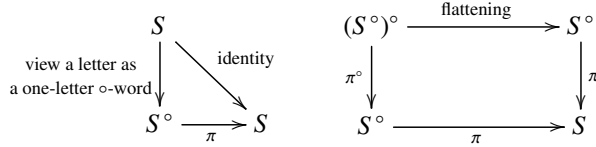
For \circ -semigroups, we take the same approach. For a set S , the flattening operation $(S^\circ)^+ \rightarrow S^\circ$ is defined in the natural way, by replacing each position with the \circ -word that is in its label (a formal definition uses a lexicographic product of labelled linear orders).

⁵ The reader might wonder why we assume countability. The reason is that the decidability theory that will be described in this section breaks down for uncountable linear orders. In fact, the mso theory of the order of real numbers $(\mathbb{R}, <)$ is undecidable, as shown [35] Shelah, “The Monadic Theory of Order”, 1975, Theorem 7

The description of \circ -semigroups in this section is based on [10] Carton, Colcombet, and Puppis, “An algebraic approach to MSO-definability on countable linear orders”, 2018 which itself is based on Shelah’s approach to countable linear orders from [35].

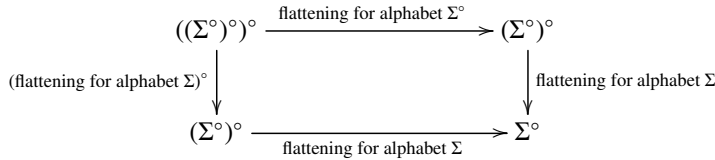
⁶ Formally speaking, this is not a set, because the linear orders form a class and not a set. However, without loss of generality we can use some fixed countably infinite set, e.g. the natural numbers, for the positions (but the order need not be the same as in the natural numbers). Under this restriction, the labelled linear orders become a set, and no isomorphism types are lost. For this reason, we can refer to Σ° as a set.

Definition 3.10. A \circ -semigroup consists of an underlying set S equipped with a product operation $\pi : S^\circ \rightarrow S$, which is associative in the sense that the following two diagrams commute:



In the above diagram, π° denotes the coordinate-wise lifting of π to \circ -words of \circ -words.

Example 9. The *free* \circ -semigroup over alphabet Σ has Σ° as its underlying set, and its product operation is flattening. To check that this product operation is associative, one needs to prove that the following diagram commutes:



To prove this formally, one uses the formal definition of flattening, in terms of lexicographic products of linear orders (see Example 19). This \circ -semigroup is called *free* for the usual reasons; a more formal description of these usual reasons will appear later in the book, when discussing monads. \square

Example 10. Recall the semigroups of size two that were discussed in Example 1.2:

$$\underbrace{(\{0, 1\}, +)}_{\text{addition mod 2}} \quad (\{0, 1\}, \min) \quad \underbrace{(\{0, 1\}, \pi_1)}_{\text{product } ab \text{ is } a} \quad \underbrace{(\{0, 1\}, \pi_2)}_{\text{product } ab \text{ is } b} \quad \underbrace{(\{0, 1\}, (a, b) \mapsto 1)}_{\text{all products are 1}}$$

Which ones can be extended to \circ -semigroups in at least one way?

The first example, i.e. the two-element group, cannot be extended in any way, because the product a of the ω -word 1^ω would need satisfy

$$a = \pi(1^\omega) = \pi(\pi(1)\pi(1^\omega)) = \pi(1a) = 1 + a.$$

The remaining semigroups can be extended to \circ -semigroups. As we will see in Example 11, the extensions are not necessarily unique. \square

We use \circ -semigroups to recognise languages of \circ -words. Define a *homomorphism of \circ -semigroups* to be a function h which makes the following diagram

commute:

$$\begin{array}{ccc} S^\circ & \xrightarrow{h^\circ} & T^\circ \\ \text{product in } S \downarrow & & \downarrow \text{product in } T \\ S & \xrightarrow{h} & T \end{array}$$

Like for semigroups, homomorphisms of \circ -semigroup can be described in terms of compositional functions. Suppose that S is a \circ -semigroup and T is a set. We say that a function $h : S \rightarrow T$ is *compositional* if there exists a function $\pi : T^\circ \rightarrow T$ which makes the following diagram commute

$$\begin{array}{ccc} S^\circ & \xrightarrow{h^\circ} & T^\circ \\ \text{product in } S \downarrow & & \downarrow \pi \\ S & \xrightarrow{h} & T \end{array}$$

Using the same proof as for Lemma 1.3, one shows that if h is a compositional and surjective, then π is necessarily associative, thus turning T into a \circ -semigroup, and furthermore h is a homomorphism.

We say that a language $L \subseteq \Sigma^\circ$ is *recognised* by a \circ -semigroup S if there is a homomorphism $h : \Sigma^\circ \rightarrow S$ which recognises it, i.e.

$$h(w) = h(w') \quad \text{implies} \quad w \in L \Leftrightarrow w' \in L \quad \text{for every } w, w' \in L.$$

We are mainly interested in languages recognised by finite \circ -semigroups.

Example 11. Consider un-labelled countable linear orders, which can be viewed as \circ -words over a one letter alphabet $\{a\}$. Consider the function

$$h : \{a\}^\circ \rightarrow \{0, 1\}$$

which sends well-founded \circ -words to 1, and the remaining \circ -words to 0. We claim that h is compositional (and therefore the language of well-founded \circ -words is recognised by a finite \circ -semigroup). Indeed, take some $v \in (\{a\}^\circ)^\circ$ which flattens to $w \in \{a\}^\circ$. To prove compositionality, need to show that $h^\circ(v)$ uniquely determines $h(w)$. This is because $h(w) = 1$ if and only if the positions of v are well-founded, and every such a position is labelled by a well-founded order. All of this information can be recovered from $h^\circ(v)$. The compositional function h induces an underlying structure of a \circ -semigroup on $\{0, 1\}$. When restricted to finite products, this \circ -semigroup is the same as $(\{0, 1\}, \min)$. Note that a symmetric \circ -semigroup can be constructed, for orders which are well-founded after reversing. The symmetric \circ -semigroup also coincides with $(\{0, 1\}, \min)$ on finite words. \square

Example 12. Consider the language $L \subseteq \{a, b, 1\}^\circ$, which contains \circ -words where some position with label a is to the left of some position with label b . Consider the following function

$$w \in \{a, b, 1\}^\circ \mapsto \begin{cases} 0 & \text{if } w \in L \\ 1 & \text{if all letters are 1} \\ b & \text{if all letters are } b \text{ or 1, and there is some } b \\ ba & \text{if both } b \text{ and } a \text{ appear, but } w \notin L \\ a & \text{otherwise} \end{cases}$$

This function is easily seen to be compositional, and therefore its image is a \circ -semigroup. The element 0 is absorbing, and 1 is a monoid identity. The language L is therefore recognised by the corresponding \circ -semigroup. \square

Monadic second-order logic on \circ -words. The *ordered model* of a \circ -word is defined in the same way as for finite words: the universe is the positions, and the relations and their meaning are the same as for finite words. We say that a language $L \subseteq \Sigma^\circ$ is definable in mso if there is an mso sentence φ , using the vocabulary of the ordered model, such that

$$w \in L \iff \text{the ordered model of } w \text{ satisfies } \varphi \quad \text{for every } w \in \Sigma^\circ.$$

Example 13. Consider the language of well-founded \circ -words that was discussed in Example 11. This language is definable in mso, by simply writing in mso the definition of well-foundedness:

$$\underbrace{\forall X}_{\substack{\text{for every} \\ \text{set of} \\ \text{positions}}} \left(\underbrace{(\exists x \in X)}_{\text{which is nonempty}} \Rightarrow \underbrace{(\exists x \in X \forall y \in X x \leq y)}_{\text{there is a least position}} \right).$$

Another example is the \circ -words which contain a sub-order that is dense:

$$\underbrace{\exists X}_{\substack{\text{exists a} \\ \text{set of} \\ \text{positions}}} \left(\underbrace{(\exists x \in X)}_{\text{which is nonempty}} \wedge \underbrace{(\forall x \in X \forall y \in Y x < y \Rightarrow \exists z \in X x < z < y)}_{\text{and dense in itself}} \right).$$

An \circ -word which violates the second property, i.e. it does not have any dense sub-order, is called *scattered*. \square

Once we have built up all the necessary ideas in the Trakhtenbrot-Büchi-Elgot Theorem for finite words, it is very easy to get the extension for \circ -words. The same proof as for finite words (using a powerset construction on \circ -semigroups) gives the following result.

Lemma 3.11. *If a language $L \subseteq \Sigma^\circ$ is definable in mso, then it is recognised by a finite \circ -semigroup.*

The above lemma seems too easy. Is there a catch? Yes: the lemma does not give any algorithm for deciding if an mso definable language is empty. In the case of finite words, we could remark that all of the constructions used in the proof (products and powersets) are effective, with finite semigroups represented by their multiplication tables. So far, we do not have any finite representation of \circ -semigroups yet, and therefore we cannot talk about effectivity. In fact, the lemma would remain true for uncountable words, and satisfiability of mso sentences for such words is undecidable, which means that for uncountable words the constructions in the lemma cannot be made effective. The issue of finite representation for \circ -semigroups will be addressed in the following section; and countability will play a crucial role.

Another interesting question is about the converse of the lemma: can one define in mso every language that is recognised by a finite \circ -semigroup? For finite words and ω -words, the answer was “obviously yes”, because one can use mso to formalise the acceptance by an automaton. Since we have no automata for \circ -words, the question is harder. However, the answer is still “yes”, and it will be given in Section 3.4.

Exercises

Exercise 72. (1) Give a formula of mso which is true in some uncountable well-founded linear order, but is false in all countable well-founded linear orders.

Exercise 73. (1) Find two countable ordinals (viewed as \circ -words over a one letter alphabet), which have the same mso theory.

Exercise 74. (1) We write ω^* for the reverse of ω . An $(\omega^* + \omega)$ -word is a \circ -word where the underlying order is the same as for the integers. Show that the following problem is decidable: given an mso sentence, decide if it is true in some bi-infinite word.

Exercise 75. (2) We say that a $(\omega^* + \omega)$ -word v is *recurrent* if every finite word $w \in \Sigma^+$ appears as an infix in every prefix of v and in every suffix of v . Show that all recurrent $(\omega^* + \omega)$ -words have the same mso theory.

Exercise 76. (2) Let Σ be an alphabet, and let $x \notin \Sigma$ be a fresh letter. For $w \in \Sigma^\circ$ and $u \in (\Sigma \cup \{x\})^\circ$, define $u[x := w] \in \Sigma^\circ$ to be the result of substituting each occurrence of variable x in u by the argument w . For a language $L \subseteq \Sigma^\circ$, define *contextual equivalence* to be the equivalence relation on Σ° defined by

$$w \sim w' \quad \text{iff} \quad u[x := w] \in L \Leftrightarrow u[x := w'] \in L \text{ for every } u \in (\Sigma \cup \{x\})^\circ.$$

Show that \sim is a \circ -congruence (which means that the function that maps w to its equivalence class is compositional) for every language recognised by some finite \circ -semigroup.

Exercise 77. (1) Give an example of a language $L \subseteq \Sigma^\circ$ where contextual equivalence is not a \circ -congruence.

Exercise 78. (1) Show that every language recognised by a finite \circ -semigroup has syntactic \circ -semigroup, but there are some languages (not recognised by finite \circ -semigroups), which do not have a syntactic \circ -semigroup.

Exercise 79. (1) Consider a binary tree (every node has either zero or two children, and we distinguish left and right children), where leaves are labelled by an alphabet Σ . The tree might have infinite branches. Define the *yield* of such a tree to be the \circ -word where the positions are leaves of the tree, the labels are inherited from the tree, and the ordering on leaves is lexicographic (for every node, its left subtree is before its right subtree). Show that every \circ -word can be obtained as the yield of some tree.

Exercise 80. (1) Show that the following problems are equi-decidable:

- given an mso sentence, decide if it is true in some \circ -word $w \in \Sigma^\circ$
- given an mso sentence, decide if its true in $(\mathbb{Q}, <)$.

Exercise 81. (1) Assume Rabin's Theorem, which says that the mso theory of the complete binary tree

$$(\{0, 1\}^*, \underbrace{x = y0}_{\substack{\text{left} \\ \text{child}}}, \underbrace{x = y1}_{\substack{\text{right} \\ \text{child}}})$$

is decidable. Show that the problems from Exercise 80 are decidable. (We will also prove this in the next section, without assuming Rabin's theorem.)

3.3 Finite representation of \circ -semigroups

The product operation in a finite semigroup can be seen as an operation of type $S^+ \rightarrow S$, or as a binary operation of type $S^2 \rightarrow S$. The binary operation has the advantage that a finite semigroup can be represented in a finite way, by giving a multiplication table. In this section, we show that a similar finite representation is also possible for \circ -semigroups. Apart from binary product, we will use two types of ω -iteration – one forward and one backward – and a shuffle operation (which inputs a set of elements, and not a tuple of fixed length).

Definition 3.12 (Läuchli-Leonard operations). For a \circ -semigroup, define its *Läuchli-Leonard operations*⁷ to be the following four operations (with their types written in red).

$$\begin{array}{cccc}
 \underbrace{ab} & \underbrace{a^\omega} & \underbrace{a^{\omega*}} & \underbrace{\{a_1, \dots, a_n\}^\eta} \\
 \text{binary} & \text{product of} & \text{product of} & \text{product of the} \\
 \text{product} & \text{aaa} \cdots & \cdots \text{aaa} & \text{shuffle of } a_1, \dots, a_n \\
 S^2 \rightarrow S & S \rightarrow S & S \rightarrow S & PS \rightarrow S
 \end{array}$$

Theorem 3.13. *The product operation in a finite \circ -semigroup is uniquely determined by its Läuchli-Leonard operations.*

Another way of stating the theorem is that if S is a finite set S equipped with Läuchli-Leonard operations, then there is at most one way of extending these operations to an associative product $S^\circ \rightarrow S$. We say at most one instead of exactly one, because the Läuchli-Leonard operations need to satisfy certain associativity laws, such as:

$$aa^\omega = a^\omega \quad (ab)^\omega = a(ba)^\omega \quad \{a_1, \dots, a_n\}^\eta = \{\{a_1, \dots, a_n\}^\eta\}^\eta$$

We do not worry too much about giving the full set of axioms⁸, because we will only consider product operations that arise from compositional functions – e.g. the product operation on mso types of given quantifier rank k – and such product operations are guaranteed to be associative.

3.3.1 Proof of Theorem 3.13

The key idea in the proof of Theorem 3.13 is that the Läuchli-Leonard operations are enough to generate all sub-algebras, as stated in the following lemma.

⁷ [24] Läuchli and Leonard, “On the elementary theory of linear order”, 1966 , p. 109

⁸ It can be found in
 [3] Bloom and Ésik, “The equational theory of regular words”, 2005 , Section 7

Lemma 3.14. *Let S be a \circ -semigroup, and let $\Sigma \subseteq S$. Then*

$$\underbrace{\{\text{product of } w : w \in \Sigma^\circ\}} \subseteq S$$

this is called the sub-algebra generated by Σ

is equal to the smallest subset of S which contains Σ and is closed under the L\"auchli-Leonard operations.

Before proving the lemma, we use it to prove Theorem 3.13.

Proof of Theorem 3.13, assuming Lemma 3.14. Suppose that S_1 and S_2 are two \circ -semigroups, which have the same underlying set, and product operations which agree on the L\"auchli-Leonard operations. We will show that the product operations are the same. Consider the product \circ -semigroup $S_1 \times S_2$, defined in the usual coordinate-wise way. Apply Lemma 3.14 to the diagonal

$$\Sigma = \{(a, a) : a \in S\} \subseteq S_1 \times S_2.$$

Since the L\"auchli-Leonard operations agree for S_1 and S_2 , it follows from the lemma that the sub-algebra generated by Σ is also on the diagonal, which shows that the product operations of S_1 and S_2 are equal. \square

The rest of Section 3.3.1 is devoted to proving Lemma 3.14. Define $L \subseteq \Sigma^\circ$ to be the \circ -words whose product can be obtained from Σ by applying the L\"auchli-Leonard operations. The lemma says that $L = \Sigma^\circ$. This follows immediately from the following lemma (which is stated slightly more generally, because it will be used again later), by taking λ to be the product operation of the \circ -semigroup.

Lemma 3.15. *Assume that $L \subseteq \Sigma^\circ$ contains Σ and is closed under binary concatenation. A sufficient condition for $L = \Sigma^\circ$ is that there exists a colouring $\lambda : \Sigma^\circ \rightarrow C$, with C a finite set of colours, such that:*

(*) *Let $w \in (\Sigma^\circ)^\circ$ be such that every position has label in L , and*

$$\underbrace{\lambda^\circ(w) = c^\omega}_{\text{for some } c \in C} \quad \text{or} \quad \underbrace{\lambda^\circ(w) = c^{\omega*}}_{\text{for some } c \in C} \quad \text{or} \quad \underbrace{\lambda^\circ(w) = \text{shuffle of } D.}_{\text{for some } D \subseteq C}$$

Then the flattening of w belongs to L .

Before proving the lemma, we fix some notation for \circ -words. Define an *interval* in a \circ -word to be any set of positions X such that is connected in the following sense:

$$\forall x \in X \forall y \in Y \forall z \quad x < z < y \Rightarrow z \in X.$$

An infix of a \circ -word is any \circ -word obtained by restricting the positions to

some interval. For example, the rationals – viewed as a \circ -word over a one letter alphabet – have uncountably many intervals, but five possible infixes. If $x < y$ are positions in a \circ -word w , then we write $w(x..y)$ for the infix corresponding to the interval $\{z : x < z < y\}$.

Proof Suppose then that L satisfies (*), as witnessed by a colouring λ . We say that $w \in \Sigma^\circ$ is *simple* if all of its infixes are in L . We will show that all of Σ° is simple, thus proving $L = \Sigma^\circ$. For the sake of contradiction, suppose that $w \in \Sigma^\circ$ is not simple. Define \sim to be the binary relation on positions in w , which identifies positions if they are equal, or $w(x..y)$ is simple (where x is the smaller position and y is the bigger position).

Claim 3.16. *The relation \sim is an equivalence relation, every equivalence class is an interval, and this interval induces a simple \circ -word.*

Proof The relation \sim is symmetric and reflexive by definition. Transitivity is because simple words are closed under binary concatenation (which itself easily follows from the fact that L contains all letters and is closed under binary concatenation). This establishes that \sim is an equivalence relation. Because simple \circ -words are closed under infixes by definition, every equivalence class of \sim is an interval.

It remains to show that every (infix induced by an) equivalence class is simple. Here we use countability and the assumption (*). Consider an equivalence class X . Choose some $x \in X$. We will show that the suffix of X that begins in x is simple. A symmetric argument will establish that the prefix leading up to x is simple, and thus X itself is simple by binary concatenation.

If X has a last position, then the suffix that begins in x is simple by definition of \sim . Suppose then that there is no last position in X . By countability, choose some sequence of positions

$$x = x_0 < x_1 < x_2 < \dots \in X$$

which is co-final, i.e. every position in X is smaller than some x_i . By the Ramsey Theorem, we can assume without loss of generality that there is some $c \in C$ such that for every $i \in \{1, 2, \dots\}$, the infix obtained from $w(x_i..x_{i+1})$ by appending position x_{i+1} has colour c under λ . Furthermore, this infix is simple by definition of \sim . Therefore, thanks to assumption (*), we know that the concatenation of all of these infixes is simple. \square

Since the equivalence classes of \sim are intervals, they can be viewed as an ordered set, with the order inherited from the original order on positions in w . Because simple words are closed under binary concatenation, the order on equivalence classes is dense, since otherwise two consecutive equivalence classes

would need to be merged into a single one. Define $w_{\sim} \in C^{\circ}$ to be the result of replacing every equivalence class of \sim by its colour under λ . By assumption that w is not simple, \sim has more than one equivalence class, and therefore the positions of w_{\sim} are an infinite dense linear order.

Claim 3.17. *Some infix of w_{\sim} is a shuffle.*

Proof Take some colour $c \in C$. If there is some infinite infix of w_{\sim} where c does not appear at all, then we can continue working in that infix (its positions are still an infinite dense linear order). Otherwise, positions with colour c are dense. By iterating this argument for all finitely many colours in C , we find an infinite infix where every colour either does not appear at all, or is dense. This infix is a shuffle. \square

By (*), the flattening of the infix from the above claim is simple. It follows that the corresponding interval should have been a single equivalence class of \sim , contradicting the assumption. \square

3.3.2 Decidability of MSO

By Theorem 3.13, a finite \circ -semigroup can be represented in a finite way, by giving its underlying set and the multiplication tables for its Lauchli-Leonard operations. We will use this representation to give decision procedure for mso on \circ -words.

A powerset construction. To decide mso, we will use a powerset construction for finite \circ -semigroups, which will correspond to set quantification in mso. We begin by describing this construction.

Definition 3.18. For a \circ -semigroup S , define the *powerset \circ -semigroup* PS as follows. The underlying set of PS is the powerset of the underlying set of S , including the empty set⁹. The product operation is defined by

$$w \in (PS)^{\circ} \mapsto \underbrace{\{\text{product of } v : v \in^{\circ} w\}}_{\text{in } S},$$

where $v \in^{\circ} w$ means that $v \in S^{\circ}$ can be obtained from $w \in (PS)^{\circ}$ by choosing

⁹ Whether or not we allow the empty set is not important for the construction.

for each position an element of its label¹⁰. We leave it as an exercise to check that the product operation defined this way is associative¹¹.

The following lemma shows that the powerset construction is computable. What is not obvious is finding the multiplication tables for the Lauchli-Leonard operations.

Lemma 3.19. *Given the multiplication tables for the Lauchli-Leonard operations in a finite \circ -semigroup S , one can compute the multiplication tables for the Lauchli-Leonard operations in the powerset \circ -semigroup \mathbf{PS} .*

Proof In the proof, we adopt the convention that elements of S are denoted by lower-case letters a, b, c , while elements of \mathbf{PS} are denoted by upper-case letters A, B, C . The multiplication table for binary product of \mathbf{PS} , namely

$$A, B \in \mathbf{PS} \quad \mapsto \quad \underbrace{\{ab\}}_{\text{product in } S} : a \in A, b \in B,$$

is easily computable using the binary product in S . The hard part is the multiplication tables for the infinitary operations, i.e. ω -power, ω^* -power and shuffles.

ω -power. We begin by clarifying some notation. For a $A \in \mathbf{PS}$, the expression A^ω can be understood in three different ways:

- (1) $A^\omega \subseteq S^\circ$ is the set of ω -words where all letters are from A ;
- (2) $A^\omega \in (\mathbf{PS})^\circ$ is the ω -word where all letters have label equal to A ;
- (3) $A^\omega \in \mathbf{PS}$ is the product of the word from item (2) in the \circ -semigroup \mathbf{PS} .

To avoid confusion, we use the red type annotation below. The Lauchli-Leonard operation in the powerset \circ -semigroup \mathbf{PS} that we are discussing in this lemma uses the third meaning of A^ω :

$$A \in \mathbf{PT} \quad \mapsto \quad A^\omega \in \mathbf{PS}.$$

To compute this operation, we will use, apart from S , one other \circ -semigroup. This other \circ -semigroup, call it T , is used to recognise the singleton language

$$\{A^\omega \in (\mathbf{PS})^\circ\} \subseteq (\mathbf{PS})^\circ.$$

¹⁰ The relation $v \in^\circ w$ can be formalised by saying that there exists a \circ -word u over alphabet

$$\{(a, A) : a \in A \subseteq S\}$$

such that v is the projection of u to the first coordinate, and w is the projection of u to the second coordinate.

¹¹ One has to a bit careful. For example, there is no such thing as a powerset group.

The elements of T are $\{\omega, +, 0\}$ and the product operation is the unique product which makes the following function h into a homomorphism:

$$w \in (\text{PS})^\circ \mapsto \begin{cases} \omega & \text{if } w = A^\omega \in (\text{PS})^\circ \\ + & \text{if } w \text{ is a finite word and all letters have label } A \\ 0 & \text{otherwise} \end{cases}$$

For the Lauchli-Leonard operations in T , all outputs are 0 with the following exceptions:

$$++ = + \quad +\omega = \omega \quad +^\omega = \omega.$$

Claim 3.20. *Let $a \in S$. Then $a \in A^\omega \in \text{PS}$ if and only if (a, ω) belongs to the sub-algebra of the \circ -semigroup $S \times T$ that is generated by $\{(b, +) : b \in A\}$.*

Proof By unravelling the definitions, (a, ω) belongs to the sub-algebra from the claim if and only if it is the product of some \circ -word u where every letter is of the form $(b, +)$ for some $b \in A$. Since the product of u has ω on the second coordinate, then u must be an ω -word. Therefore, if we project u onto the first coordinate, we get a word in $A^\omega \subseteq S^\circ$ whose product is a , thus proving the right-to-left implication. The left-to-right implication reverses this reasoning. \square

By the above claim, in order to compute $A^\omega \in \text{PS}$, it is enough to compute the sub-algebra from the claim. Thanks to Lemma 3.14, this sub-algebra is the closure of $\{(b, +) : b \in A\}$ under the Lauchli-Leonard operations of $S \times T$, which are simply the coordinate-wise liftings of the Lauchli-Leonard operations in S and T . Therefore, $A^\omega \in \text{PS}$ can be computed.

Shuffle power. The argument for ω^* -power is symmetric to the one above, so we are left with the shuffle power in the \circ -semigroup PS :

$$\{A_1, \dots, A_n\} \mapsto \{A_1, \dots, A_n\}^\sharp$$

We use a similar argument as for the ω -power, except that we need to define a \circ -semigroup which describes the singleton language $\{w\}$ where

$$w \stackrel{\text{def}}{=} \underbrace{\text{shuffle of } \{A_1, \dots, A_n\}}_{\text{a } \circ\text{-word over alphabet PS}}.$$

Such a \circ -semigroup T and a recognising homomorphism

$$h : (\text{PS})^\circ \rightarrow T$$

are left as an exercise for the reader (see Exercise 82). By definition of the powerset \circ -semigroup,

$$\{A_1, \dots, A_n\}^\eta = \underbrace{\{\text{product of } v : v \in^\circ w\}}_{\text{in } S}.$$

Let $a \in S$. The same proof as for Claim 3.20 shows that

$$a \in \{A_1, \dots, A_n\}^\eta$$

holds if and only if $(a, h(w))$ belongs to the sub-algebra of $S \times T$ that is generated by elements of the form

$$\{(b, h(A_i)) : i \in \{1, \dots, n\}, b \in A_i\}.$$

This sub-algebra can be computed, as in the case of ω -power. □

Decidability of MSO. Using the powerset construction on \circ -semigroups, we can now prove decidability of MSO over \circ -words.

Theorem 3.21. *The following problem is decidable:*

Input. An MSO sentence φ , which defines a language $L \subseteq \Sigma^\circ$.

Question. Is the language L nonempty?

The rest of Section 3.3.2 is devoted to proving the above theorem. As in Section 2.1, instead of using MSO in the ordered model, it will be easier to use first-order logic over (the extension for \circ -words of) the set model, see Definition ???. By induction on formula size, for every first-order formula over the set model we will construct a homomorphism into a finite \circ -semigroup that recognises the language of the formula. The finite \circ -semigroup will be represented, according to Theorem 3.13, by giving the underlying set and the multiplication tables for its Lauchli-Leonard operations.

For the induction, we need to deal with formulas with free variables. Consider a first-order formula

$$\varphi(\underbrace{x_1, \dots, x_n}_{\substack{\text{the free variables range} \\ \text{over sets of positions}}})$$

over the vocabulary of the set model (over some fixed input alphabet Σ .) We define the *language of φ* to be the set of \circ -words over alphabet $\Sigma \times \{0, 1\}^n$, such that φ is true in the projection to the Σ coordinate, assuming that variable x_i is set to the set of positions that have 1 on the i -th coordinate from $\{0, 1\}$.

Lemma 3.22. *Let Σ be an input alphabet. Given a formula $\varphi(x_1, \dots, x_n)$ of first-order logic over the vocabulary of the set model, we can compute a finite \circ -semigroup S , a (not necessarily surjective) homomorphism*

$$h : (\Sigma \times \{0, 1\}^n)^\circ \rightarrow S,$$

and an accepting set $F \subseteq S$ such that the language of φ is exactly $h^{-1}(F)$. The \circ -semigroup is represented by multiplication tables for its Lauchli-Leonard operations, and the homomorphism is represented by its restriction to one-letter words.

Once we have proved the lemma, Theorem 3.21 follows immediately. We simply need to check if the image of h contains some accepting element. There is a slightly subtle point: since h is not necessarily surjective, the accepting set could be nonempty but disjoint with the image of h . Therefore, we need to compute the image of h , which is done by closing the images of the one-letter words under the Lauchli-Leonard operations.

It remains to prove the lemma.

Proof of Lemma 3.22 Induction on the size of $\varphi(x_1, \dots, x_n)$.

Atomic formulas. The atomic formulas are $x \subseteq y$, $x < y$, $x = \emptyset$ and $x \subseteq a$.

With the exception of $x < y$, all of the atomic formulas are of the form “the label of every position has some property”. Therefore, for every atomic formula except for $x < y$, the corresponding language is recognised by a homomorphism into the semigroup $\{0, 1\}$ with product defined by

$$w \in \{0, 1\}^\circ \mapsto \begin{cases} 0 & \text{if some position has label 0} \\ 1 & \text{if all positions have label 1.} \end{cases}$$

For $x < y$, the appropriate \circ -semigroup is the one from Example 12.

Boolean combinations. For negation $\neg\varphi$, we use the same homomorphism as for φ , and we complement the accepting set. For conjunction $\varphi_1 \wedge \varphi_2$ and disjunction $\varphi_1 \vee \varphi_2$, we use the product $S_1 \times S_2$ of the inductively obtained \circ -semigroups, with a naturally defined homomorphism¹². Since $S_1 \times S_2$ is defined coordinate-wise, the multiplication tables for its Lauchli-Leonard operations can be computed using those from S_1 and S_2 .

Quantification. We are left with quantification. Consider an existentially quantified formula (for universal quantification, the reasoning is the same):

$$\exists x_{n+1} \varphi(x_1, \dots, x_{n+1}).$$

¹² The homomorphism needs to account for the possibility that φ_1 and φ_2 use different subsets of the free variables in $\varphi_1 \wedge \varphi_2$.

Apply the induction assumption, yielding a homomorphism

$$h : (\Sigma \times \{0, 1\}^{n+1})^\circ \rightarrow S$$

which recognises the language of φ . Consider the function

$$H : (\Sigma \times \{0, 1\}^n)^\circ \rightarrow \underbrace{PS}_{\text{powerset } \circ\text{-semigroup}},$$

such that $H(w)$ is the set of all values $h(v) \in S$, where v ranges over \circ -words over alphabet $\Sigma \times \{0, 1\}^{n+1}$ such that w can be obtained from v by erasing the last bit from each letter. It is not hard to see that H is a homomorphism. The homomorphism H recognises the language of the existentially quantified formula; the accepting set consists of those subsets of S which intersect the accepting set for φ . The multiplication tables for the L\"auchli-Leonard operations in PS can be computed thanks to Lemma 3.19.

□

Exercises

Exercise 82. (1) Let Σ be a finite alphabet, and let w be the shuffle of all letters in Σ . Show a finite \circ -semigroup which recognises the singleton language $\{w\}$.

Exercise 83. (2) A \circ -word w is called *regular* if the singleton language $\{w\}$ is recognised by a finite \circ -semigroup. Show that w is regular if and only if it can be constructed from the letters by using the L\"auchli-Leonard operations.

Exercise 84. (1) Show that every nonempty mso definable language $L \subseteq \Sigma^\circ$ contains some regular \circ -word.

Exercise 85. (1) Show that if w is a regular \circ -word, then $\{w\}$ is mso definable (without invoking Theorem 3.23).

Exercise 86. (2) Show that for every finite alphabet Σ there exists a \circ -word $w \in \Sigma^\circ$ such that

$$h(wvw) = h(w) \quad \text{for every } \underbrace{h : \Sigma^\circ \rightarrow S}_{\substack{\text{homomorphism into} \\ \text{a finite } \circ\text{-semigroup}}} \text{ and } v \in \Sigma^\circ.$$

Exercise 87. (2) For a countable linear order X , let $\{a, b\}^X \subseteq \{a, b\}^\circ$ be the set of \circ -words with positions X . We can equip this set with a probabilistic measure, where for each position $x \in X$, the label is selected independently, with a and b both having probability half. We say that X has a zero-one law if for every mso definable language L , the probability of $\varphi \cap \{a, b\}^X$ is either zero or one. For which of the following $X = \mathbb{N}, \mathbb{Z}, \mathbb{Q}$ is there a zero-one law?

Exercise 88. (2) A countable linear order can be viewed as a \circ -word over a one-letter alphabet. Among these, we can distinguish the countable linear orders that are regular, i.e. generated by the L\"auchli-Leonard operations, see Exercise 83. Give an algorithm, which inputs a countable linear order that is regular in the above sense, and decides if it has a zero-one law (in the sense of Exercise 87).

Exercise 89. (2) Show that every mso definable language of \circ -words belongs to the least class of languages which:

- contains the following two languages over alphabet $\{a, b, c\}$:

$$\underbrace{\exists x a(x)}_{\text{some } a} \quad \underbrace{\exists x \exists y a(x) \wedge b(y) \wedge x < y}_{\text{a before b}}$$

- is closed under Boolean combinations;
- is closed under images and inverse images of letter-to-letter homomorphisms.

Exercise 90. (2) We say that a binary tree (possibly infinite) is *regular* if it has finitely many non-isomorphic sub-trees. Show that a \circ -word is regular (in the sense of Exercise 83) if and only if it is the yield (in the sense of Exercise 79) of some regular tree.

Exercise 91. (91) Consider the embedding ordering (Higman ordering) $w \hookrightarrow v$ on \circ -words. Show that for every \circ -word w there is a regular \circ -word v such that $w \hookrightarrow v$ and $v \hookrightarrow w$. Hint: use Lemma 3.15.

Exercise 92. (1) Suppose that we are given a language $L \subseteq \Sigma^\circ$, represented by a finite \circ -semigroup S , a homomorphism $h : \Sigma^\circ \rightarrow S$, and an accepting set $F \subseteq S$. Give an algorithm which computes the syntactic \circ -semigroup (which exists by Exercise 78).

Exercise 93. (2) Let \mathcal{L} be a class of languages, such that \mathcal{L} satisfies the following conditions:

- every language in \mathcal{L} is recognised by a finite \circ -semigroup;
- \mathcal{L} is closed under Boolean combinations;
- \mathcal{L} is closed under inverse images of homomorphisms $h : \Sigma^\circ \rightarrow \Gamma^\circ$;
- Let $L \subseteq \Sigma^\circ$ be a language in \mathcal{L} . For every $w, w_1, \dots, w_n \in \Sigma^\circ$, \mathcal{L} contains the inverse image of L under the following operations:

$$v \mapsto wv \quad v \mapsto vw \quad v \mapsto v^\omega \quad v \mapsto v^{\omega*} \quad v \mapsto \text{shuffle of } \{w_1, \dots, w_n, v\}.$$

Show that if L belongs to \mathcal{L} , then the same is true for every language recognised by its syntactic \circ -semigroup.

Exercise 94. (2) Let Σ be an alphabet and let $c \notin \Sigma$ be a fresh letter. We say that $L \subseteq \Sigma^\circ$ is definable in $\text{LTL}[F]$ if there is a formula of $\text{LTL}[F]$ which defines the language cL , see Exercise 41. Give an algorithm which inputs the finite syntactic \circ -semigroup of a language $L \subseteq \Sigma^\circ$, and answers if the language is definable in $\text{LTL}[F]$. Hint: the \circ -semigroup must be suffix trivial, but this is not sufficient.

Exercise 95. (3) Give an algorithm which inputs the finite syntactic \circ -semigroup of a language $L \subseteq \Sigma^\circ$, and answers if the language is definable in two-variable first-order logic FO^2 . Hint: the \circ -semigroup must be in DA , but this is not sufficient.

Exercise 96. (2) Show that aperiodicity is not sufficient for first-order definability for \circ -words: give an example of a language $L \subseteq \Sigma^\circ$ that is recognised by a finite aperiodic \circ -semigroup, but which is not definable in first-order logic.

3.4 From \circ -semigroups to MSO

In Theorem 3.11, and again in Lemma 3.22, we have shown that if a language of \circ -words is definable in MSO , then it is recognised by a finite \circ -semigroup. We now show that the converse implication is also true.

Theorem 3.23. *If a language of \circ -words is recognised by a finite \circ -semigroup, then it is definable in MSO ¹³ which says that*

¹³ This theorem was first shown in
 [10] Carton, Colcombet, and Puppis, “An algebraic approach to MSO-definability on countable linear orders”, 2018, Theorem 5.1.
 The proof presented here is different, and it is based on the proof in
 [33] Schützenberger, “On finite monoids having only trivial subgroups”, 1965, p. 192
 which shows that every aperiodic monoid recognises a star-free language. We use the different

As mentioned before in this chapter, the theorem would be easy if there was an automaton model, which would assign states to positions, and where the acceptance condition could be formalised in mso. Unfortunately, no such automaton model is known. Therefore, we need a different proof for the theorem. The rest of Section 3.4 is devoted to such a proof.

We begin by defining regular expressions for \circ -words. For a finite family \mathcal{L} of languages of \circ -words, define the shuffle of \mathcal{L} to be the \circ -words which can be partitioned into intervals so that: (a) every interval induces a word from L for some $L \in \mathcal{L}$; (b) the order type on the intervals is that of the rational numbers; and (c) for every $L \in \mathcal{L}$, the intervals from L are dense.

Lemma 3.24. *Languages definable in mso are closed under Boolean combinations and the following kinds of concatenation:*

$$LK \quad L^+ \quad L^\omega \quad L^{\omega*} \quad \text{shuffle of } \underbrace{\mathcal{L}}_{\substack{\text{a finite family} \\ \text{of languages}}}$$

Proof For the Boolean operations, there is nothing to do, since Boolean operations are part of the logical syntax. For the concatenations, we observe that mso can quantify over factorisations, as described below.

Define a *factorisation* of a \circ -word to be a partition of its positions into intervals, which are called *blocks*. For a factorisation, define a *compatible colouring* to be any colouring of positions that uses two colours, such that all blocks are monochromatic, and if a block has a successor, then the successor has a different colour. A compatible colouring always exists (there could be uncountably many choices). A factorisation can be recovered from any compatible colouring: two positions are in the same block if and only if the interval connecting them is monochromatic. A compatible colouring can be represented using a single set – namely the positions with one of the two colours. This representation can be formalised in mso, i.e. one can write an mso formula $\varphi(x, y, X)$ which says that positions x and y are in the same block of the factorisation which arises from the compatible colouring represented by set X .

Using the above representation, we show closure of mso under the concatenations in the lemma. For LK , we simply say that there exists a factorisation with two blocks, where the first block is in L and the second block is in K . (To say that a block is in L or K , we observe that mso sentences can be relativised to a given interval.) For L^+ , we say that there exists a factorisation with finitely many blocks, where all blocks are in L . Here is how we express that there are finitely many blocks: there are first and last blocks, and there is

proof because, after suitable modifications, it allows us to characterise star-free languages of \circ -words, see Exercise 102.

no proper subset of positions that contains the first block and is closed under adding successor blocks. For L^ω , we do the same, except that there is no last block. For $L^{\omega*}$, we use a symmetric approach. For the shuffle, we say that the blocks are dense and there is no first and last block. \square

In the proof of Theorem 3.23, we will only use the closure properties of mso from the above lemma. In particular, it will follow that every language recognised by a finite \circ -semigroup can be defined by a regular expression which uses single letters and the closure operations from the lemma (which include intersection and complementation).

To prove Theorem 3.23, we will show that the product operation of every finite \circ -semigroup can be defined in mso, in the following sense. Let S be a finite \circ -semigroup. We will show that for every $a \in S$, the language

$$L_a = \{w \in S^\circ : w \text{ has product } a\}$$

is mso definable. This will immediately imply that every language recognised by a homomorphism into S is mso definable, thus proving the theorem.

The proof is by induction on the infix ordering. Fix for the rest of this section an infix class $J \subseteq S$. We partition S into two parts:

$$\underbrace{\text{easy elements}}_{\text{proper prefixes of } J} \cup \underbrace{\text{hard elements}}_{\text{the rest}}.$$

The induction hypothesis says L_a is mso definable for every easy $a \in S$. We need to prove the same thing for every $a \in J$.

We begin with an observation about smooth products, which follows from the Ramsey argument that was used in Theorem 3.2. We say that $w \in S^\circ$ is J -smooth if every finite infix of w has a product in J . This is a lifting to infinite words of the notion of smoothness that was used in Section 1.3. In particular, by Claim 1.22 from that section, a \circ -word is J -smooth if and only if all of its infixes of length at most two are J -smooth. The following lemma describes the products of certain J -smooth words.

Lemma 3.25. *For every idempotent $e \in J$ the following holds. Let $w \in J^\circ$ be J -smooth. If w is an ω -word, then its product is ae^ω , where a depends only on e and the first letter in w . If w is an $(\omega* + \omega)$ -word, then its product is $e^{\omega*}e^\omega$.*

Since the lemma is true for every choice of idempotent e , it follows that $e^{\omega*}e^\omega$ does not depend on the choice of e .

Proof The main observation is the following claim.

Claim 3.26. *If w is J -smooth and has first letter e , then its product is e^ω .*

Proof By Lemma 3.4, the product of w is equal to af^ω , for some a, f . Since w is J -smooth, a and f belong to J . Since the first letter of w is e , we have $ea = a$. Since f is infix equivalent to e , it admits a decomposition $f = xee y$. Therefore

$$af^\omega = ea(xee y)^\omega = \underbrace{eaxe}_g (\underbrace{eyxe}_h)^\omega.$$

We now continue as in the proof of Lemma 3.5: because g, h, e are in the same group, then $g^\omega = e^\omega = h^\omega$, and therefore $gh^\omega = e^\omega$. \square

The claim immediately proves the lemma. Indeed, consider a J -smooth ω -word with first letter a . The first letter admits a decomposition axe for some x , because every prefix class intersects the suffix class of e . By the above claim, the product of every J -smooth ω -word that begins with a is equal to axe^ω . A similar argument works when the positions are ordered as the integers: every J -smooth $(\omega^* + \omega)$ -word has the same product as a smooth $(\omega^* + \omega)$ -word with an infix ee , and the latter has product $e^{\omega^*} e^\omega$ thanks to the claim and its symmetric version for ω^* . \square

In the rest of the proof, we will use the following terminology. We say that a colouring (a function from \circ -words to a finite set of colours) is mso definable if for every colour, its inverse image is an mso definable language. We say that a colouring λ is mso definable on a subset L of inputs if there exists an mso definable colouring that agrees with λ on inputs from L .

The strategy for the rest of the proof is as follows. Define $L_J \subseteq S^\circ$ to be the \circ -words that have product in J . We first show in Lemma 3.27 that the colouring

$$w \in S^\circ \quad \mapsto \quad \text{prefix class of the product of } w$$

is mso definable on L_J . Next, we use this result about prefixes and a symmetric one for suffixes, to show in Lemma 3.29 that the product operation is mso definable on L_J . Finally, in Lemma 3.32 we show that the language L_J is definable in mso. We can then conclude as follows: a \circ -word has product $a \in J$ if and only if it belongs to L_J , and the colouring from Lemma 3.29 step maps it to a . It remains to prove the lemmas.

Lemma 3.27. *The following colouring is mso definable on L_J :*

$$w \in S^\circ \quad \mapsto \quad \text{prefix class of the product of } w.$$

Proof We write $H \subseteq S^\circ$ for the \circ -words with a hard product. This language is definable in mso, as the complement of the easy products that are definable by induction assumption. For an interval in w , define its product to be the product

of the infix of w that is induced by the interval. An interval is called *easy* if its product is easy, otherwise it is called *hard*. By the induction assumption, we can check in MSO if an interval is easy or hard. An interval is called *almost easy* if all of its proper sub-intervals are easy.

Claim 3.28. *The product operation is MSO definable on easy intervals.*

Proof If there is a last position, then the product can be easily computed: remove the last position, compute the product, and then add the last position. Otherwise, if there is no last position, then we can use Lemma 3.4 to see that an almost easy interval has product $b \in S$ if and only if it belongs to

$$L_a(L_e)^\omega \quad \text{for some easy } a, e \text{ such that } ae^\omega = a.$$

The above condition can be formalised in MSO thanks to the induction assumption and Lemma 3.24. \square

Define a *prefix interval* to be a nonempty downward closed interval. We will compute in MSO the prefix class of some hard prefix interval; if the \circ -word is in L_J then this hard prefix has the same prefix class as w . We do a case disjunction, depending on whether or not there is an easy prefix interval (which can clearly be checked in MSO).

- Suppose first that there is no easy prefix interval. This means that either w has a first letter which is in J , or $w \in H^{\omega^*}$. In the first case, the first letter uniquely determines the prefix class of the product of w . In the second case, when $w \in H^{\omega^*}$, then under the assumption of $w \in L_J$, we can use Lemma 3.25 to conclude that the product of w is in the prefix class of e^{ω^*} , where e is some arbitrarily chosen idempotent from J .
- Suppose next that there is some easy prefix interval. Let X be the union of all easy prefix intervals. This is an almost easy interval, and therefore its product $a \in S$ can be computed thanks to Claim 3.28. If a is hard, then we know the prefix class of w . Otherwise, if a is easy, it follows that after removing X , we get a \circ -word as in the first case, and we can use that case to compute the prefix class.

\square

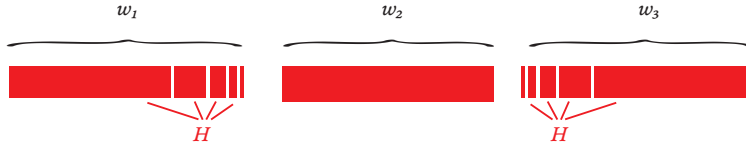
Lemma 3.29. *The product operation of S is MSO definable on L_J .*

Proof Let $w \in L_J$. We use the terminology about intervals from the proof of Lemma 3.27.

Claim 3.30. *There exists a factorisation $w = w_1w_2w_3$ such that:*

- w_1 is either empty or in H^ω ;
- w_2 is a finite concatenation of almost easy \circ -words;
- w_3 is either empty or in $H^{\omega*}$.

Here is a picture of the factorisation:



Proof Define a *limit prefix* of w to be any prefix interval which induces a \circ -word in H^ω . Limit prefixes are closed under (possibly infinite) unions. If there is a limit prefix, then there is a maximal one, namely the union of all limit prefixes (if there are not limit prefixes, we define the maximal limit prefix to be empty). Define $w_1 \in H^\omega$ to be the maximal limit prefix of w (if no limit prefix exists, then w_1 is empty). Remove the prefix w_1 , and to the remaining part of the word apply a symmetric process, yielding a suffix $w_3 \in H^{\omega*}$ and a remaining part w_2 . This is the factorisation in the statement of the claim.

It remains to show that w_2 is a finite concatenation of almost easy \circ -words. By construction, the remaining part w_2 does not have any prefix in H^ω , nor does it have any suffix in $H^{\omega*}$. Take the union of all easy prefixes of w_2 (this union exists, because w_2 has no suffix in $H^{\omega*}$, and it is almost easy), and cut it off. After repeating this process a finite number of times, we must exhaust all of w_2 , since otherwise there would be a prefix in H^ω . Therefore, w_2 is a finite concatenation of almost easy \circ -words. \square

Let w_1, w_2, w_3 be as in the above claim. By Lemma 3.25, the product of w_1 is uniquely determined by its prefix class (under the assumption that the entire \circ -word belongs to L_J). Therefore, thanks to Lemma 3.27, we can compute in mso the product of the w_1 . Symmetrically, we can compute the product of w_3 . It remains to compute the product of w_2 . This is done in the following claim.

Claim 3.31. *If a \circ -word is a finite concatenation almost easy \circ -words, then its product can be computed in mso.*

Proof By the Kleene theorem about regular expressions being equivalent to finite automata, the set of finite concatenations of almost easy intervals can be described using a regular expressions, where the atomic expressions describe almost easy words of given product. Such a regular expression can be formalised in mso thanks to Lemma 3.24 \square

□

Lemma 3.32. *The language L_J is mso definable.*

Proof Define $I \subseteq S$ to be the hard elements which are not in J . This is an ideal in the \circ -semigroup S , i.e. if $w \in S^\circ$ has at least one letter in I , then its product is in I . Define L_I to be the \circ -words with product in I . Again, this is an ideal, this time in the free \circ -semigroup S° . We will show how to define L_I in mso; it will follow that L_J is mso definable as

$$L_J = H - L_I.$$

The key is the following characterisation of L_I . Define an *error* to be a \circ -word in S° which satisfies at least one of the following conditions:

- binary error: belongs to $L_a L_b$ for some $a, b \in S - I$ such that $ab \in I$;
- ω -error: belongs to $(L_a)^\omega$, for some $a \in S - I$ such that $a^\omega \in I$;
- ω^* -error: belongs to $(L_a)^{\omega^*}$, for some $a \in S - I$ such that $a^{\omega^*} \in I$;
- shuffle error: is in the shuffle of $\{L_a\}_{a \in A}$ for some $A \subseteq S - I$ such that $A^\eta \in I$.

Note that in the above definition, we can use languages L_a for $a \in J$. These languages are not yet known to be definable in mso.

Claim 3.33. *A \circ -word belongs to L_I if and only if it has an error infix.*

Proof Clearly every error is in L_I , and since L_I is an ideal, it follows that L_I contains every \circ -word with an error infix. We are left with the converse implication: every \circ -word in L_I contains an error infix. To prove this implication, we will show that the language

$$L = \{w \in S^\circ : \text{if } w \in L_I \text{ then } w \text{ has an error infix}\}$$

satisfies the assumptions of Lemma 3.15, with λ being the product operation in S . The conclusion of Lemma 3.15 will then say that L is equal to S° , thus showing that every \circ -word in L_I has an error infix.

The first assumption of Lemma 3.15 says that L is closed under binary concatenation. Suppose that $u, v \in L$. We need to show that $uv \in L$. Suppose that $uv \in L_I$. If $u \in L_I$, then it has an error infix by assumption on $u \in L$, and therefore also uv has an error infix. We argue similarly if $v \in L_I$. Finally, if both u, v have products in $S - I$, then uv is a binary error.

The remaining assumptions of Lemma 3.15 are checked the same way. □

As remarked before Claim 3.33, the definition of errors refers to languages L_a with $a \in J$, which are not yet known to be definable in mso. We deal with this issue now. By Lemma 3.29, for every $a \in J$ there an mso definable language

which contains all \circ -words that have product a , and does not contain any \circ -words that have product in $J - \{a\}$. By removing the \circ -words with easy products from that language, we get an mso definable language K_a with

$$L_a \subseteq K_a \subseteq L_a \cup L_I.$$

Define a *weak error* in the same way as an error, except that K_a is used instead of L_a for $a \in J$. Since K_a is obtained from L_a by adding some words from the ideal L_I , it follows from Claim 3.33 that a \circ -word is in L_I if and only if it has an infix that is a weak error. Finally, weak errors can be defined by an expression which uses mso definable languages and the closure operators from Lemma 3.24, and therefore weak errors are mso definable. It follows that L_I is mso definable, and therefore L_J is mso definable. \square

As we have already remarked when describing the proof strategy, the above lemma completes the proof of the induction step in Theorem 3.23. Indeed, a \circ -word has product $a \in J$ if and only if it belongs to L_J and it is assigned a by the colouring from Lemma 3.29.

Exercises

Exercise 97. (2) The syntax of star-free expression for \circ -words is the same as for finite words, except that the complementation operation is interpreted as $\Sigma^\circ - L$ instead of $\Sigma^* - L$. Define a \circ -star-free language to be a language $L \subseteq \Sigma^\circ$ that is defined by a star-free expression. Show that if L is \circ -star-free, then its syntactic \circ -semigroup is aperiodic, but the converse implication fails.

Exercise 98. (1) What is the modification for \circ -star-free expressions that is needed to get first-order logic (over the ordered model)?

Exercise 99. (2) Show that if $L \subseteq \Sigma^\circ$ is \circ -star-free, then the same is true for every language recognised by its syntactic \circ -semigroup.

Exercise 100. (2) Show that if $L \subseteq \Sigma^\circ$ is \circ -star-free, then the same is true for L^ω .

Exercise 101. (2) Show that if S is aperiodic, then the constructions from Lemmas 3.27 and 3.29 can be done using \circ -star-free expressions.

Exercise 102. (2) Show that $L \subseteq \Sigma^\circ$ is \circ -star-free if and only if its syntactic \circ -semigroup is finite, aperiodic and satisfies¹⁴:

$$e^{\omega^*} = e = e^\omega \quad \Rightarrow \quad e = \{e\}^\eta \quad \text{for every idempotent } e.$$

Hint: use Exercises 100 and 101.

Exercise 103. (1) Show that languages of \circ -words definable in first-order logic (in the ordered model) are not closed under concatenation LK .

Exercise 104. (2) We say that a product operation $\pi : S^\circ \rightarrow S$ is regular-associative if it satisfies the associativity condition from Definition 3.10, but with the diagrams restricted so that only

$$S^\bullet = \{w \in S^\circ : w \text{ is regular}\}$$

is used instead of S° . Show that if S finite and $\pi : S^\circ \rightarrow S$ is mso definable and regular-associative, then π is associative.

Exercise 105. (2) Show that if S is finite and $\pi : S^\bullet \rightarrow S$ is regular associative, then it can be extended uniquely to an associative product $\bar{\pi} : S^\circ \rightarrow S$. Hint: the mso formulas defined in the proof of Theorem 3.23 depend only on the Lauchli-Leonard operations of the \circ -semigroup S .

¹⁴ This exercise is based on
 [12] Colcombet and Sreejith, "Limited Set Quantifiers over Countable Linear Orderings",
 2015, Theorem 2, item 2.

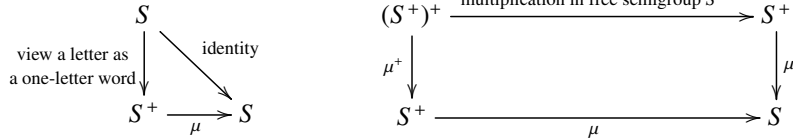
PART TWO

MONADS

4

Monads

As discussed in Chapter 1, instead of viewing a semigroup as having a binary multiplication operation, one could think of a semigroup as a set S equipped with a multiplication operation $\mu : S^+ \rightarrow S$, which is associative in the sense that the following two diagrams commute:



The same is true for monoids, with $*$ used instead of $+$, and for \circ -semigroups, with \circ used instead of $+$. In this chapter, we examine the common pattern behind this constructions, which is that they are the Eilenberg-Moore algebras for the monads of $+$ -words, $*$ -words and \circ -words, respectively.

From the perspective of this book, the idea behind monads is the following. Instead of first defining not necessarily free algebras (e.g. semigroups) and then defining free algebras (e.g. the free semigroup) as a special case, an opposite approach is used. We begin with the free algebra (which is the monad), and then other, not necessarily free, algebras are defined as a derived notion (which is the Eilenberg-Moore algebras of the monad). This opposite approach is useful for less standard algebras such as graphs: axiomatising the not necessarily free algebras is possible but tedious and not intuitive, while the free algebra is very natural, because it consists of graphs with a certain substitution structure.

4.1 Monads and their Eilenberg-Moore algebras

This section, presents the basic definitions for monads and their algebras. These notions make sense for arbitrary categories. However, for simplicity we use the category of sets and functions, because this is where most of our examples live. In later chapters we will consider multi-sorted sets (e.g. sets with sorts $\{+, \omega\}$ for ω -semigroups, or sets with sorts $\{0, 1, \dots\}$ for hypergraphs), but this is as far as we go with respect to the choice of categories.

Definition 4.1 (Monad). A *monad* in the category of sets¹ consists of the following ingredients:

- *Structures*: for every set X , a set TX ;
- *Substitution*: for every function $f : X \rightarrow Y$, a function $Tf : TX \rightarrow TY$;
- *Unit and free multiplication*: for every set X , two functions

$$\underbrace{\text{unit}_X : X \rightarrow TX}_{\text{the unit of } X} \quad \underbrace{\text{mult}_X : TTX \rightarrow TX}_{\text{free multiplication on } X}.$$

These ingredients are subject to six axioms

$$\begin{array}{ccc} \underbrace{(4.1)} & \underbrace{(4.2) \quad (4.3)} & \underbrace{(4.4) \quad (4.5) \quad (4.6)} \\ T \text{ is a functor} & \text{unit and multiplication} & TX \text{ with free multiplication} \\ & \text{are natural transformations} & \text{is an Eilenberg-Moore algebra,} \\ & & \text{and one more associativity axiom} \end{array}$$

which will be described later in this section.

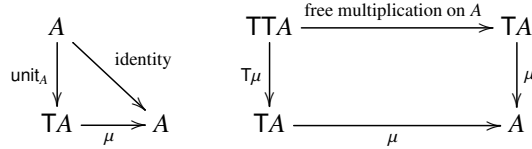
Before describing the monad axioms, we discuss some examples, and define Eilenberg-Moore algebras. The purpose of the monad axioms is to ensure that Eilenberg-Moore algebras are well-behaved. Therefore it is easier to see the monad axioms after the definition of Eilenberg-Moore algebras.

Example 4.2 (Monad of finite words). The monad of finite words is defined as follows. The structures are defined by $TX = X^*$. For a function f , its corresponding substitution Tf is defined by applying f to every letter in the input word. The unit operation maps a letter to the one-letter word consisting of this letter. Free multiplication flattens a word of words into a word. The monad of \circ -words is defined as in the same way, except that one uses \circ -words.

For this book, the key notion about monads is their Eilenberg-Moore algebras. The idea is that TX describes the free algebra, while the Eilenberg-Moore algebras are the algebras that are not necessarily free.

¹ The same definition can be applied to any other category, by using “object” instead of “set”, and “morphism” instead of “function”.

Definition 4.3 (Eilenberg-Moore algebras). An Eilenberg-Moore algebra in a monad T , also called a T -algebra, consists of an *underlying set* A and a *multiplication operation* $\mu : TA \rightarrow A$, subject to the following associativity axioms:



The reader will recognise, of course, the diagrammatic definitions of semi-groups, monoids and ω -semigroups. Also, the diagrammatic definition of ω -semigroup will fall under the scope of the above definition, if we think about ω -semigroups as living in the category of sets with two sorts $\{+, \omega\}$.

By abuse of notation, we use the same letter to denote a T -algebra and its underlying set, assuming that the multiplication operation is clear from the context. Also, if the monad T is clear from the context, we will say algebra instead of T -algebra.

Example 14. [Group monad] The *free group over a set X* is defined to be

$$\underbrace{(X + X)^*}_{\substack{\text{two copies of } X, \\ \text{one blue, and one red}}}$$

modulo the identities

$$xx = xx = \varepsilon \quad \text{for every } x \in X,$$

where ε represents the empty word, while x and x represent the blue and red copies of x . The identities can be applied in any context. For example,

$$zx \quad zyx \quad zzzxyy,$$

represent the same element of the free group. Define T to be the monad where TX is the free group over X , and the remaining monad structure is defined similarly as for finite words, except that we have the two copies of the alphabet, and the identities. The unit operation maps an element to its blue copy.

An algebra over this monad is the same thing as a group. Indeed, if G is an algebra over this monad, with multiplication μ , then the group structure is recovered as follows:

$$\underbrace{1 \stackrel{\text{def}}{=} \mu(\varepsilon)}_{\text{group identity}} \quad \underbrace{x^{-1} \stackrel{\text{def}}{=} \mu(x)}_{\text{group inverse}} \quad \underbrace{x \cdot y \stackrel{\text{def}}{=} \mu(xy)}_{\text{group operation}}$$

The axioms of a group are easily checked, e.g.

$$\begin{aligned}
 x \cdot x^{-1} &= && \text{(definition of inverse)} \\
 x \cdot \mu(x) &= && \text{(unit followed by multiplication is the identity, i.e. axiom 4.4)} \\
 \mu(x) \cdot \mu(x) &= && \text{(definition of the group operation)} \\
 \mu(\mu(x)\mu(x)) &= && \text{(associativity of multiplication, i.e. axiom 4.5)} \\
 \mu(xx) &= && \text{(equality in the free group)} \\
 \mu(\varepsilon) &= && \text{(definition of group identity)} \\
 &1
 \end{aligned}$$

For the converse, we observe that for every group G , its group multiplication can be extended uniquely to an operation of type $TG \rightarrow G$, and the resulting operation will be associative in the sense required by Eilenberg-Moore algebras. \square

4.1.1 Axioms of a monad

We now describe the axioms of a monad.

Functoriality. The first group of axioms says that the first two ingredients (the structures and substitutions) of a monad are a functor in the sense of category theory. This means that substitutions preserve the identity and composition of functions. Preserving identity means that if we apply T to the identity function on X , then the result is the identity function on TX . Preserving composition means that the composition of substitutions is the same as the substitution of their composition, i.e. for every functions $f : X \rightarrow Y$ and $g : Y \rightarrow Z$, the following diagram commutes

$$\begin{array}{ccc}
 TX & \xrightarrow{Tf} & TY \\
 & \searrow T(g \circ f) & \downarrow Tg \\
 & & TZ
 \end{array} \tag{4.1}$$

Naturality. The naturality axioms say that for every function $f : X \rightarrow Y$, the following diagrams commute.

$$\begin{array}{ccc}
 X & \xrightarrow{f} & Y \\
 \text{unit}_X \downarrow & & \downarrow \text{unit}_Y \\
 TX & \xrightarrow{Tf} & TY
 \end{array} \tag{4.2}$$

$$\begin{array}{ccc}
 \mathbb{T}\mathbb{T}X & \xrightarrow{\mathbb{T}\tau_f} & \mathbb{T}\mathbb{T}Y \\
 \downarrow \text{free multiplication on } X & & \downarrow \text{free multiplication on } Y \\
 \mathbb{T}X & \xrightarrow{\tau_f} & \mathbb{T}Y
 \end{array} \tag{4.3}$$

In the language of category theory, this means that the unit and free multiplication are natural transformations. Also, as we will see later on, the second naturality axiom (naturality of free multiplication) says that the substitution τ_f is a homomorphism between the free algebras $\mathbb{T}X$ and $\mathbb{T}Y$.

Associativity. We now turn to the most important monad axioms, which ensure the the Eilenberg-Moore algebras are well behaved. The main associativity axiom says that for every set X , the set $\mathbb{T}X$ equipped with free multiplication on X is a \mathbb{T} -algebra (we call this the free \mathbb{T} -algebra over X , or simply free algebra if the monad is clear from the context). By unravelling the definitions, this means that the following two diagrams commute:

$$\begin{array}{ccc}
 \mathbb{T}X & & \\
 \text{unit}_{\mathbb{T}X} \downarrow & \searrow \text{identity} & \\
 \mathbb{T}\mathbb{T}X & \xrightarrow{\text{free multiplication on } X} & \mathbb{T}X
 \end{array} \tag{4.4}$$

$$\begin{array}{ccc}
 \mathbb{T}\mathbb{T}\mathbb{T}X & \xrightarrow{\text{free multiplication on } \mathbb{T}X} & \mathbb{T}\mathbb{T}X \\
 \mathbb{T}(\text{free multiplication on } X) \downarrow & & \downarrow \text{free multiplication on } X \\
 \mathbb{T}\mathbb{T}X & \xrightarrow{\text{free multiplication on } X} & \mathbb{T}X
 \end{array} \tag{4.5}$$

Apart from the above two, there is one more associativity axiom, namely:

$$\begin{array}{ccc}
 \mathbb{T}X & & \\
 \mathbb{T}(\text{unit}_X) \downarrow & \searrow \text{identity} & \\
 \mathbb{T}\mathbb{T}X & \xrightarrow{\text{free multiplication on } X} & \mathbb{T}X
 \end{array} \tag{4.6}$$

This completes the axioms of a monad, and the definition of a monad.

Exercises

Exercise 106. Consider a monad T in the category of sets. For a binary relation R on a set X , define

$$R^T \subseteq (TX) \times (TX)$$

to be the binary relation on TX that is defined by

$$R^T = \{((T\pi_1)(t), (T\pi_2)(t)) : t \in TR\}.$$

Does transitivity of R imply transitivity of R^T ?

4.1.2 Homomorphisms and recognisable languages

A homomorphism between two T -algebras is any function between their underlying sets which is consistent with the multiplication operation.

Definition 4.4 (Homomorphism). Let T be a monad. A T -homomorphism is a function $h : A \rightarrow B$ on the underlying sets of two T -algebras A and B , which makes the following diagram commute

$$\begin{array}{ccc} TA & \xrightarrow{Th} & TB \\ \text{multiplication in } A \downarrow & & \downarrow \text{multiplication in } B \\ A & \xrightarrow{h} & B \end{array}$$

When the monad is clear from the context, we simply write homomorphism, instead of T -homomorphism. Again, the reader will recognise the notion of homomorphism for semigroups, monoids and \circ -semigroups.

In the rest of this section, we describe some basic properties of homomorphisms.

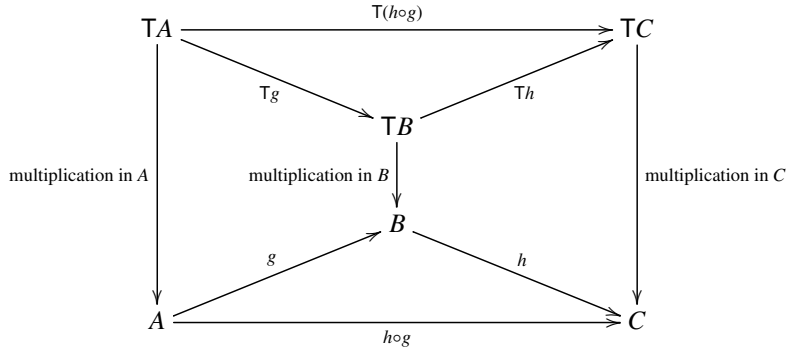
Lemma 4.5. *Homomorphisms are closed under composition.*

Proof Consider two homomorphisms

$$A \xrightarrow{g} B \xrightarrow{h} C.$$

Saying that the composition $h \circ g$ is a homomorphism is the same as saying

that the perimeter of the following diagram commutes:



The upper face commutes because of the functoriality axioms (substitutions are compatible with composition). The left and right faces commute by assumption that g and h are homomorphisms, and the lower face commutes by definition. \square

Recall that we defined the *free algebra over X* to be X equipped with free multiplication on X . The monad axioms say that this is indeed an algebra. It is called free because of the universal property given in the following lemma.

Lemma 4.6 (Free Algebra Lemma). *For every set X , the free algebra TX has the following universal property:*

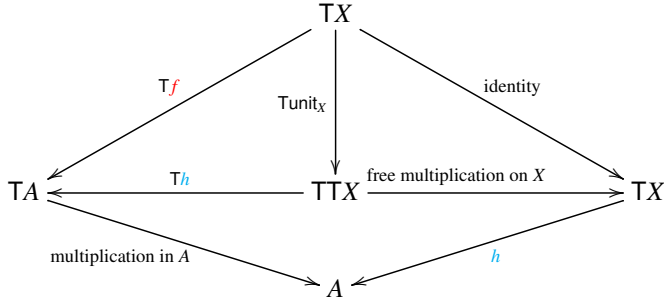
$$\forall f! \quad X \begin{array}{l} \xrightarrow{\text{function on sets}} A \\ \searrow \text{unit}_X \quad \nearrow \text{homomorphism of } T\text{-algebras} \\ \quad \quad \quad TX \end{array} \tag{4.7}$$

Proof We begin by showing that there is at least one blue homomorphism h for each red function f ; later we show that this homomorphism is unique. Let $\mu : TA \rightarrow A$ be multiplication in the algebra A . Define h to be the composition of the following functions:

$$TX \xrightarrow{Tf} TA \xrightarrow{\mu} A.$$

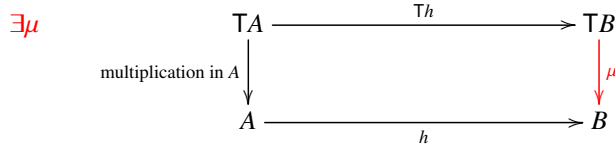
The axiom on naturality of free multiplication says that Tf is a homomorphism from the free algebra TX to the free algebra TA . The associativity axiom in the definition of an Eilenberg-Moore algebra says that multiplication μ is a homomorphism from the free algebra TA to the algebra A . Therefore, h is a homomorphism, as the composition of two homomorphisms Tf and μ .

We now show uniqueness – every homomorphism h which makes the diagram must be equal to the one described above. Consider the following diagram:



The upper left face commutes by applying T to the assumption that h extends f . (Applying T preserves commuting diagrams, because of the functoriality axioms.) The upper right face commutes by the first associativity axiom. The lower face commutes, because it says that h is a homomorphism. Therefore, the perimeter of the diagram commutes. The perimeter says that h must be equal to Tf followed by multiplication in A , and therefore h is unique. \square

Compositional functions. Fix a monad T . Suppose that A is an algebra, while B is a set, which is not (yet) equipped with a multiplication operation. We say that a function $h : A \rightarrow B$ is *compositional* if



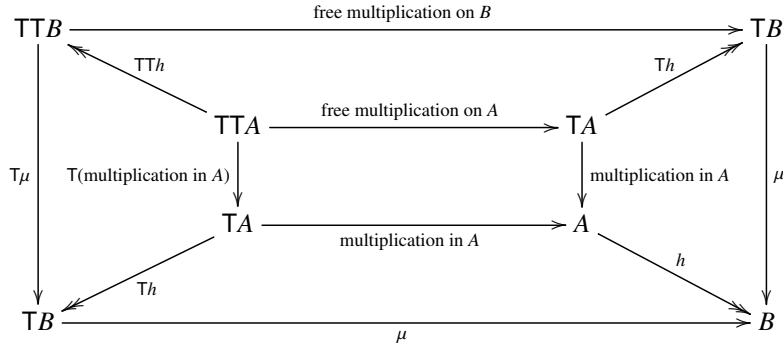
This is the same notion of compositionality as was used for monoids, semigroups and \circ -semigroups in part I of the book. For the same reason as before, surjective compositional functions are equivalent to surjective homomorphisms, as stated in the following lemma.

Lemma 4.7. *If A is an algebra, B is a set, and $h : A \rightarrow B$ is compositional and surjective, then there is a (unique) multiplication operation on B which turns it into an algebra and h into a homomorphism.*

Proof The multiplication operation – no surprises here – is μ from the definition of a compositional function. The diagram in the definition of a compositional function is the same diagram as in the definition of a homomorphism, and therefore if B equipped with μ is an algebra, then h is a homomorphism.

It remains to show that B equipped with μ is indeed an algebra. We only prove the more interesting of the two associativity diagrams.

We first observe that \mathbb{T} preserves surjectivity of functions. Indeed, if a function $h : A \rightarrow B$ is surjective, then it has a one-sided inverse, i.e. a function $h^{-1} : B \rightarrow A$ such that $h \circ h^{-1}$ is the identity on B . By the functoriality axioms, $\mathbb{T}h^{-1}$ is a one-sided inverse for $\mathbb{T}h$, and therefore $\mathbb{T}h$ is also surjective. This argument justifies the surjectivity annotation (double-headed arrows) in the following diagram.



The central face commutes by the assumption that A is an algebra. The upper face commutes by naturality of free multiplication. The right and lower faces commute by definition of a compositional function, and the left face commutes by the same definition with \mathbb{T} applied to it. It follows that all paths that begin in $\mathbb{T}TA$ and end in B denote the same function. Since h is surjective, it follows that the perimeter of the diagram commutes. This proves the second of the associativity diagrams in the definition of an Eilenberg-Moore algebra. \square

Recognisable colourings and languages. In this book, we are most interested in the Eilenberg-Moore algebras as recognisers of languages. A language is a subset L of a free algebra $\mathbb{T}\Sigma$. (Typically we are interested in the case where the alphabet Σ is finite, but this assumption does not seem to play a role in the results we care about, so we omit it.) A language is called *recognisable* if it is recognised by a finite algebra, as explained in the following definition (which uses a slightly more general notion of language, called colourings).

Definition 4.8 (Recognisable colourings). Fix a monad \mathbb{T} . An *algebra colouring* is defined to be any function from an algebra to a set of colours². A *finite*

² For some monads, it will be more useful to deviate from this definition. For example, in the monad from Example 23 that deals with vector spaces, a more useful notion of colouring is a linear map to the underlying field. Therefore, one could think of a parametrised notion of recognisability, where the notion of “algebra colouring” is taken as a parameter.

algebra is an algebra where the underlying set is finite³. An algebra colouring $\lambda : A \rightarrow U$ is called *recognisable* if it factors through a homomorphism into a finite algebra, as expressed in the following diagram:

$$\exists \quad \begin{array}{ccc} A & \xrightarrow{\lambda} & U \\ & \searrow \text{homomorphism} & \uparrow \text{colouring} \\ & \text{into a finite algebra} & B \end{array}$$

Note that a recognisable colouring will necessarily use finitely many colours.

A language can be viewed as the special case of an algebra where the algebra is a free algebra and there are two colours “yes” and “no”. For languages, we prefer set notation, e.g. we can talk about the complement of a language, or order languages by inclusion. The above definition is easily seen to coincide with the notions of recognisability for semigroups, monoids and \circ -semigroups that were discussed in the first part of this book. In the next section, we give more examples.

Exercises

Exercise 107. (2) For an algebra A with multiplication operation $\mu : TA \rightarrow A$, define its powerset as follows: the underlying set is the powerset PA , and multiplication is defined by

$$t \in TPA \quad \mapsto \quad \{\mu(s) : s \in {}^T t\},$$

where $\in {}^T$ is defined as in Exercise 106. Show an example of a monad T where this construction does not yield an algebra.

Exercise 108. (2) Does the group monad satisfy the following implication:

- (*) If $L \subseteq T\Sigma$ is recognisable, and $h : T\Sigma \rightarrow T\Gamma$ is a homomorphism, then $h(L)$ is recognisable.

What about surjective homomorphisms?

³ Like for algebra colourings, sometimes this notion of finite algebra is not the right one. In the monad from Example 23, the more useful notion is that a finite algebra is one where the underlying set is a vector space of finite dimension. Again, one could think of the notion of “finite algebra” as being a parameter.

Exercise 109. (2) Consider the implication in the previous exercise. Show that even if we restrict h to functions of the form $\top f$ for some surjective $f : \Sigma \rightarrow \Gamma$, then the implication can still be false in some monads.

4.2 A zillion examples

Monads have an abundance of interesting examples. This section is devoted to a collection of such examples, with an emphasis on the algebras arising from the monads, and the languages recognised by the finite algebras.

Example 15. [Finite multisets] Define $\top X$ to be the finite multisets over X , i.e. functions of type $X \rightarrow \mathbb{N}$ which have finite support (all but finitely many elements are mapped to 0). Functions are lifted to multisets point-wise, e.g.

$$\underbrace{\{x_1, \dots, x_n\}}_{\text{red brackets indicate multisets}} \xrightarrow{\top f} \{f(x_1), \dots, f(x_n)\}.$$

Another perspective on finite multisets is that they are finite words modulo commutativity $xy = yx$. The unit is $x \mapsto \{x\}$, and free multiplication is simply removing nested brackets, e.g.

$$\{\{x, y\}, \{z\}\} \mapsto \{x, y, z\}.$$

This is a monad. An algebra over this monad is the same thing as commutative monoid. Recognisable languages over this monad are the same things are regular languages – in the usual sense – which are commutative, see Exercise 11.

If we lift the restriction on finite supports, then we do not get a monad. The problem is with the substitutions: if $f : X \rightarrow \{a\}$ is the constant function with an infinite domain, then there is no way to define

$$(\top f)\{\underbrace{x_1, x_2, \dots}_{\text{infinitely many distinct elements}}\}.$$

The problem is that the above multiset should contain a infinitely many times. To overcome this problem, we would need to allow the multisets to have infinitely many copies of an element. \square

Example 16. [Idempotent finite words] Define $\top X$ to be finite words X^* , modulo the equation $ww = w$. This equation can be applied to arbitrary words, e.g.

$$abcabc = abc.$$

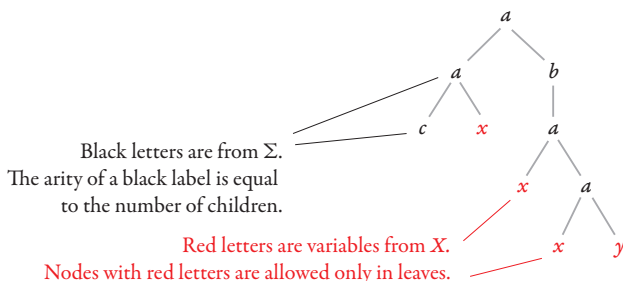
The remaining ingredients of the monad are defined in the natural way. An algebra over this monad is the same thing as an idempotent monoid, i.e. a monoid where all elements are idempotent. Green and Rees show that if X is a finite set, then TX is finite⁴. It follows that for every finite alphabet, there are finitely many languages over this alphabet, and all of them are recognisable. \square

Example 17. [Powersets] The *powerset monad*, and its variant the *finite powerset monad*, are defined in the same way as the multiset monad, except that we use sets (or finite sets) instead of multisets. The substitutions are defined via images (in the language of category theory, we use the covariant powerset functor):

$$A \subseteq X \quad \xrightarrow{Tf} \quad \{f(x) : x \in A\} \subseteq Y.$$

The algebras over the finite powerset monad are the same thing as monoids that are commutative and idempotent. If X is a finite set, then both powerset monads generate finite sets; and therefore all languages over finite alphabets are recognisable. \square

Example 18. [Terms] Fix a ranked set Σ , i.e. a set where every element has an associated arity in $\{0, 1, \dots\}$. Define TX to be the terms over Σ with variables X , i.e. an element of TX is a tree that looks like this:



The unit operation maps $x \in X$ to a term which consists only of x . The substitution Tf is defined in the natural way, by applying f to the variables and leaving the remaining part of the term unchanged. Finally, free multiplication replaces each variable with the corresponding term. It is a simple exercise to check that a T -algebra is the same thing as an *algebra of type Σ* , in the sense

⁴ [20] Green and Rees, “On semi-groups in which $x^r = x$ ”, 1952, p. 35

of universal algebra⁵. In the terminology of automata theory, both of these notions are the same as deterministic bottom-up tree automata over finite trees, where Σ is the input alphabet⁶. From the above observation it follows that a language $L \subseteq TX$ is recognisable in the sense of Definition 4.8 if and only if it is a regular tree language in the sense of automata theory⁷, where the input alphabet is obtained from Σ by adding one letter of arity 0 for each element of X . If the ranked set Σ contains only letters of arity exactly one, then a T -algebra can be seen as a deterministic word automaton with input alphabet Σ , without distinguished initial and final states. \square

Example 19. [Linear orders] Define a *chain* over a set X to be a linear order with positions labelled by X , modulo isomorphism of labelled linear orders. For an infinite cardinal κ , a monad T_κ is defined as follows. The set $T_\kappa X$ consists of chains over X , which have cardinality strictly less than κ . The monad structure is defined in the same way as for \circ -words (which are the special case of this construction for where κ is the first uncountable cardinal, which is the same as \aleph_1 assuming the Continuum Hypothesis). Nevertheless, we give a more exact description below.

For a function f , the corresponding substitution $T_\kappa f$ is defined in the natural way, by applying f to the labels in a chain and leaving the positions and ordering unchanged. The unit maps a letter to the one letter chain with that letter. The free multiplication operation is defined using lexicographic products, as follows. Suppose that $w \in T_\kappa T_\kappa X$. The positions in the free multiplication of w are pairs (i, j) such that i is a position of w and j is a position in the label $w(i) \in T_\kappa X$ of position i in the chain w . The label of such a position is inherited from j , and the ordering is lexicographic. The cardinality of the resulting chain is at most κ , since every infinite cardinal satisfies $\kappa = \kappa^2$.

⁵ An algebra of type Σ is defined to be an underlying set A , together with an interpretation which maps each n -ary symbol from Σ to an operation of type $A^n \rightarrow A$. See

[32] Sankappanavar and Burris, "A course in universal algebra", 1981, Definition 1.3

⁶ Thatcher and Wright, "Generalized Finite Automata Theory with an Application to a Decision Problem of Second-Order Logic", 1968, Section 2

⁷ This monad describes finite trees. Finding an algebraic account of languages of infinite trees remains an open problem This problem is discussed in the following papers:

[4] Blumensath, "Regular Tree Algebras", 2018

[6] Bojańczyk and Klin, "A non-regular language of infinite trees that is recognizable by a sort-wise finite algebra", 2019

We only prove one of the monad axioms:

$$\begin{array}{ccc}
 T_\kappa T_\kappa T_\kappa X & \xrightarrow{\text{free multiplication on } T_\kappa X} & T_\kappa T_\kappa X \\
 \downarrow T_\kappa(\text{free multiplication on } X) & & \downarrow \text{free multiplication on } X \\
 T_\kappa T_\kappa X & \xrightarrow{\text{free multiplication on } X} & T_\kappa X
 \end{array}$$

Let $w \in T_\kappa T_\kappa T_\kappa X$. If, in the diagram above, we first go right and then down, then the resulting linear order will have positions of the form $((i, j), k)$, where i is a position in w , j is a position in $w(i)$, and k is a position in $w(i)(j)$. If, in the diagram, we first go down and then right, then we get positions of the form $(i, (j, k))$, where i, j, k satisfy the same conditions as above. In both cases, the tuples of positions are ordered lexicographically, and the label is inherited from k . Therefore

$$((i, j), k) \mapsto (i, (j, k))$$

is an isomorphism of labelled linear orders, and hence the two outcomes are equal as chains.

If we take $\kappa = \aleph_0$, then we get the monad of finite words. If we take κ to be the first uncountable cardinal, then we get the monad corresponding to ω -words.

If we take κ to be the first cardinal bigger than \mathfrak{c} , then T_κ describes chains of cardinality at most \mathfrak{c} . In this case, we have the following phenomenon. Recall the powerset construction that was described in Section 3.3. This construction also makes sense for chains of size at most \mathfrak{c} . Define \mathcal{X} to be the least class of T_κ -algebras which contains the syntactic algebra of the language “every a is before every b ”, and which is closed under products and the powerset construction. Using the same proof as for Lemma 3.22, one can show that every mso definable language $L \subseteq T_\kappa \Sigma$ is recognised by an algebra from \mathcal{X} . Since satisfiability for mso over the reals is undecidable⁸, it follows that there is no finite way of representing algebras from \mathcal{X} . This means that the powerset construction over finite T_κ -algebras is not computable. \square

Example 20. Consider a class \mathcal{X} of linear orders which have size bounded by some cardinal κ , and which are closed under free multiplication as defined in the previous example, when viewed as chains over a one letter alphabet. If we restrict the monad from the previous example to chains where the underlying linear order is in \mathcal{X} , then we also get a monad. This construction yields the

⁸ [35] Shelah, “The Monadic Theory of Order”, 1975, Theorem 7

following monads (in all cases, we assume some fixed upper bound on κ on the cardinality, e.g. we can require countability):

- well-founded words (the class of well-founded linear orders);
- scattered words (the class of scattered orders, i.e. those into which one cannot embed the rational numbers)⁹.
- dense words (the class which contains two orders: a singleton order for units, and the rational numbers).

□

Example 21. [ω -semigroups] We now describe a monad that corresponds to ω -semigroups, see Definition 3.7. Since an ω -semigroup has two sorts, we leave the category of sets, and use instead the category

$$\mathbf{Set}^{(+, \omega)}$$

of sets with two sorts $+$ and ω . An object in this category is a set, where every element is assigned exactly one of two sorts, called $+$ and ω . A morphism in this category is any sort-preserving function between sorted sets. We use the name *sorted set* for the objects. We write the objects and morphisms of this category in red, to distinguish them from usual sets and functions. We also use the following notation for sorts:

$$\underbrace{X}_{\text{a sorted set}} = \underbrace{X[+]_{\text{elements of sort } +}} \cup \underbrace{X[\omega]_{\text{elements of sort } \omega}}.$$

We define a monad \mathbf{T} over this category as follows. For a sorted set X , the sorted set $\mathbf{T}X$ is defined by:

$$(\mathbf{T}X)[+] = X[+]^+ \quad \cup \quad (\mathbf{T}X)[\omega] = (X[+])^*(X[\omega]) \cup (X[+])^\omega.$$

For a morphism $f : X \rightarrow Y$, the substitution morphism $\mathbf{T}f$ is defined in the natural way, by applying f to every letter. The unit and free multiplication are defined in the natural way as well. (An element of sort $+$ in $\mathbf{T}X$ is simply a finite nonempty word of finite nonempty words over $X[+]$, and we can use free multiplication from the monad of finite nonempty words. On sort ω , there are more cases to consider, but the definition is natural as well.) An Eilenberg-Moore algebra over this monad is the same thing as an ω -semigroup, as defined at the end of Section 3.1. □

⁹ Algebras for the monad of countable scattered words are studied in
 [31] Rispal and Carton, “Complementation of Rational Sets on Countable Scattered Linear Orderings”, 2005

Example 22. [Vector spaces] Define $\mathbb{T}X$ to be the vector space, over the field of rational numbers, where the basis is X . In other words, elements of $\mathbb{T}X$ are finite linear combinations of elements from X with rational coefficients. For example,

$$3x + 7y - 0.5z \in \mathbb{T}\{x, y, z\}.$$

The action of \mathbb{T} on functions is defined by

$$q_1x_1 + \cdots + q_nx_n \quad \xrightarrow{\mathbb{T}f} \quad q_1f(x_1) + \cdots + q_nf(x_n).$$

The unit operation maps $x \in X$ to the corresponding basis vector, and free multiplication is defined in the natural way, as illustrated in the following example:

$$3(4x + 0.5y) - 0.2(5x - 0.1y) \quad \mapsto \quad 12x + 1.5y - x + 0.02y = 11x + 1.52y.$$

An algebra A over this monad, with multiplication μ , is also equipped with the structure of a vector space, because we can add elements

$$a + b \stackrel{\text{def}}{=} \mu(a + b)$$

and multiply them by scalars $q \in \mathbb{Q}$

$$qa \stackrel{\text{def}}{=} \mu(qa).$$

If $B \subseteq A$ is a basis for the vector space A , then the algebra A is isomorphic to $\mathbb{T}B$. Therefore, over this monad, every algebra is isomorphic to a free algebra. \square

Example 23. [Algebra over a field] Define $\mathbb{T}X$ to be finite linear combinations of words in X^* , with rational coefficients. For example,

$$2xyx + -2xx + 0.5xyz \in \mathbb{T}\{x, y, z\}.$$

In other words, $\mathbb{T}X = \mathbb{T}_{\text{vec}}X^*$, where \mathbb{T}_{vec} is the monad of vector spaces from Example 22 and X^* is the monad of finite words¹⁰. On functions, the monad acts as follows

$$q_1x_1 + \cdots + q_nx_n \quad \xrightarrow{\mathbb{T}f} \quad q_1f^*(x_1) + \cdots + q_nf^*(x_n),$$

where f^* is the substitutions in the monad of finite words. The unit maps x to the linear combination which has the one-letter word x with coefficient 1.

¹⁰ This is an example of a composite monad that arises via a distributive law of two monads.

This type of construction was first described in

[1] Appelgate et al., *Seminar on triples and categorical homology theory*, 1969, Chapter on distributive laws

The free multiplication is defined like for polynomials, but the variables are not commuting, e.g.:

$$3(4x - 2y)(2xy + yy) \mapsto 24x.xy + \underbrace{12xyy - 12yxy - 6yyy}_{\text{this is not 0}}.$$

Every algebra over this monad has the structure of a vector space over the rationals, but there is more structure (e.g. one can multiply two elements of the algebra)¹¹.

What is a recognisable colouring over this monad? In the context of this monad (and also the simpler monad of vector spaces from Example 22), it is more useful to work with different notions of “finite algebra” and “algebra colouring”: instead of finite algebras, one should consider finite dimensional algebras (i.e. those where the underlying vector space has finite dimension), and instead of algebra colourings one should consider linear maps to vector spaces. Under these adapted definitions, the algebra colourings recognised by finite algebras are exactly those which are recognised by weighted automata, see Exercise 122. \square

Exercises

Exercise 110. (2) Consider the monad T_Σ from Example 18, where Σ is some ranked set (possibly infinite). Let X be some possibly infinite set of variables, and consider a set

$$\mathcal{E} \subseteq \underbrace{(T_\Sigma X) \times (T_\Sigma X)}_{\text{elements of this set will be called identities}}.$$

elements of this set will be called identities

For a set Y , define \sim to be the least congruence on $T_\Sigma Y$ that satisfies

$$(T_\Sigma f)(t_1) \sim (T_\Sigma f)(t_2) \quad \text{for every } (t_1, t_2) \in \mathcal{E} \text{ and } f : X \rightarrow Y.$$

(This congruence can be obtained by intersecting all congruences with the above property.) Define a new monad as follows: TY is equal to $T_\Sigma Y$ modulo \sim , and the remaining components of the monad are defined in the natural way. Show that this is a monad.

¹¹ Algebras over this monad are known as “algebras over the field of rational numbers”, but we avoid this terminology due to the over-loading of “algebra over”.

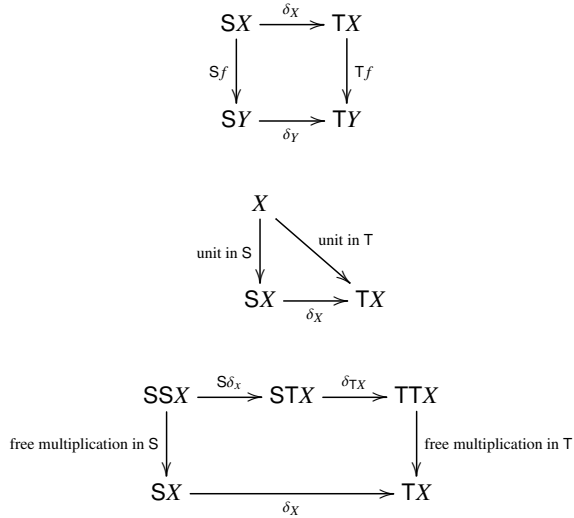


Figure 4.1 These three diagrams should commute for all sets X and all functions $f : X \rightarrow Y$. The top diagram says that $\{\delta_X\}_X$ is a natural transformation, while the bottom two diagrams say that it is compatible with unit and free multiplication.

Exercise 111. (2) For monads S and T , define a monad morphism from S to T to be a family of functions

$$\{\delta_X : SX \rightarrow TX\}_{X \text{ is a set}}$$

which is subject to the axioms in Figure 4.1. Using the monads from Section 4.2, give five examples of monad morphisms, and five examples of pairs of monads which do not allow a monad morphism.

Exercise 112. (1) We say that $w \in T\Sigma$ is *regular* if $\{w\}$ is a recognisable language. Find a monad where there are no regular elements. (Hint: it appears in this section.)

Exercise 113. (1) What is the monad for rings (commutative and non-commutative)? Semirings?

Exercise 114. (2) Consider the following monad T . The set TX is the set of ω -words X^ω , and the substitution $Tf : TX \rightarrow TY$ is defined coordinate-wise.

The unit is $x \mapsto x^\omega$, and free multiplication is defined by

$$w \in \text{TTX} \quad \mapsto \quad (i \mapsto \underbrace{(w[i])[i]}_{\substack{i\text{-th letter of} \\ \text{of } i\text{-th letter of } w}}).$$

Show that this is a monad. Also, show that a language $L \subseteq \text{T}\Sigma$ is recognisable if and only if it is clopen in the sense of Exercise 65.

Exercise 115. (1) Let \mathbb{T} be the monad of countable well founded. Show that a finite algebra with universe S is uniquely determined by the operations

$$\begin{array}{cc} \underbrace{ab}_{\substack{\text{binary} \\ \text{product}}} & \underbrace{a^\omega}_{\substack{\text{product of} \\ aaa \dots}} \\ S^2 \rightarrow S & S \rightarrow S \end{array}$$

Exercise 116. (2) Consider the monad from Example 115. Show that a language $L \subseteq \text{T}\Sigma$ is definable in first-order logic (in the ordered model) if and only if it is recognised by a finite \mathbb{T} -algebra S where the underlying semigroup is aperiodic.

Exercise 117. (2) Consider the monad from Example 115. Consider regular expressions defined by the usual operators, plus L^ω . Show that these expressions do not describe all recognisable languages.

Exercise 118. (2) Consider the monad and regular expressions from Exercise 117. Given an effective condition on finite algebras which corresponds exactly to the Boolean combinations of regular expressions.

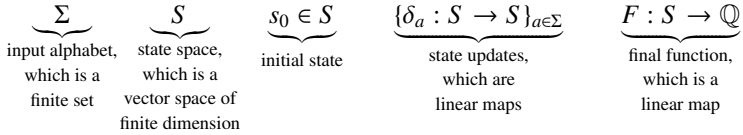
Exercise 119. (1) Let \mathbb{T} be the monad of countable scattered chains. Show that a finite algebra with universe S is uniquely determined by the operations

$$\begin{array}{ccc} \underbrace{ab}_{\substack{\text{binary} \\ \text{product}}} & \underbrace{a^\omega}_{\substack{\text{product of} \\ aaa \dots}} & \underbrace{a^{\omega*}}_{\substack{\text{product of} \\ \dots aaa}} \\ S^2 \rightarrow S & S \rightarrow S & S \rightarrow S \end{array}$$

Exercise 120. (2) Consider the monads from Examples 115 and 119. In which of these monads is first-order logic (over ordered models) equivalent to star-free expressions?

Exercise 121. (2) Consider the monad from Example 119. Which class of languages corresponds to aperiodicity (of the semigroup underlying the T-algebra)?

Exercise 122. (2) A weighted automaton over the rationals consists of:



The semantics of this automaton is a function of type $\Sigma^* \rightarrow \mathbb{Q}$ defined as follows. Given an input word, do the following: start with the initial state, apply the state update for the first letter, then the state update for the second letter, and so on for all letters, and at the end apply the final function. The semantics of a weighted automaton can also be naturally extended to finite linear combinations of words over Σ , i.e. to elements of the $T\Sigma$ as in Example 23.

Show that $L : T\Sigma \rightarrow \mathbb{Q}$ is recognised by a weighted automaton if and only if it is recognised by a finite dimensional T-algebra.

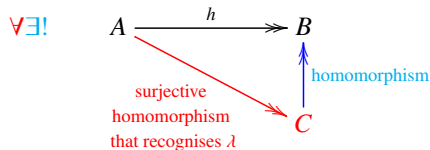
4.3 Syntactic algebras

In this section, we show that every recognisable algebra colouring has a syntactic homomorphism, i.e. a recognising homomorphism that stores the minimal amount of information. This can be viewed as a monad version of the Myhill-Nerode theorem.

Definition 4.9 (Syntactic homomorphism). The *syntactic homomorphism* of an algebra colouring $\lambda : A \rightarrow U$ is any surjective homomorphism

$$h : A \rightarrow B$$

which recognises λ and which is minimal in the sense explained in the following quantified diagram



The algebra used by the syntactic homomorphism is called the *syntactic algebra*. The syntactic algebra, if it exists, is unique up to isomorphism of algebras. Also the syntactic homomorphism is unique in the following sense: every two syntactic homomorphisms will have the same kernel (equivalence relation on A that identifies two elements with the same homomorphic image). This kernel is called the *syntactic congruence* of λ .

We are mainly interested in the case where λ describes a language, i.e. A is a free algebra, and there are two colours.

There are two main results in this section. The first one, Theorem 4.11, says that if an algebra colouring is recognisable, then it has a syntactic homomorphism. The second one, Theorem 4.16, says that a monad is finitary (roughly speaking, this means that all structures described by the monad are finite) if and only if every (not necessarily recognisable) algebra colouring has a syntactic homomorphism. To illustrate these theorems, we begin with an example of an algebra colouring that does not have a syntactic homomorphism. By Theorem 4.11, this colouring is necessarily not recognisable.

Example 24. Consider the monad of \circ -words. Define L to be the set of \circ -words over a one letter alphabet which contain every finite word as an infix. More formally,

$$L = \{w \in \{a\}^\circ : a^n \text{ is an infix of } w \text{ for every } n \in \{0, 1, \dots\}\},$$

which is a language over alphabet a , and therefore a special case of an algebra colouring. This language is not recognisable, because all finite words must have different images under any recognising homomorphism (we leave this as an exercise for the reader). We will show that L does not have a syntactic homomorphism. For $n \in \{1, 2, \dots\}$ define

$$w_n = (\text{shuffle of } \{a\}) \cdot a^n \cdot (\text{shuffle of } \{a\})$$

and define h_n to be the function

$$w \in \{a\}^\circ \quad \mapsto \quad \begin{cases} w_n & \text{if } w = w_{n+1} \\ w & \text{otherwise.} \end{cases}$$

This function is compositional for every n , and therefore it can be viewed as a homomorphism. If there would be a syntactic homomorphism, then it would need to factor through h_n . Since h_n gives the same result for w_n and w_{n+1} , therefore the same would have to be true for the syntactic homomorphism. Therefore, the syntactic homomorphism h , if it existed, would need to satisfy

$$h(w_1) = h(w_2) = \dots$$

By associativity, we would have

$$h(\underbrace{w_1 w_1 w_1 \cdots}_{\notin L}) = h(\underbrace{w_1 w_2 w_3 \cdots}_{\in L}).$$

Therefore the syntactic homomorphism does not exist. \square

As we will see later in this section, the monad of \circ -words from the above example is not finitary.

4.3.1 Terms, polynomials and congruences

To prove the existence of syntactic homomorphisms, we will use classical concepts from universal algebra such as a terms, polynomials and congruences. We begin by defining these notions.

Terms and polynomials. If X is a set, then a *term over variables X* is defined to be any element of TX . Given an algebra A , a term is interpreted as the following operation

$$\eta \in A^X \quad \mapsto \quad \text{multiply } (\text{T}\eta)(t) \text{ in } A.$$

We write t^A for the above operation; and we use the name *variable valuation* for η . An operation of this form is called a *term operation* in the algebra A . We distinguish between a term (which can be viewed as syntax) and the term operation that it generates in a given algebra (which can be viewed as semantics). If a term uses a finite set of variables $\{x_1, \dots, x_n\}$, then for $a_1, \dots, a_n \in A$ we write

$$t^A(a_1, \dots, a_n)$$

for the result of applying t^A to the variable valuation $x_i \mapsto a_i$.

Example 25. Consider the monad of finite words. The word xy is a term, and the term operation induced it in an algebra (which is the same thing as a monoid) is binary multiplication. The operation induced by the term ε , which has an empty set of variables, is the constant that represents the monoid identity. Another example of a term operation is squaring, which is given by the term xx . A non-example is the idempotent power operation $a \mapsto a^!$. This is not a term operation, because the number $!$ depends on the algebra at hand (also, this number does not exist in some infinite algebras).

In the monad of \circ -chains, the Lauchli-Leonard operations x^ω and $x^{\omega*}$ are also term operations. To model the remaining Lauchli-Leonard operation, namely

shuffling, we use an infinite family of terms, with the n -th one being the shuffle $\{x_1, \dots, x_n\}$. \square

A *polynomial* in an algebra is like a term except that some variables are fixed elements of the algebra. More formally, a *polynomial with variables X in an algebra A* is defined to be any element $t \in T(A + X)$. The semantics of such a polynomial is the operation

$$\eta \in A^X \quad \mapsto \quad t^A(\underbrace{\text{id}_A + \eta}_{\eta \text{ extended by the identity on } A}).$$

Such an operation is called a *polynomial operation* in A . We will be mainly interested in unary polynomial operations, which are functions of type $A \rightarrow A$ that arise from polynomials with one variable.

Congruences. Consider an equivalence relation \sim on the underlying set in an algebra A . It is not hard to see that the following conditions are equivalent:

- \sim is the kernel of some homomorphism from A to some algebra B ;
- the function which maps $a \in A$ to its equivalence class is compositional;
- \sim is compatible with every term operation $f : A^X \rightarrow A$, which means:

$$\underbrace{\eta_1 \sim \eta_2}_{\substack{\eta_1(x) \sim \eta_2(x) \\ \text{for every } x \in X}} \quad \Rightarrow \quad f(\eta_1) \sim f(\eta_2) \quad \text{for every } \eta_1, \eta_2 \in A^X.$$

Equivalence of the first two conditions is Lemma 4.7, and equivalence of the second two conditions follows by unravelling the definition of compositionality. We say that \sim is *congruence* if it satisfies any of the above conditions. Every congruence induces a quotient algebra, where the universe is equivalence classes.

Contextual equivalence. To construct the syntactic homomorphism, we will use contextual equivalence, which is defined similarly as in the usual proof of the Myhill-Nerode Theorem, see e.g. Theorem 1.7. The only difference is that instead of two-sided contexts (as in monoids and semigroups), we use unary polynomial operations.

Definition 4.10 (Contextual equivalence). Define *contextual equivalence* of an algebra colouring $\lambda : A \rightarrow U$ to be the equivalence relation on A , which identifies $a, a' \in A$ if they have the same image under $\lambda \circ f$ for every unary polynomial operation $f : A \rightarrow A$.

The following example shows that, in general, contextual equivalence is not a congruence.

Example 26. Consider the language L from Example 24. We show that contextual equivalence for this language (viewed as a colouring with values “yes” and “no”) is not a congruence. A unary polynomial operation in $\{a\}^\circ$ corresponds to a \circ -word over alphabet $\{a, x\}$. From the point of view of L , there are two kinds of unary polynomial operations. If f is a unary polynomial such that the corresponding \circ -word over $\{a, x\}$ contains a^n as an infix for every n , then $f(w) \in L$ for every $w \in \{a\}^\circ$. Otherwise, if the corresponding \circ -word does not contain a^n as an infix for some n , then $f(w) \in L \Leftrightarrow w \in L$. These observations imply that contextual equivalence for L has two equivalence classes: namely L and its complement. In particular, contextual equivalence is not a congruence, since otherwise L would be recognisable. \square

Exercises

Exercise 123. (1) Fix a monad. Let A be an algebra, and let $p : A^X \rightarrow A$ be a polynomial operation. Show that for every homomorphism $h : A \rightarrow B$ there is a polynomial operation $p^h : B^X \rightarrow B$ which commutes with the homomorphism in the following sense:

$$\begin{array}{ccc} A^X & \xrightarrow{p} & A \\ \downarrow h^X & & \downarrow h \\ B^X & \xrightarrow{p^h} & B \end{array}$$

Exercise 124. (2) Let T be a monad in the category of sets, and consider a set of terms \mathcal{B} , each one using finitely many variables. We say that \mathcal{B} is a *term basis* if for every finite algebra A and subset $\Gamma \subseteq A$, the sub-algebra generated by Γ is equal to the least subset of A that contains Γ and which is closed under applying terms from \mathcal{B} . Show that if \mathcal{B} is a term basis, then a finite algebra A is uniquely determined by its \mathcal{B} -multiplication tables, which is the family of term operations $\{t^A\}_{t \in \mathcal{B}}$.

Exercise 125. (1) Let T be a monad which has a finite term basis \mathcal{B} . Show

that given the \mathcal{B} -multiplication tables in an algebra A , one can compute the \mathcal{B} -multiplication tables of the powerset algebra PA (as defined in Exercise 107, assuming that PA is indeed an algebra).

Exercise 126. (2) Find a notion of *computable term basis* which generalises the previous exercise so as to capture the L\"auchli-Leonard operations in the monad of \circ -words.

4.3.2 Syntactic homomorphisms for recognisable colourings

We now state the first result about syntactic homomorphisms, which says that they always exist for algebra colourings that are recognisable.

Theorem 4.11. *Let T be a monad in the category of sets. Every recognisable algebra colouring has a syntactic homomorphism.*

If \sim is a equivalence relation on the underlying set of an algebra A , then we say that \sim recognises an algebra colouring $\lambda : A \rightarrow U$ if all equivalence classes are monochromatic, i.e.

$$a_1 \sim a_2 \quad \text{implies} \quad \lambda(a_1) = \lambda(a_2).$$

An algebra colouring is recognisable if and only if there is equivalence relation that recognises it, is a congruence, and has finitely many equivalence classes.

We say that an equivalence relation \approx refines an equivalence relation \sim if $a \approx b$ implies $a \sim b$. If we view equivalence relations as sets of pairs, this is the same as set inclusion. The following lemma shows that contextual equivalence refines every recognising congruence. (Note that the lemma does not say that contextual equivalence is a congruence. This will be proved later, using the assumption on recognisable.)

Lemma 4.12. *For every algebra colouring λ , contextual equivalence of λ is an equivalence relation that refines every congruence that recognises λ .*

Proof Let f be a unary polynomial operation in the underlying algebra. Let \approx be a congruence that recognises λ . As a congruence, \approx is compatible with every term operation, and therefore it is compatible with every unary polynomial operation (see Exercise 123). This means that if $a \approx b$, then also $f(a) \approx f(b)$, and therefore $f(a)$ has the same colour under λ as $f(b)$, by assumption that \approx is a recognising congruence. By arbitrary choice of f , we have established that $a \approx b$ implies $a \sim b$. \square

We now use the assumption on recognisability to prove that contextual equivalence is a congruence.

Lemma 4.13. *If an algebra colouring λ is recognisable, then its contextual equivalence is a congruence.*

Proof Let $\lambda : A \rightarrow U$ be an algebra colouring, and let \sim be its contextual equivalence relation. We will show that \sim is a congruence in three steps.

- (1) In the first step, we show that \sim is compatible with all unary polynomial operations. The key observation is that unary polynomial operations are closed under composition. Indeed, consider two unary polynomial operations

$$A \xrightarrow{f} A \xrightarrow{g} A.$$

Take the unary polynomial in $\mathbb{T}(A + \{x\})$ that defines f , and for the variable x substitute the unary polynomial in $\mathbb{T}(A + \{x\})$ that defines g . The resulting unary polynomial defines the composition $g \circ f$.

Equipped with this observation, we prove that \sim is compatible with all unary polynomial operations. Let f be a unary polynomial operation. We need to show

$$a_1 \sim a_2 \quad \text{implies} \quad f(a_1) \sim f(a_2).$$

The assumption of the implication says that a_1, a_2 have the same image under all functions of the form $\lambda \circ g$, where g is a unary polynomial operation, while the conclusion of the implication says that a_1, a_2 have the same image under all functions of the form $\lambda \circ g \circ f$, where g is a unary polynomial operation. Because $g \circ f$ is a unary polynomial operation, the implication follows.

- (2) We now show that \sim is compatible with all term operations that have finitely many variables. Consider an n -ary term operation $f : A^n \rightarrow A$. We need to show

$$\bigwedge_{i \in \{1, \dots, n\}} a_i \sim b_i \quad \text{implies} \quad f(a_1, \dots, a_n) \sim f(b_1, \dots, b_n). \quad (4.8)$$

Assume the assumption of the above implication. For $i \in \{1, \dots, n\}$ consider the following unary polynomial operation

$$f_i(x) = f(a_1, \dots, a_{i-1}, x, b_{i+1}, \dots, b_n).$$

Because $a_i \sim b_i$ and \sim is compatible with all unary polynomial operations,

$$f(a_1, \dots, a_{i-1}, b_i, \dots, b_n) = f_i(a_i) \sim f_i(b_i) = f(a_1, \dots, a_i, b_{i+1}, \dots, b_n).$$

By iterating the above for $i = 1, \dots, n$, we get the conclusion of (4.8).

- (3) Finally, we show that \sim is compatible with all term operations, possibly with infinitely many variables. In this step, we use the assumption that λ is recognisable. Consider a term $t \in TX$. We need to show that every valuations $\eta_1, \eta_2 \in A^X$ satisfy

$$\eta_1 \sim \eta_2 \quad \text{implies} \quad t^A(\eta_1) \sim t^A(\eta_2).$$

Assume the assumption of the above implication. Because λ is recognisable, there is some congruence \approx on A which recognises λ and has finitely many equivalence classes. Consider a choice function $\tau : A \rightarrow A$ which maps each \approx -equivalence class in A to a chosen element in that equivalence class. Let $i \in \{1, 2\}$. Because the choice function respects \approx -equivalence, we have

$$\eta_i \approx \tau \circ \eta_i. \quad (4.9)$$

Since \approx is a congruence,

$$t^A(\eta_i) \approx t^A(\tau \circ \eta_i). \quad (4.10)$$

Since each of the valuations $\tau \circ \eta_i$ for $i = 1, 2$ has finite image, we can find a finite set of variables Y (think of pairs of equivalence classes of \approx) and functions δ and γ_i which make the following diagram commute:

$$\begin{array}{ccc} & X & \\ \tau \circ \eta_1 \swarrow & \downarrow \sigma & \searrow \tau \circ \eta_2 \\ A & Y & A \\ \gamma_1 \longleftarrow & & \longrightarrow \gamma_2 \end{array} \quad (4.11)$$

Define $s \in TY$ to be the term that results from t by the substitution $T\sigma$. By definition,

$$s^A(\gamma_i) = t^A(\gamma_i \circ \sigma) = t^A(\tau \circ \eta_i) \stackrel{(4.10)}{\approx} t^A(\eta_i).$$

By the first step of the proof, \approx refines \sim . Therefore, to finish the proof, it remains to show

$$s^A(\gamma_1) \sim s^A(\gamma_2).$$

Since s is a term with finitely many variables, and we have already established that \sim is compatible with terms that have finitely many variables, it is enough to show $\gamma_1 \sim \gamma_2$. This is because

$$\tau \circ \eta_1 \sim \eta_1 \sim \eta_2 \sim \tau \circ \eta_2.$$

In the above, the first and last equivalence are due to (4.9) and the fact that \approx refines \sim . Since the first valuation in the above is equal to $\gamma_1 \circ \sigma$, and the last one is equal to $\gamma_2 \circ \sigma$, it follows that $\gamma_1 \sim \gamma_2$.

We have proved that \sim is a congruence. □

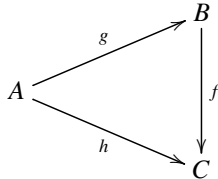
We now complete the proof of Theorem 4.11. Let \sim be contextual equivalence of a recognisable algebra colouring $\lambda : A \rightarrow U$. By Lemma 4.13, \sim is a congruence. Because \sim is a congruence, the quotient function, call it h , is a homomorphism. Since the identity function is a unary polynomial operation, it follows that equivalence classes of \sim are monochromatic, and therefore h recognises λ .

It remains to show that every other recognising homomorphism factors through h . Suppose that

$$g : A \rightarrow B$$

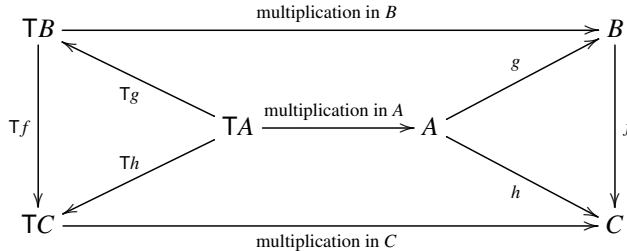
is a surjective homomorphism that recognises λ . The kernel of g is a congruence that recognises λ , and therefore the kernel of g refines \sim . In other words, g factors through h , i.e. there is some function f such that $g = f \circ h$. The last thing to show is that f is in fact a homomorphism. This is shown in the following lemma, with C being the quotient A/\sim .

Lemma 4.14. *Let A, B, C be algebras, let g, h be surjective homomorphisms, and let f be a function which makes the following diagram commute.*



Then f is a homomorphism.

Proof Consider the following diagram:



The upper and lower faces commute because g and h are homomorphisms, and the left and right faces commute by definition of f . Since all arrows in the diagram are surjective, it follows that the perimeter of the diagram commutes, which means that f is a homomorphism. □

Exercises

Exercise 127. (2) Consider the monad from Example 23. Show that if $\lambda : T\Sigma \rightarrow \mathbb{Q}$ is recognised by a weighted automaton, then the same is true for the syntactic homomorphism.

4.3.3 Finitary monads

In the monad of finite words, every language – even not necessarily recognisable – has a syntactic homomorphism. In the monad of \circ -words, this is no longer true, as witnessed by Example 24. What is the difference?

The difference is that every finite word uses only a finite subset of the alphabet, which is no longer true for \circ -words. This is made precise by the following definition.

Definition 4.15 (Finitary elements and monads.). Let T be a monad in the category of sets. We say that an element $t \in TX$ is *finitary*

$$t = (Tf)(t) \quad \text{for some } f : X \rightarrow X \text{ with finite image.}$$

We say that T is finitary if for every X , all elements of TX are finitary.

For example, the monad of finite words is finitary, while the monads of \circ -words is not. The following theorem shows that finitary monads are exactly those monads where all algebra colourings have syntactic homomorphisms.

Theorem 4.16. *Let T be a monad in the category of sets. Then T is finitary if and only if every algebra colouring has a syntactic homomorphism.*

Proof We begin with the left-to-right implication. Suppose that T is finitary, and let $\lambda : A \rightarrow U$ be some algebra colouring, not necessarily recognisable. We use the same proof as in Theorem 4.11. The only place where the proof used the assumption on recognisability was in step (3) of Lemma 4.13, where contextual equivalence was shown to be compatible with all term operations. If the monad is finitary, then every term operation uses finitely many variables, and therefore step (3) of the proof follows immediately from step (2), without any need for invoking recognisability.

We now prove the converse implication. Fix some infinite set X . We will show that all elements of X are finitary. Let X be a disjoint copy of X . For a finite subset $Y \subseteq X$, define

$$f_Y : X + X \rightarrow X + X$$

to be the function which maps all elements of Y to their corresponding black copies, and which is the identity on the remaining elements. Define \sim_Y to be the equivalence relation on $T(X + X)$ which is the kernel of the homomorphism

$$\mathbb{T}f_Y : T(X + X) \rightarrow T(X + X)$$

The mapping $Y \mapsto \sim_Y$ is monotone with respect to inclusion:

$$Y \subseteq Z \text{ implies } \sim_Y \subseteq \sim_Z .$$

In the conclusion of the above implication, we view equivalence relations as sets of pairs; under this view smaller equivalence relations refine bigger ones. Viewing equivalences relations as sets of pairs, define \sim to be the union of the equivalences \sim_Y , ranging over all finite subsets $Y \subseteq X$. This is an equivalence relation; transitivity follows from the monotonicity property above.

Consider the two natural injections

$$X + X \xleftarrow{f} X \xrightarrow{f} X + X,$$

such that f maps all elements to their black copies, and f maps all elements to their red copies. By assumption on the monad, there is a syntactic homomorphism for the quotient function of \sim , call it h . For every $x \in X$, the units of x and x are equivalent under the congruence $\sim_{\{x\}}$. Since the latter is a congruence that recognises \sim , it follows that the units of x and x must have the same image under h . By arbitrary choice of x , we see that

$$h \circ \mathbb{T}f = h \circ \mathbb{T}f. \quad (4.12)$$

We now prove that every $t \in TX$ is finitary. Fix t . Since h recognises \sim , the equality (4.12) implies that

$$(\mathbb{T}f)(t) \sim (\mathbb{T}f)(t).$$

By definition of \sim , there must be some finite Y , which depends on t , such that

$$(\mathbb{T}f)(t) \sim_Y (\mathbb{T}f)(t) \quad (4.13)$$

Choose an element of Y and let

$$g : X + X \rightarrow X$$

be the function which is the identity on X and sends all red elements to the

chosen element of Y . We have

$$\begin{aligned}
 t &= \text{(because } g \circ f \text{ is the identity on } X\text{)} \\
 (\mathbb{T}(g \circ f))(t) &= \text{(because } f_Y \circ f = f\text{)} \\
 (\mathbb{T}(g \circ f_Y \circ f))(t) &= \text{((4.13))} \\
 (\mathbb{T}(g \circ f_Y \circ f))(t) &
 \end{aligned}$$

The image of the function $g \circ f_Y \circ f$ is contained in Y , and therefore we have established that t is finitary. \square

Exercises

Exercise 128. (1) Give an example of a monad which is not finitary, but where every language $L \subseteq T\Sigma$ with a finite alphabet Σ has a syntactic homomorphism.

Exercise 129. (2) Consider an algebra colouring $\lambda : A \rightarrow U$. Consider the family of congruences on A that recognise U , ordered by inclusion. Show that this family is a lattice, i.e. every two elements have a greatest lower bound and also a least upper bound.

Exercise 130. (1) Let \mathbb{T} be a monad in the category of sets. Show that if every algebra colouring with two colours has a syntactic homomorphism, then every algebra colouring with an arbitrary number of colours has a syntactic homomorphism.

Exercise 131. (2) Give an example of a monad \mathbb{T} which is not finitary, but such that all elements of $\mathbb{T}X$ are finitary for countable X .

Exercise 132. (2) Let S be a finite set of sort names, and consider the category

$$\text{Set}^S$$

of S -sorted sets with sort-preserving functions. Prove Theorem 4.11 for monads over this category.

Exercise 133. (2) Consider a category of sorted sets, as in the previous exercise, but with infinitely many sort names. Define a finite algebra to be one that is finite on every sort. Show that Theorem 4.11 fails.

Exercise 134. (2) Show that a monad in the category of sets is finitary if and only if it arises as a result of the construction described in Exercise 110.

Exercise 135. Recall the notion of *regular elements* from Exercise 112. Show that if t is regular, then it is finitary.

Exercise 136. (2) Assume that the regular elements, as considered in the previous exercise, are closed under free multiplication in the following sense: if $t \in \mathbb{T}X$ is a regular term operation, and $\eta : X \rightarrow TY$ is a valuation of its variables that uses only regular elements, then $t^{\mathbb{T}Y}(\eta)$ is a regular element. Under these assumptions, define a monad of regular elements.

4.4 The Eilenberg Variety Theorem

In Chapter 2, we proved several theorems of the kind

$$\text{class of languages} \quad \sim \quad \text{class of semigroups.}$$

For example, a language of finite words is definable in first-order logic if and only if it is recognised by an aperiodic semigroup. In this section, we prove a result, originally due to Eilenberg, which says that every class of languages with good closure properties will correspond to a class of algebras with good closure properties. Eilenberg proved the theorem for monoids and semigroups, but the proof, as presented below, is the same in the abstract setting of monads.

The classes with good closure properties will be called *varieties*, in analogy with the varieties that appear in Birkhoff's theorem from universal algebra. There will be two kinds of varieties: for algebras and for languages. We begin with the algebras.

Definition 4.17 (Algebra variety). Let \mathbb{T} be a monad in the category of sets. An *algebra variety* is a class \mathcal{A} of finite \mathbb{T} -algebras with the following closure properties:

- *Quotients.* If \mathcal{A} contains A , then it contains every quotient of A .
- *Sub-algebras.* If \mathcal{A} contains A , then it contains every sub-algebra of A .
- *Products.* If \mathcal{A} contains A and B , then it contains $A \times B$.

Example 27. Consider the monad of finite words, where algebras are monoids. Examples algebra varieties include: finite groups, finite aperiodic monoids, finite infix trivial monoids, or finite prefix trivial monoids. \square

Example 28. Here is a non-example. Consider the monad of nonempty finite words, where algebras are semigroups. The class of monoids (i.e. semigroups which have an identity element) is not an algebra variety, because it is not closed under sub-algebras. \square

Example 29. Consider a monad T . Define an *identity* to be a pair of terms $s, t \in TX$ over a common set of variables. An algebra A is said to satisfy the identity if

$$s^A(\eta) = t^A(\eta) \quad \text{for every } \eta \in A^X.$$

The class of finite algebras that satisfy a given identity (more generally, all identities in a given set of identities) is easily seen to be an algebra variety. For example, the algebra variety of commutative semigroups arises from the identity

$$xy = yx.$$

Unlike in Birkhoff's theorem, some algebra varieties do not arise this way. For example, varieties discussed in Example 27 do not arise from (even possibly infinite sets of) identities¹² \square

We now describe language varieties. In Eilenberg's original formulations, this is a class of regular languages that is closed under Boolean combinations, inverse images of homomorphisms, and inverse images of operations of the form

$$w \in \Sigma^+ \mapsto v_1 w v_2 \in \Sigma^+ \quad \text{for fixed } v_1, v_2 \in \Sigma^*.$$

In the more abstract setting of monads, the role of these operations will be played by unary polynomial operations, as described in the following definition.

In the following definition, by recognisable languages we mean recognisable subsets of free algebras.

Definition 4.18 (Language variety). Let T be a monad in the category of sets. A *language variety* is a class of recognisable languages with the following closure properties:

¹² If we use the more general form of *profinite identities*, then every variety can be axiomatised equationally, which was shown in:

[30] Reiterman, "The Birkhoff theorem for finite algebras", 1982, Theorem 3.1.

Another answer is given in:

[16] Eilenberg and Schützenberger, *On pseudovarieties*, 1975, Theorem 1,

where it is shown that every variety can be axiomatised as follows: one gives an infinite set of identities, but requires only that all but finitely many identities are satisfied.

- *Boolean combinations.* \mathcal{L} is closed under Boolean combinations, including complementation.
- *Inverses of homomorphisms.* If $h : \mathbb{T}\Sigma \rightarrow \mathbb{T}\Gamma$ is a homomorphism of free algebras, then \mathcal{L} is closed under inverse images of h .
- *Inverses of unary polynomials.* If $f : \mathbb{T}\Sigma \rightarrow \mathbb{T}\Sigma$ is a unary polynomial operation in a free algebra $\mathbb{T}\Sigma$, then \mathcal{L} is closed under inverse images of f .

Example 30. Consider the monad of finite words, where algebras are monoids. The languages definable in first-order logic are a language variety. Closure under Boolean combinations is immediate, because we are dealing with a logic. Closure under inverse images of homomorphism or unary polynomials can be proved using Ehrenfeucht-Fraïssé games: if f is either a homomorphism or a unary polynomial operation, then a strategy copying argument shows

$$\begin{array}{ccc} \text{Duplicator wins the} & \text{implies} & \text{Duplicator wins the} \\ k \text{ round game on } w \text{ and } w' & & k \text{ round game on } f(w) \text{ and } f(w'). \end{array}$$

This implies that first-order definable languages are closed under inverse images of homomorphisms and unary polynomial operations. The same is true for first-order logic on \circ -words. \square

Example 31. Consider again the monad of finite words, where algebras are monoids. The definite languages from Example 19 are not a variety, because it is not closed under inverse images of the homomorphisms. Indeed, the language

$$a\{a, b\}^* \subseteq \{a, b\}^*$$

is definite. If we take the inverse image under the homomorphism

$$h : \{a, b, c\}^* \rightarrow \{a, b\}^*,$$

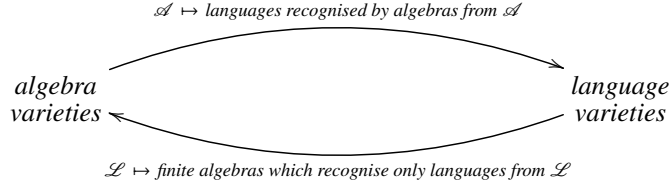
which erases the c letters, then we get the language

$$c^* a\{a, b, c\}^* \subseteq \{a, b, c\}^*,$$

which is not definite. The problem is with homomorphism that erase letters. If we would consider the same class of languages but in the monad of nonempty finite words, where algebras are semigroups, then we would get a variety. \square

The Eilenberg Pseudo-Variety Theorem says that the two notions of variety are two sides of the same coin.

Theorem 4.19 (Eilenberg Variety Theorem). *Let \mathbb{T} be a monad in the category of sets. The following maps are mutually inverse bijections*



Proof Let us write L for the left-to-right map, and A for the right-to-left map. We first show that each of these two maps take varieties to varieties, and then we show that the two maps are mutual inverses.

- (1) We first show that if \mathcal{A} is a class of finite algebras closed under products (which covers the special case of algebra varieties), then $L\mathcal{A}$ is a language variety. We begin with Boolean combinations. If L is recognised by an algebra $A \in \mathcal{A}$, then its complement is recognised by the same algebra. If furthermore K is recognised by $B \in \mathcal{A}$, then then $L \cup K$ and $L \cap K$ are both recognised by the product $A \times B$, which belongs to \mathcal{A} by closure under products. Consider now the inverse images. Let L be recognised by a homomorphism

$$h : \mathbb{T}\Sigma \rightarrow A \in \mathcal{A}.$$

If $g : \mathbb{T}\Sigma \rightarrow \mathbb{T}\Gamma$ is a homomorphism, then the inverse image $g^{-1}(L)$ is recognised by the composition $h \circ g$. Consider now a unary polynomial operation $f : \mathbb{T}\Sigma \rightarrow \mathbb{T}\Sigma$. As we have remarked in the proof of Lemma 4.12, congruences are compatible with unary polynomial operation, which means that

$$h(w_1) = h(w_2) \quad \text{implies} \quad h(f(w_1)) = h(f(w_2)).$$

This, in turn, means that h also recognises the inverse image $h^{-1}(L)$.

- (2) Next we show that if \mathcal{L} is a class of recognisable languages closed under unions and intersections (which covers the case of language varieties), then $A\mathcal{L}$ is an algebra variety¹³. Every language recognised by a sub-algebra of A is also recognised by A , and the same is true for quotients, and therefore

¹³ The first two steps of this proof establish that the maps L and A form what is known as a Galois connection, between

- classes of finite algebras closed under products; and
- classes of recognisable languages closed under unions and intersections.

In the terminology of Galois connections, the varieties are the closed sets, with respect to this Galois connection.

\mathcal{AL} is closed under sub-algebras and quotients of A . Consider now products. Suppose that a language L is recognised by a homomorphism

$$h : \mathbb{T}\Sigma \rightarrow A \times B \quad \text{with } A, B \in \mathcal{AL}.$$

For every $a \in A$, the inverse image

$$L_a = h^{-1}(\{a\} \times B)$$

is recognised by the homomorphism

$$h_A : \mathbb{T}\Sigma \rightarrow A,$$

which is the composition of h with the projection to A . Since the latter homomorphism has domain A , it follows that $L_a \in \mathcal{L}$. For similar reasons, if $b \in B$ then \mathcal{L} contains the language

$$L_b = h^{-1}(A \times \{b\}).$$

The intersection $L_a \cap L_b$ is the inverse image under h of the pair (a, b) . Every language recognised by h is a finite union of such languages; and therefore it belongs to \mathcal{L} by closure under unions and intersections.

- (3) We now show that the maps \mathbf{A} and \mathbf{L} are mutual inverses. We first show that every algebra variety \mathcal{A} satisfies

$$\mathcal{A} = \mathbf{AL}\mathcal{A}.$$

This is the same as showing that $A \in \mathcal{A}$ if and only if

- (*) every language recognised by A is recognised by some algebra in \mathcal{A} .

Clearly every algebra $A \in \mathcal{A}$ satisfies (*). We now prove the converse implication. Suppose that an algebra A satisfies (*). The multiplication operation

$$\mu : \mathbb{T}A \rightarrow A$$

in the algebra A is a homomorphism from the free algebra $\mathbb{T}A$ to A . By the assumption that A satisfies (*), every language recognised by this homomorphism is recognised by some algebra from \mathcal{A} . In particular, for every $a \in A$ the language $\mu^{-1}(a)$ is recognised by some homomorphism

$$h_a : \mathbb{T}A \rightarrow B_a \in \mathcal{A}.$$

Consider the product homomorphism

$$h : \mathbb{T}A \rightarrow \prod_{a \in A} B_a \quad t \mapsto (h_a(t))_{a \in A}$$

Define B to be the image of h . This is a sub-algebra of the a product of algebras from \mathcal{A} , and therefore it belongs to \mathcal{A} . From now on we view h as

surjective homomorphism onto its image B . This homomorphism recognises all languages $\mu^{-1}(a)$, and therefore μ factors through h , i.e. there is some function f which makes the following diagram commute:

$$\begin{array}{ccc} \mathsf{T}A & \xrightarrow{\mu} & A \\ & \searrow h & \uparrow f \\ & & B \end{array}$$

By Lemma 4.14, f is not just a function but also a homomorphism of algebras. This means that A is the image of B under a surjective homomorphism. In other words, A is a quotient of B , and therefore $A \in \mathcal{A}$.

(4) In the final step, we show that every language variety \mathcal{L} satisfies

$$\mathcal{L} = \mathsf{L}\mathcal{A}\mathcal{L}.$$

This is the same as showing that $L \in \mathcal{L}$ if and only if

(*) L is recognised by an algebra that only recognises languages from \mathcal{L} .

Clearly (*) implies $L \in \mathcal{L}$, so we focus on the converse implication. Suppose that $L \in \mathcal{L}$. Since L is recognisable, it has a syntactic homomorphism

$$h : \mathsf{T}\Sigma \rightarrow A$$

thanks to Theorem 4.11. To prove (*), we will show that all languages recognised by the syntactic algebra A belong to \mathcal{L} .

By the proof of Theorem 4.11, the kernel of the syntactic homomorphism h is the same as contextual equivalence for L . Therefore, by definition of contextual equivalence, we know that every $v, w \in \mathsf{T}\Sigma$ satisfy

$$h(v) = h(w) \quad \text{iff} \quad \underbrace{p(w) \in L \Leftrightarrow p(v) \in L}_{\text{for every unary polynomial operation } p : \mathsf{T}\Sigma \rightarrow \mathsf{T}\Sigma}.$$

Unary polynomial operations are compatible with congruences, which means that for every unary polynomial operation p there some function $p^A : A \rightarrow A$ which makes the following diagram commute

$$\begin{array}{ccc} \mathsf{T}\Sigma & \xrightarrow{h} & A \\ p \downarrow & & \downarrow p^A \\ \mathsf{T}\Sigma & \xrightarrow{h} & A \end{array}$$

Since p^A can be chosen in finitely many ways, and $p^{-1}(L)$ depends only on

p^A , it follows that there is a finite set \mathcal{P} of unary polynomial operations in $\mathbb{T}\Sigma$ such that

$$h(v) = h(w) \quad \text{iff} \quad \underbrace{p(v) \in L \Leftrightarrow p(w) \in L}_{\text{for every unary polynomial operation } p \in \mathcal{P}}. \quad (4.14)$$

We now show that \mathcal{L} contains every language recognised by A . Let then

$$g : \mathbb{T}\Gamma \rightarrow A$$

be some homomorphism. By surjectivity of the syntactic homomorphism and the universal property of the free algebra $\mathbb{T}\Gamma$, we can choose some homomorphism f which makes the following diagram commute

$$\begin{array}{ccc} \mathbb{T}\Gamma & & \\ \downarrow f & \searrow g & \\ \mathbb{T}\Sigma & \xrightarrow{h} & A \end{array}$$

We will show that for every $a \in A$, the language $g^{-1}(a)$ belongs to in \mathcal{L} . Since \mathcal{L} is closed under finite unions, it will follow that every language recognised by g belongs to \mathcal{L} , thus completing the proof of (*) for L .

Fix some $a \in A$. Since the syntactic homomorphism is surjective, one can choose some $v \in \mathbb{T}\Sigma$ such that $h(v) = a$. By (4.14), every $w \in \mathbb{T}\Gamma$ satisfies

$$g(w) = h(f(w)) = a = h(v) \quad \text{iff} \quad \underbrace{p(f(w)) \in L \Leftrightarrow p(v) \in L}_{\text{for every unary polynomial operation } p \in \mathcal{P}}.$$

If v is fixed, then the right side of the above is a finite Boolean combination of constraints of the form

$$p(f(w)) \in L \quad \text{which is the same as} \quad w \in \underbrace{f^{-1}(p^{-1}(L))}_{\text{a language in } \mathcal{L}}.$$

Summing up, we can express the constraint $g(w) = a$ as a finite Boolean combination of constraints $w \in K$, where K is some language in \mathcal{L} . Therefore, $g^{-1}(a)$ belongs to \mathcal{L} .

□

Exercises

Exercise 137. (1) Consider the monad of finite words. Show that a class of languages \mathcal{L} is a variety if and only if it is closed under Boolean combinations, inverse images under homomorphisms, and inverse images of unary polynomial operations of the form:

$$w \mapsto v_1 w v_2 \quad \text{for every choice of parameters } v_1, v_2 \in \Sigma^*.$$

Exercise 138. (2) Consider the monad of finite words. Show that there are uncountably many algebra varieties. In particular, for some algebra varieties, the membership problem $A \stackrel{?}{\in} \mathcal{A}$ is undecidable.

Exercise 139. (1) Consider the monad of \circ -words. Show that a class of languages \mathcal{L} is a variety if and only if it is closed under Boolean combinations, inverse images under homomorphisms, and inverse images of unary polynomial operations of the following forms:

$$w \mapsto v_1 w v_2 \quad \text{for every choice of parameters } v_1, v_2 \in \Sigma^\circ$$

$$w \mapsto w^\omega$$

$$w \mapsto w^{\omega*}$$

$$w \mapsto \text{shuffle of } \{w, v_1, \dots, v_n\} \quad \text{for every choice of parameters } v_1, \text{dots}, v_n \in \Sigma^\circ$$

Exercise 140. (2) Let S be a finite set, and consider the category

$$\text{Set}^S$$

of S -sorted sets with sort-preserving functions. State and prove the Eilenberg Pseudo-Variety Theorem for monads over this category.

Exercise 141. (2) Consider the monad from Example 23, which corresponds to weighted automata. We adapt to varieties to the weighted setting as follows. Define an algebra variety to be class of finite-dimensional algebras which is closed under sub-algebras, quotients and products. Define a language variety to be a class \mathcal{L} of linear maps $\mathbb{T}\Sigma \rightarrow \mathbb{Q}$, recognised by finite-dimensional algebras, which is closed under inverse images of homomorphisms and polynomials, and which is closed under combinations in the following sense: if \mathcal{L} contains

$$\{\lambda_i : \mathbb{T}\Sigma \rightarrow U_i\}_{i \in \{1,2\}},$$

and $f : \mathbb{Q}^2 \rightarrow \mathbb{Q}$ is a linear map, then \mathcal{L} contains also

$$w \mapsto f(\lambda_1(w), \lambda_2(w)).$$

Show that the Eilenberg Pseudo-Variety Theorem holds for varieties understood in this way.

Exercise 142. (2) Consider the weighted varieties from the previous example. What is the weighted analogue of star-free languages? Hint: consider the concatenation of two linear maps

$$\lambda_1, \lambda_2 : \mathbb{T}\Sigma \rightarrow U$$

to be the linear map which is defined as follows on Σ^*

$$(\lambda_1 \cdot \lambda_2)(a_1 \cdots a_n) = \sum_{i \in \{0, \dots, n\}} \lambda_1(a_1 \cdots a_i) \cdot \lambda_2(a_{i+1} \cdots a_n),$$

and which is extended to $\mathbb{T}\Sigma$ by linearity.

Exercise 143. (2) Consider the monad \mathbb{T} from Example 18, where $\mathbb{T}X$ describes terms over a fixed ranked set Σ with variables X . We view term $t \in \mathbb{T}\Gamma$ as a model, where the elements are the nodes of the corresponding tree, there is a binary ancestor relation $x \leq y$, and for every $\sigma \in \Sigma + \Gamma$ there is a unary relation $\sigma(x)$ which selects nodes with label σ . Show that the class of languages definable in first-order logic is not a variety. Hint: read the exercises in Chapter 5.

PART THREE

TREES AND GRAPHS

5

Forest algebra

In this chapter, we present a monad that models trees. The trees are going to be finite, node labelled, unranked (no restriction on the number children for a given node), and without a sibling order. Other kinds of trees can be modelled by other monads.

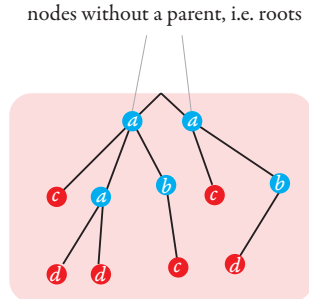
5.1 The forest monad

In fact, the monad will represent a slightly more general object, namely forests, i.e. multisets of trees. The algebras are going to be two-sorted, with the sort names being “forest” and “context”. For the rest of this chapter, define a *two-sorted set* to be a set together with a partition into elements of forest sort and elements of context sort. We use a convention where forest-sorted elements are written in **red**, context-sorted elements are written in **blue**, and black is used for elements whose sort is not known or which come from a set without sorts.

A *forest* over a two-sorted set Σ consists of:

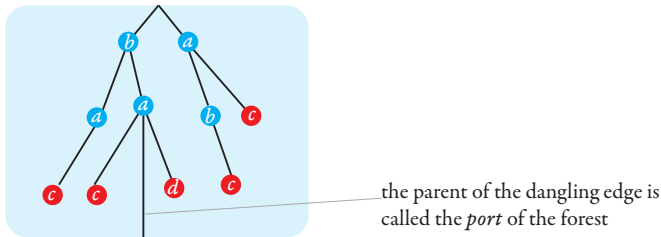
- a set of nodes;
- a parent function, which is a partial function from nodes to nodes;
- a labelling form nodes to Σ .

The parent function must be acyclic, and the labelling function must respect the following constraint: leaves (nodes that are not parents of any other node) have labels of sort “forest”, while non-leaves have labels of sort “context”. We assume that forests are nonempty, i.e. there is at least one node.



We use the usual tree terminology, such as root (a node without a parent), ancestor (transitive reflexive closure of the parent relation), child (opposite of the parent relation), descendant (opposite of ancestor) and sibling (children of a common parent).

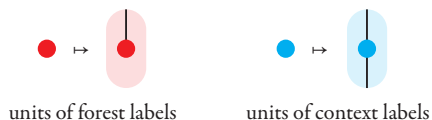
Apart from forests, the forest monad will also talk about *contexts*, which are forests with an extra dangling edge that is attached to a node with a context label, as in the following picture:



The forest monad. We now define a monad structure on forests and contexts. The underlying category is two-sorted sets, where objects are two-sorted sets and the morphisms are sort-preserving functions between two-sorted sets.

Definition 5.1. Define the *forest* monad as follows.

- The category is two-sorted sets, with the sorts called “forest” and “context”.
- For a two-sorted set Σ , the forest-sorted elements in $F\Sigma$ are forests over Σ , while the context-sorted elements are contexts over Σ . A sort-preserving function $f : \Sigma \rightarrow \Gamma$ is lifted to a sort-preserving function $Ff : F\Sigma \rightarrow F\Gamma$ by applying f to the label of every node and leaving the rest of the structure unchanged.
- The unit operation maps a label $a \in \Sigma$ to the unique forest/context that has one node with label a , as in the following pictures:



- Free multiplication is the operation of type $FF\Sigma \rightarrow F\Sigma$ that is illustrated in Figure 5.1. More formally, the free multiplication of $t \in FF\Sigma$ is defined as follows. The nodes are pairs (u, v) such that u is a node of t and v is a node

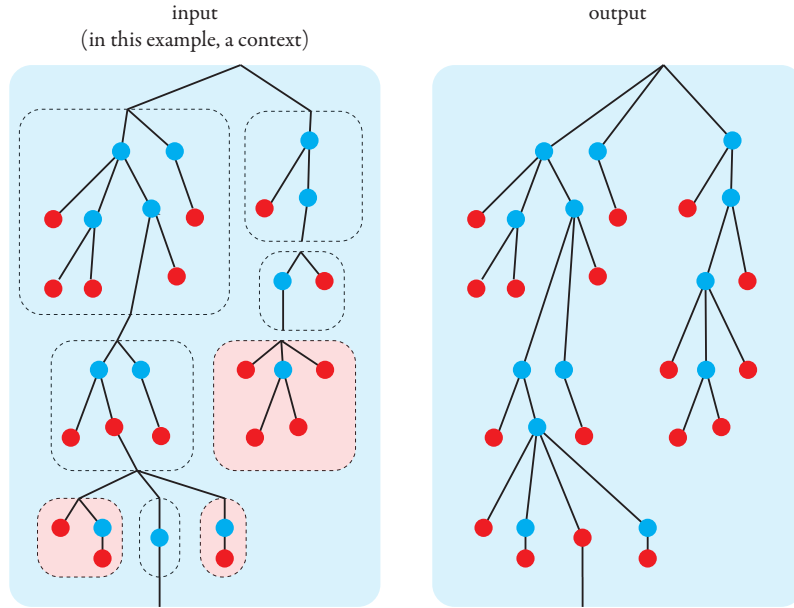


Figure 5.1 Free multiplication in the forest monad

in the label of v . The label is inherited from v , while the parent function is defined as follows (in the following $t_u \in \mathbb{F}\Sigma$ is the label of node u in t):

$$\begin{cases} (u, t_u\text{-parent of } v) & \text{if } v \text{ is not a root in } t_u; \\ (t\text{-parent of } u, \text{port of } t\text{-parent of } u) & \text{if } v \text{ is a root in } t_u \text{ and } u \text{ is not a root in } t; \\ \text{undefined} & \text{otherwise} \end{cases}$$

If t is a context, then the port in the free multiplication is defined to be the port of the context that labels the port of t .

We leave it as an exercise for the reader to check that the monad axioms are satisfied by the above definition. We use the name *forest algebra* for Eilenberg-Moore algebras over this monad.

5.2 Recognisable languages

The rest of this chapter is devoted to a study of the languages recognised by forest algebras. We care mainly about languages recognised by finite forest

algebras, which are forest algebras that have finitely many elements on both sorts. We begin with some examples.

Example 32. Let $\Gamma \subseteq \Sigma$ be two-sorted alphabets. We can view

$$F\Gamma \subseteq F\Sigma$$

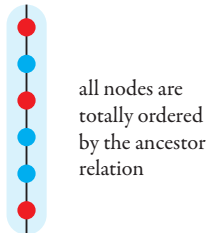
as a language, which only contains those forests and context over alphabet Σ where all labels are from Γ . This language is recognised by homomorphism

$$h : F\Sigma \rightarrow A$$

into a finite forest algebra A defined as follows. The underlying sorted set A is two copies of $\{0, 1\}$, one on the forest sort and one on the context sort. The multiplication operation in the algebra A maps $t \in FA$ to the maximal number from $\{0, 1\}$ that appears as a label in t , with the number cast into the appropriate sort. It is not hard to see that this is indeed a forest algebra, i.e. the operation μ is associative in the sense required of Eilenberg-Moore algebras. The homomorphism h maps an input to 1 (in the appropriate sort) if it belongs to the language, and to 0 otherwise. \square

Example 33. Let Σ be a sorted alphabet, let $n \in \{1, 2, \dots\}$. Consider the function h which inputs a forest or context in $F\Sigma$, and outputs the following information: (a) is it a forest or context; (b) what is the number of nodes modulo n . This function is compositional, in the sense of Section ?? . Lemma 4.7 is also true for monads in the category of sorted sets, and therefore h can be viewed as a homomorphism of forest algebras. This homomorphism recognises the language of forests/contexts where the number of nodes is divisible by n . \square

Example 34. Consider an alphabet Σ where all letters have context type. In this case, there are no forests over Σ , because there can be no leaves. For the same reason, every context over Σ looks like this:



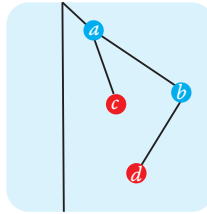
In other words, $F\Sigma$ is empty on the forest sort, and is isomorphic to the free semigroup Σ^+ on the context sort. Since the monad structure of the free semigroup agrees with the monad structure of the forest monad, it follows that a forest algebra with an empty forest sort is the same thing as a semigroup. \square

Exercises

Exercise 144. (2) Show that recognisable languages in the forest monad are closed under images of (not necessarily letter-to-letter) homomorphisms

$$h : F\Sigma \rightarrow F\Gamma.$$

Exercise 145. (2) Consider a variant of the forest monad, where we allow contexts where the port is a root, like in the following example:



Show that in this variant, recognisable languages are not closed under images of homomorphisms, but are closed under images of letter-to-letter homomorphisms.

5.2.1 A finite representation

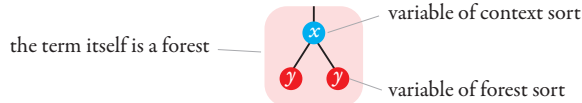
As usual with the monad approach, one needs to explain how the algebras can be finitely represented. Even if the underlying sorted set is finite, the multiplication operation

$$\mu : FA \rightarrow A$$

is in principle an infinite object. We show below a finite representation for the multiplication operation, in analogy to semigroups, where one only needs to define the multiplication operation for inputs of length two. When discussing this finite representation, we use as much as possible the abstract language of

monads; this will allow us to see analogies with other finite representations in this book.

A term basis. Like for any monad, a *term* in the forest monad is defined to be an element of FX for some set of variables X . Here is a picture of a term:



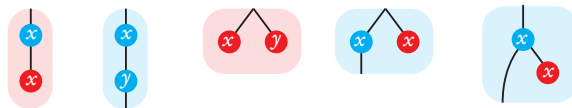
A difference with respect to terms for monads in the category of sets is that in the forest monad – which lives in the category of two-sorted sets – the variables are sorted, which means that there are forest variables, and context variables. Also, the term itself has a sort (call this the output sort). When interpreted in an algebra A , a term $t \in TX$ induces an operation

$$\eta \in A^X \mapsto \text{multiplication in } A \text{ applied to } (F\eta)(t).$$

The input to this operation, which is denoted by t^A and called a *term operation* in A , is a sort-preserving valuation of the variables, while the output is an element of the algebra whose sort is the output sort of the term. For example, the term in the picture above induces an operation, which inputs a context sorted x and a forest-sorted y , and outputs a forest sorted element. Note that term operations are not morphisms in the category of two-sorted sets, if only because there is no clear way of assigning a sort to the input valuation.

We distinguish the following terms in forest algebra.

Definition 5.2. Define the *basic forest algebra terms* to be the following terms:



(These happen to be all terms with exactly two nodes.) The *basic operations* in a forest algebra A are defined to be the term operations that are induced in A by these terms. We also use the following notation for the basic operations, listed in the order from the picture above (the colour of an operator is the colour of

| | | | |
|----------------------|-----|-----------------------|---|
| $x+(y+z)$ | $=$ | $(x+y)+z$ | (F1) forests with $+$ are a semigroup |
| $x+y$ | $=$ | $y+x$ | (F2) the forest semigroup is commutative |
| $x\cdot(yz)$ | $=$ | $(xy)\cdot z$ | (F3) contexts with \cdot are a semigroup |
| $x\cdot(y\cdot z)$ | $=$ | $(xy)\cdot z$ | (F4) \cdot is an action of contexts on forests |
| $x+(y+z)$ | $=$ | $(x+y)+z$ | (F5) $+$ is an action of forests on contexts |
| $x\oplus(y+z)$ | $=$ | $(x\oplus y)\oplus z$ | (F6) \oplus is an action of forests on contexts |
| $(x\cdot y)+z$ | $=$ | $(x+z)\cdot y$ | (F7) compatibility of the actions |
| $(x\oplus y)\cdot z$ | $=$ | $x\cdot(y+z)$ | (F8) compatibility of the actions |
| $(x\oplus y)\cdot z$ | $=$ | $x\cdot(z+y)$ | (F9) compatibility of the actions |

Figure 5.2 Axioms of forest algebra. The colour of the brackets indicates the sort of the bracket, and the colour of the equality sign indicates the sort of the compared elements.

the output sort):

| | | | | |
|--|--|---|---|---|
| $x\cdot x$ | $x\cdot y$ | $x+y$ | $x+x$ | $x\oplus x$ |
| inputs a context x and forest x and outputs a forest | inputs a context x and context y and outputs a context | inputs a forest x and forest y and outputs a forest | inputs a forest x and context x and outputs a context | inputs a forest x and context x and outputs a context |

Theorem 5.3. *The multiplication operation in a forest algebra is uniquely determined by the basic operations.*

Proof Every forest or context can be constructed from the units by applying the basic operations. □

The forest algebra in the above theorem does not need to be finite. If the forest algebra is finite, then it can be finitely represented by giving the multiplication tables for the basic operations.

We can also give simple list of axioms forest algebra, see Figure 5.2. These axioms are sound (they are satisfied by the basic operations in every forest algebra) and complete (if one gives five operations on a two sorted set A that satisfy the axioms, then these operations can be extended to a forest algebra operation $\mu : FA \rightarrow A$). Using this axiomatisation, we can effectively check if a finite representation of a forest algebra is correct, i.e. it comes from some forest algebra.

Exercises

Exercise 146. (2) Show that for every $t \in \mathbf{F}\Sigma$ there is a decomposition

$$t = f(t_1, \dots, t_n)$$

such that f is a term of size at most 4 (and therefore the number of arguments n is at most 4), and all arguments t_1, \dots, t_n have at most half the size (number of nodes) of t .

Exercise 147. (2) Fix some language $L \subseteq \Sigma$ that is recognised by a finite forest algebra. Suppose that we begin with some forest $t \in \mathbf{F}\Sigma$ and then we receive a stream of updates and queries. Each update changes a label of some node (the remaining forest structure is not changed). Each query asks if the current forest belongs to L . Show that one can compute in linear time a data structure (at the beginning, when the first forest t is given), such updates can be processed in logarithmic time and queries can be processed in constant time.

Exercise 148. (1) Prove completeness for the axioms (F1)–(F6).

5.2.2 Syntactic algebras

We now discuss syntactic algebras must necessarily exist in the forest monad. This is proved by a minor adaptation of the results from Section 4.3.

The notion of algebra colouring from Definition 4.8 makes sense also for the forest monad. In fact, this notion makes sense for any monad in any category, not just the category of sets: an algebra colouring is a morphism

$$\lambda : A \rightarrow C$$

from the underlying object in an algebra to some object in the category. Also the notion of a syntactic homomorphism from Definition 4.9 makes sense for monads in any category, if one interprets “surjective functions” as “epimorphisms”.

In the forest monad, where the underlying category is two-sorted sets, an algebra colouring is a sort-preserving function from the underlying two-sorted set in a forest algebra to some two-sorted set of colours. A subset $L \subseteq A$ can be seen as special case of algebra colouring which uses four colours

$$\underbrace{\{\text{yes}, \text{no}\}}_{\text{forest sort}} \cup \underbrace{\{\text{yes}, \text{no}\}}_{\text{context sort}}.$$

For the category of two-sorted sets, surjective functions are those which are surjective on both sorts.

The results on existence of syntactic homomorphisms from Section 4.3 can be easily adapted to the forest monad – more generally, to every monad in every category of sorted sets – as explained in the following theorem and its proof.

Theorem 5.5. *Let \mathbb{T} be a monad in a category of sorted sets (there could be more than two sorts, even infinitely many).*

- (1) *If \mathbb{T} is finitary, then every algebra colouring has a syntactic homomorphism;*
- (2) *If the monad is not necessarily finitary, but there are finitely many sort names, then every algebra colouring recognised by a finite algebra (finite on every sort) has a syntactic homomorphism.*

Proof For item (1) we use the same proof as in the left-to-right implication for Theorem 4.16, while for item (2) we use the same proof as in Theorem 4.11. In both cases the proofs use contextual equivalence. The only difference is that the definition of unary polynomial operations used for contextual equivalence from Definition 4.10 needs to be adapted to the sorted setting. In the sorted setting, a unary polynomial operation has an input sort (the sort of the variable) and an output sort (the sort of the polynomial). In item (2), the assumption on finitely many sort names is used¹ to conclude that the equivalence relation \approx in step (3) of Lemma 4.13 has finitely many equivalence classes; if there would be infinitely many sorts then there would be at least one equivalence class for each sort. Apart from this difference, the rest of the proof is the same. \square

In particular, since the forest monad is finitary, it follows that every language $L \subseteq \mathbb{F}\Sigma$ in the forest monad has a syntactic algebra. As discussed in the exercises, the syntactic algebra for a language recognised by a finite forest algebra can be computed.

Also, the Eilenberg variety theorem holds for the forest monad. In the statement, the unary polynomial operations are the sorted version that is described in the proof of Theorem 5.5, apart from this change the statement of the theorem and its proof are the same as in Section 4.4. More generally, the Eilenberg variety theorem works for every monad in every category of sorted sets, assuming that there are finitely many sorts. When generalising the proof of the Eilenberg Variety Theorem to multi-sorted algebras, we use the assumption on finitely many sorts in step (4) of the proof, to show that there are finitely many possible unary polynomial operations in a finite algebra.

¹ This assumption is indeed necessary, which can be proved using ideas from [6] Bojańczyk and Klin, “A non-regular language of infinite trees that is recognizable by a sort-wise finite algebra”, 2019

Exercises

Exercise 149. (2) Show that the syntactic algebra can be computed for a language $L \subseteq F\Sigma$ that is recognised by a finite forest algebra. The input to the algorithm is a homomorphism

$$h : F\Sigma \rightarrow A$$

into a finite forest algebra, together with an accepting set $F \subseteq A$. The forest algebra is represented using its basic operations, as in Theorem 5.3, and the homomorphism is represented by its values on the units.

5.2.3 Infinite trees

Define a monad F_∞ in the same way as the monad F , except that we allow the forests and contexts to be infinite. A node might have infinitely many children, and there might be infinite branches. This monad is no longer finitary.

We do not discuss this monad in more detail, apart from the following example, which shows that it is not clear what a “finite algebra” should be for this monad.

Example 35. Consider two-sorted alphabet $\Sigma = \{a, b\}$. Even though there are not forest-sorted letters, it is still possible to construct an infinite forest over this alphabet, because there is no need for leaves. Let $L \subseteq \Sigma^\omega$ be some prefix-independent language of ω -words, not necessarily regular. For example

$$L = \{b^{n_1} ab^{n_2} a \cdots : \text{the sequence } n_1, n_2, \dots \text{ contains infinitely many primes}\}.$$

Define a *branch* in a forest to be a set of nodes that is linearly ordered by the descendant relation, and which is maximal inclusion-wise for this property. An L -branch is a branch where the sequence of labels, starting from the root, belongs to L .

Define $L' \subseteq F_\infty \Sigma$ to be the set of infinite forests where every node has at least two children and every node belongs to some L -branch. We claim that L' is recognised by a finite algebra, with at most 36 elements, regardless of the choice of L (as long as it is prefix independent). Since there are uncountably many possible choices for L , it follows that there is no finite way of representing algebras in this monad that have at most 36 elements. In particular, “finite on every sort” is not a reasonable choice of “finite algebra” for this monad.

Define a function h from $F\Sigma$ to a finite two-sorted set as follows. For forests, the function h gives the answers to the following questions:

1. is the forest in K ?
2. are there are at least two roots?

For contexts, the function h gives the answers to the following questions:

1. is it possible to fill the port with some forest so that the result is in K ?
2. are there are at least two roots?
3. does the port have a sibling?
4. is the context equal to the unit of a ?
5. is the context equal to the unit of b ?

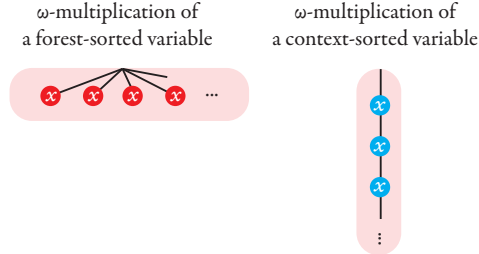
The red questions have at most 4 possible answers, and the blue questions have at most 32 possible answers hence the number 36. In fact, this number can easily be reduced; for example in case of a “no” answer to question 1, there is not need to store the answers for the remaining questions. We leave it as an exercise for the reader to check that the function h is compositional. It follows that the image of the function h , call it A , is a finite algebra for the monad F_∞ . \square

Exercises

Exercise 150. (1) Consider the monad F_∞ . We say that a forest/context in this monad is *thin* if it has countably many branches. Define $F_{\text{thin}}\Sigma \subseteq F_\infty\Sigma$ to be thin forests/contexts. Show that this is a monad.

Exercise 151. (2) Show that a forest/context is thin, in the sense of Exercise 150, if and only if one can assign countable ordinal numbers to its children so that if a node is labelled by ordinal number α , then all of its children are labelled by ordinal numbers $\leq \alpha$, and at most one child is labelled by α .

Exercise 152. (2) Consider the monad F_{thin} from Exercise 150. Show that a finite algebra over this monad is determined uniquely by its forest algebra operations (as in Theorem 5.3) plus the following two term operations:



5.3 Logics for forest algebra

In this section, we discuss languages in the forest monad that can be defined using logic, and the structural properties of the forest algebras that define them.

5.3.1 Monadic second-order logic

We begin with monadic second-order logic. The idea is as usual: to each forest/context we associate a model, and then we can use monadic second-order logic to describe properties of that model. There will be one twist: because siblings in a forest/context are not ordered, we will need to extend mso with modulo counting in order to make it expressively complete for all recognisable languages.

Definition 5.7. Define the *ordered model* of a forest or context as follows: the universe is the nodes, and it is equipped with the following relations:

$$\underbrace{x \leq y}_{\text{ancestor}} \quad \underbrace{a(x)}_{x \text{ has label } a \in \Sigma} \quad \underbrace{port(x)}_{x \text{ is the port}}$$

The arguments to the relations are x, y , while the letter a is a parameter. Each choice of $a \in \Sigma$ gives a different relation.

By using different logics on the ordered model, we get different classes of languages.

We begin with monadic second-order logic. A language $L \subseteq F\Sigma$ is called *mso definable* if it can be defined by a formula of monadic second-order logic using the ordered model. The logic mso is not enough to define all recognisable languages, because it cannot count the number of nodes modulo two (or three, etc.). The problem is that there is no order on the siblings, so if we get a forest

$$a + \cdots + a$$

that consists of n nodes that are both roots and leaves, then we cannot use the usual trick of selecting even-numbered nodes to count parity. (A more formal argument will be given below.) For these reasons, we extend mso with modulo counting. In this extension – called *counting* mso – for every set variable X and numbers $n \in \{2, 3, \dots\}$ and $\ell \in \{0, \dots, n - 1\}$ we can write a formula

$$|X| \equiv \ell \pmod{n}$$

which says that the size of the set X is congruent to ℓ modulo n . This extension is expressively complete for the recognisable languages, as shown in the following theorem.

Theorem 5.8. *A language $L \subseteq \text{F}\Sigma$ is recognised by a finite forest algebra if and only if it is definable in counting mso.*

Proof Both implications in the theorem are proved in a similar way as for finite words, so we only give a proof sketch.

From counting mso to a finite forest algebra. Same proof as for finite words: we remove the first-order variables (by coding them as singleton sets), and then we show by induction that for every formula of mso (possibly with free variables), its corresponding language is recognised by a finite forest algebra. In the induction steps we use products and powersets, both of which are finiteness preserving constructions for forest algebras.

From a finite forest algebra to counting mso. Suppose that $L \subseteq \text{F}\Sigma$ is recognised by a homomorphism

$$h : \text{F}\Sigma \rightarrow A$$

into a finite forest algebra. We will show that for every $a \in A$, the inverse image $h^{-1}(a)$ is definable in counting mso. By taking finite unions of such inverse images, we can get any language recognised by h , in particular L .

The idea is the same as for finite words: the defining formula inductively computes the value under h for every subtree in the input tree/context. The induction corresponds to a bottom-up pass through the input². Define the *type* of a node to be the image under h of its subtree, with subtrees defined in the same way as in Section ???. The following claim shows that the type of a node can be inferred from the types of its children using counting.

² This idea works for objects such as finite words or forest algebra, because they have a canonical way of parsing (left-to-right for words, or bottom-up for forest algebra), which can be defined in mso. For \circ -words, we do not know any simple parsing method, which is the reason why the implication from finite algebras to logic in Section 3.4 was hard. A similar phenomenon will appear for graphs, which will be discussed in the next chapter, which also do not have any simple definable canonical way of parsing.

Claim 5.9. For every $t \in \mathbf{F}\Sigma$ and every node x in t , the type of x depends only on the answers to the following questions:

- what is the label of node x ?
- are there exactly n children of x type a ?
- does n divide the number of children of x with type a ?

where a ranges over elements of the forest algebra A and $n \in \{0, \dots, |A|\}$.

Proof Let a_1, \dots, a_m be the types of the subtrees of the children of x . At most one of these types is a context, because there is at most one port. We only consider the case where all of the types are forests (and hence they will be written in red below); the case when one type is a context is treated similarly. If a is the label of node x , then the type of x is equal to

$$a \cdot (a_1 + \dots + a_m).$$

The label a is known, while the red part is multiplication in the forest semigroup of A , which is a commutative semigroup. In a commutative semigroup, the result of multiplication depends only on the number of times each argument is used. Furthermore, since the forest semigroup has size at most $|A|$, then the number of times an argument is used needs to be remembered only up to threshold $|A|$ and modulo some number that is at most $|A|$, see Exercise 11. \square

Consider some enumeration $A = \{a_1, \dots, a_n\}$ of the elements in the algebra. Using the above claim, we can write a formula

$$\varphi(X_1, \dots, X_n)$$

of counting mso which holds if and only if for every $i \in \{1, \dots, n\}$, the set X_i is exactly the set of nodes with type a_i . The formula simply checks that the types for each node are consistent with the types of its children, as described in the above claim. Finally, the image $h(t)$ of a forest/context can be computed in counting mso by guessing the sets X_1, \dots, X_n that satisfy the formula φ above, and then inferring $h(t)$ from the types of the root nodes (with the same argument as in the above claim). \square

The construction of an algebra in the above theorem is effective: given a sentence of mso, we can construct a recognising homomorphism

$$h : \mathbf{F}\Sigma \rightarrow A$$

into a finite forest algebra (and compute an accepting set $F \subseteq A$). The finite forest algebra is represented by its basic operations, as discussed in the previous section, and the homomorphism is represented by its images for the units.

As discussed above, counting is needed to define all recognisable languages. The exact role of counting is explained in the following theorem.

Theorem 5.10. *A language $L \subseteq F\Sigma$ is definable in mso (without counting) if and only if it is recognised by a forest algebra where the forest semigroup (forests equipped with $+$) is aperiodic.*

A corollary of this theorem is that modulo counting is needed to define the language “even number of nodes”, since this language cannot be defined by a forest algebra with an aperiodic forest semigroup.

Proof For the left-to-right implication, we use the same proof as in the left-to-right implication of Theorem 5.8. The only difference is that in Claim 5.9 we do not need modulo counting. This is because for every commutative aperiodic semigroup, the outcome of multiplication depends only on the number of times that each argument is used up to some finite threshold, without modulo counting.

Consider now the right-to-left implication, which says that if a language is definable in mso without counting, then it is recognised by a finite forest algebra with an aperiodic forest semigroup. Here, again, we use the same proof as in Theorem 5.8, where a recognising forest algebra is constructed by starting with some atomic forest algebras, and then applying products and the powerset construction. Since we do not need the relation

$$|x| \equiv \ell \pmod n$$

from the set model, all of the atomic forest algebras have forest semigroups that are aperiodic. Products clearly preserve aperiodicity of the forest semigroup, and the same is true powersets, as explained in the following lemma.

Lemma 5.11. *If S is a commutative³ aperiodic semigroup, then the same is true for its powerset semigroup PS .*

Proof In this proof, we use multiplicative notation for the semigroup operation. By aperiodicity of S , there is some $! \in \{1, 2, \dots\}$ such that every element of $b \in S$ satisfies $b^! = b^!b$. To establish aperiodicity of the powerset semigroup, we will show that every element $A \subseteq S$ of the power set semigroup satisfies

$$A^n = A^{n+1}$$

³ Commutativity is important in the proof of the lemma. For example, when proving the equivalence of mso and finite semigroups in the first part of the book, we started out with atomic semigroups that are aperiodic, and all other semigroups (including groups) could be obtained by applying products and powersets.

where n is the size of S times $! + 1$. We only show the inclusion $A^{n+1} \subseteq A^n$, the same proof can be used to establish the opposite inclusion. Let

$$a = a_1 \cdots a_{n+1} \in A^{n+1}.$$

By the pigeon-hole principle and choice of n , some $b \in A$ must appear at least $! + 1$ times in the sequence a_1, \dots, a_{n+1} . By commutativity and aperiodicity of S , one extra occurrence of b can be eliminated from the product, proving $A^{n+1} \subseteq A^n$. \square

\square

Corollary 5.12. *A language is definable in mso without counting if and only if its syntactic forest algebra is finite and has an aperiodic forest semigroup.*

Proof Aperiodicity of the forest semigroup is preserved when taking subalgebras and quotients (images under surjective homomorphisms). Since the syntactic forest algebra can be obtained from any recognising forest algebra by taking a subalgebra and then a quotient, the result follows from Theorem 5.10. \square

Since the syntactic forest algebra can be computed for recognisable languages, it follows that given a sentence of counting mso, we can decide if there is a sentence of mso which does not use counting and which is equivalent on forests and contexts.

5.3.2 First-order logic

For finite words, the king of algebraic characterisations was the Shützenberger-McNaughton-Papert-Kamp Theorem, which described the languages of finite words that can be defined in first-order logic (using the ordered model). Unfortunately, finding a generalisation of this theorem to forest algebra (or any other algebra modelling trees) remains an open problem. Therefore, our discussion of first-order logic in the forest monad is limited to some remarks and one example.

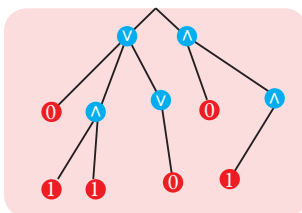
As discussed in Section 5.2.2, The Eilenberg Variety Theorem works also for the forest monad. One can show that, in the forest monad, the class of languages definable in first-order logic is a language variety, see Exercise 153. Therefore, from the Eilenberg Variety Theorem it follows that whether or not a language $L \subseteq F\Sigma$ is definable in first-order logic depends only on the syntactic algebra of the language. However, it is not known if the corresponding property

of syntactic algebras is decidable. Here is an example which shows that aperiodicity – which characterised the syntactic algebras for first-order definable languages of finite words – is not enough for forest algebra.

Example 36. Consider an alphabet

$$\Sigma = \underbrace{\{\vee, \wedge\}}_{\text{context sort}}, \underbrace{\{0, 1\}}_{\text{forest sort}}.$$

A forest over this alphabet is the same thing as a multiset of positive Boolean formulas, as in the following picture:



We define the *value* of a node in a forest over this alphabet to be the value of the Boolean formula in the subtree of the node. Consider the language

$$L = \{t \in \mathbf{F}\Sigma : t \text{ is a forest where all roots have value } 1\}$$

If we look at the syntactic forest algebra of this language, then both the forest semigroup and the context semigroup are aperiodic (in fact, they are idempotent). Nonetheless, the language is not definable in first-order logic, see Exercise 155. \square

Exercises

Exercise 153. (2) Prove that first-order logic, as discussed in Section 5.3.2, is a variety in the sense of the Eilenberg variety theorem.

Exercise 154. (2) Show that the language of forests where some leaf has even depth is not definable in first-order logic.

Exercise 155. (2) Show that the language from Example 36 is not definable in first-order logic.

Exercise 156. (2) Consider a two-sorted alphabet

$$\Sigma = \{\text{left, right, left, right}\}.$$

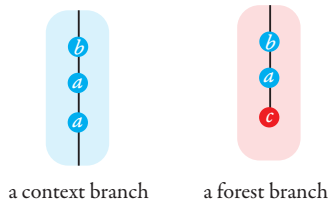
A *binary tree* over this alphabet is a tree where every node is either a leaf, or it has exactly two children, with labels “left” and “right” in the appropriate sort. There are no constraints on the root label. Define $L \subseteq F\Sigma$ to be the set of binary trees where all leaves are at even depth. Show that this language is first-order definable. Hint: show first that there is a first-order language which separates L from the set of binary trees where all leaves are at odd depth.

Exercise 157. (2) Define *anti-chain logic* to be the variant of mso where set quantification is restricted to anti-chains, i.e. sets of nodes that are pairwise incomparable with respect to the descendant relation. Show that anti-chain logic can define all recognisable languages that contain only binary trees, as defined in Example 156.

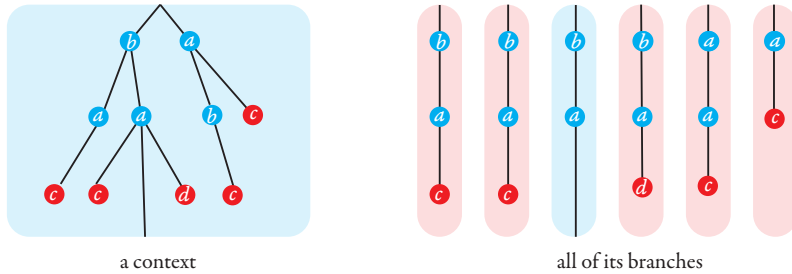
Exercise 158. (2) A unary node in a forest/context is a node with exactly one child. Show that anti-chain logic with modulo counting can define every recognisable language where every element has no unary nodes.

5.3.3 Branch languages

In this section, we discuss languages which are defined only by looking at branches in a forest/context. We say that a forest/context is a *branch* if all nodes are linearly ordered by the ancestor relation. Here is a picture:



A branch can be viewed as a word, consisting of the labels of the nodes in the branch, listed in root-to-leaf order. For a forest/context t , define a *branch of t* to be any branch that can be obtained from selecting some x which is either the port or a leaf, and restricting t to the ancestors of x . The branch is a context if x is the port, otherwise the branch is a forest. Here is a picture:



The following theorem gives a characterisation of languages that are determined by only looking at branches.

Theorem 5.13. *For every language $L \subseteq F\Sigma$, not necessarily recognisable, the following conditions are equivalent*

- (1) membership $t \in L$ depends only on the set of branches in t ;
- (2) the syntactic forest algebra of L satisfies the identities

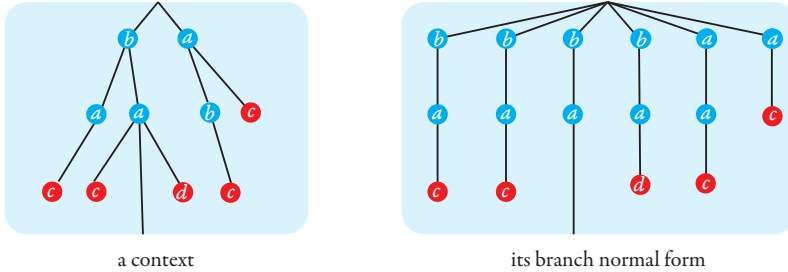
$$\underbrace{a \cdot (b + c) = (a \cdot b) + (a \cdot c)}_{\text{distributivity}} \quad \text{and} \quad \underbrace{b + b = b}_{\substack{\text{idempotence of} \\ \text{the forest semigroup}}} \quad \text{for every } a, b, c \in A.$$

If L is recognisable, then the above conditions are also equivalent to:

- (3) L is a finite Boolean combination of languages of the form “for some branch, the corresponding word is in $K \subseteq \Sigma^*$ ”, where K is regular:

Proof

- (1) \Rightarrow (2) Applying the identities in the free algebra $F\Sigma$ does not affect the set of branches.
- (2) \Rightarrow (1) For $t \in F\Sigma$, define its *branch normal form* to be the forest/context that is the union of all branches in t , as described in the following picture:



If the syntactic algebra satisfies the distributivity identity in the theorem, then a forest/context has the same image under the syntactic homomorphism as its branch normal form. Since the branch normal form is determined uniquely by the multiset of branches, it follows that the image under the syntactic homomorphism depends only on the multiset of branches⁴. Thanks to the idempotence identity, it is only the set of branches that matters, and therefore the language must be branch testable.

- (3) \Leftrightarrow (1) for recognisable languages. Clearly (3) implies (1). Consider now the converse implication. Let h be the syntactic homomorphism of a recognisable language. By condition (1) and the definition of a syntactic homomorphism, membership $t \in L$ depends only on the set

$$H(t) = \{h(s) : s \text{ is a branch in } t\}.$$

For every a in the syntactic algebra, define $K_a \subseteq \Sigma^*$ to be the words that correspond to branches which have value a under the syntactic homomorphism. This language is recognised by a finite semigroup (which is easily constructed from the forest algebra used by h), and therefore it is regular. Finally, $a \in H(t)$ if and only if for some branch the corresponding word is in K_a . Therefore, $H(t)$ can be computed using a finite Boolean combination of languages of the form K_a .

□

Condition (2) in the above theorem can be effectively checked given the syntactic algebra. Since the syntactic algebra can be computed for recognisable languages, it follows that one can decide if a recognisable language satisfies any of the conditions in the above theorem.

⁴ One could think that the distributivity identity alone (without the one for idempotence) characterises exactly the languages where membership depends only on the multiset of branches. This is not true, see Exercise 159.

Exercises

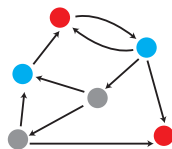
Exercise 159. (2) Give an example of a language $L \subseteq F\Sigma$ where membership depends only on the multiset of branches, but where the syntactic algebra violates the distributivity identity from Theorem 5.13.

Exercise 160. (2) Give an algorithm which decides if a recognisable language $L \subseteq F\Sigma$ is of the form: “for some branch, the corresponding word is in $K \subseteq \Sigma^*$ ”, for some regular K .

5.3.4 Modal logic

In this section, we discuss the tree variants of some of the temporal logics that were discussed in Chapter 2. When working with trees and forests, we use the terminology of modal logic, described as follows.

Define a *Kripke model* to be a directed graph with vertices labelled by some alphabet Σ . Here is a picture of a Kripke model:



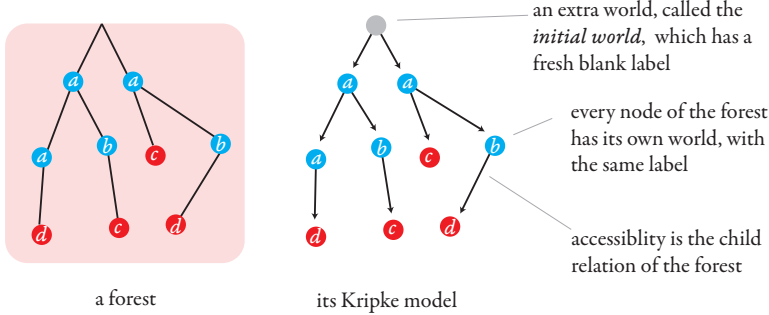
Vertices of the Kripke model are called *worlds*, and the edge relation is called *accessibility*. Accessibility does not need to be a transitive relation. In this section, we will only study Kripke models where accessibility is acyclic. To express properties of worlds in Kripke models we use modal logic, whose formulas are constructed as follows:

| | | | |
|--|---|--|--|
| \underbrace{a} | $\underbrace{\diamond\varphi}$ | $\underbrace{\square\varphi}$ | $\underbrace{\neg\varphi \quad \varphi \wedge \psi \quad \varphi \vee \psi}_{\text{Boolean combinations}}$ |
| the current world has label $a \in \Sigma$ | some accessible world satisfies φ | every accessible world satisfies φ | |

We use the following notation for the semantics of modal logic:

$$\underbrace{M}_{\text{Kripke model}}, \underbrace{v}_{\text{world of } M} \models \underbrace{\varphi}_{\text{formula of modal logic}}.$$

We will use modal logic to define properties of forests, by assigning a Kripke model to each forest⁵, as explained in the following picture:



Definition 5.15 (Forest languages definable in modal logic). We say that a formula of modal logic is true in a forest if it is true in the initial world of its Kripke model. A forest language is called *definable in modal logic* if there is a formula of modal logic that is true in exactly the forests from the language.

The following theorem characterises modal logic in terms of two identities. A corollary of the theorem is that one can decide if a language is definable in transitive modal logic; this is because it suffices to check if the identities hold in the syntactic algebra of a language.

Theorem 5.16. *Let $L \subseteq F\Sigma$ be a language that contains only forests. Then L is definable in modal logic if and only if its syntactic forest algebra is finite and satisfies the identities*

$$a + a = a \quad c^! \cdot a = c^! \cdot b,$$

where $! \in \{1, 2, \dots\}$ is the idempotent exponent of the context semigroup.

Proof The rough idea is that the identities say that the membership in the language is invariant under bisimulation (the first identity) and depends only on nodes at constant depth (the second identity). These are exactly the properties that characterise modal logic. A more detailed proof is given below.

Define the *modal rank* of a formula to be the nesting depth of the modal operators \diamond and \square . Here are some examples:

$$\underbrace{a}_{\text{modal rank 0}} \quad \underbrace{(\diamond a) \wedge (\diamond b)}_{\text{modal rank 1}} \quad \underbrace{(\diamond(a \wedge \square b) \wedge (\diamond b))}_{\text{modal rank 2}}.$$

⁵ One could also assign a Kripke model to a context, by doing the same construction, except with a special marker for the port node. We choose not to do this, without any deeper reasons, and therefore in what follows we only discuss languages that contain only forests.

When the alphabet is finite and fixed, then there are finitely many formulas of given modal rank, up to logical equivalence. To prove the theorem, we use a slightly more refined result, in the following claim, which characterises the expressive power of modal logic of given modal rank.

Claim 5.17. *A forest language can be defined by a formula of modal rank $n \in \{0, 1, \dots\}$ if and only if its syntactic algebra satisfies the identities*

$$a + a = a \quad c^n \cdot a = c^n \cdot b$$

Proof We say that a Kripke model is tree-shaped if the accessibility relation gives a finite tree, with edges directed away from the root (this is the case for the Kripke models that we assign to forests). We say that two tree-shaped Kripke models are *bisimilar* if one can be transformed into the other by applying the identity $a + a = a$, i.e. duplicating or de-duplicating identical sibling subtrees⁶. For $n \in \{0, 1, \dots\}$, we say that two tree-shaped Kripke models are *n-bisimilar* if, after removing all worlds that are separated by more than n edges from the root, they are bisimilar. By induction on n one shows that every equivalence class of *n-bisimilarity* can be defined by a formula of modal logic with modal rank n ; and conversely formulas of modal rank n are invariant under *n-bisimilarity*. The identities in the statement of the claim say that the forest language is invariant under *n-bisimilarity*, and hence the claim follows. \square

The theorem follows immediately from the above claim. Indeed, if the identities in the theorem are satisfied, then the language can be defined by a formula of modal logic with modal rank n . Conversely, if the language is defined by a formula of nesting depth n , then membership in the language is not affected by nodes which are more than n edges away from the root, and therefore the syntactic algebra must satisfy

$$\begin{aligned} c^n a &= && \text{(because } c^n \text{ is idempotent)} \\ c^{2n} a &= && \text{(because } a \text{ is more than } n \text{ edges away from the root)} \\ c^{2n} b &= && \text{(because } c^n \text{ is idempotent)} \\ c^n b & && \end{aligned}$$

\square

Transitive modal logic. A formula of modal logic can only talk about nodes that are at some constant distance from the root. We now discuss a variant of modal logic which can talk about arbitrarily deep nodes.

⁶ For tree-shaped Kripke models this notion coincides with the usual notion of bisimulation for general Kripke models.

For a forest, define its *transitive Kripke model* of a forest in the same way as the Kripke model, except that the accessibility relation now models the transitive closure of the child relation. In other words, accessibility now represents the proper descendant relation.

Definition 5.18 (Forest languages definable in transitive modal logic). A language that contains only forests is called *definable in transitive modal logic*⁷ if there is a formula of modal logic that is true in (the initial world) of exactly the forests from the language.

The following theorem characterises transitive modal logic in terms of two identities. A corollary of the theorem is that one can decide if a language is definable in transitive modal logic.

Theorem 5.19. *Let $L \subseteq \mathbf{F}\Sigma$ be a language that contains only forests. Then L is definable in transitive modal logic if and only if it is recognised by a finite forest algebra which satisfies the identities*

$$a + a = a \quad c \cdot a = (c \cdot a) + a.$$

Proof It is easy to see that the identities must be true in the syntactic algebra of every language definable in transitive modal logic. The first identity says that the language must be invariant under bisimulation, which is clearly true for transitive modal logic. For the second identity, we observe that going from $c \cdot a$ to $(c \cdot a) + a$ does not affect the transitive Kripke model, up to bisimulation.

The rest of this proof is devoted to the right-to-left implication. Let

$$h : \mathbf{F}\Sigma \rightarrow A$$

be a homomorphism into an algebra A that satisfies the identities. By induction on the size of A , we will show that for every forest-sorted $a \in A$, the inverse image $h^{-1}(a)$ is definable in transitive temporal logic (we say that such a is definable in the rest of the proof). This immediately yields the right-to-left implication. In the rest of the proof, we define the *type* of an element of $\mathbf{F}\Sigma$ to be its image under h .

In the proof, we use a reachability ordering on the algebra A which is defined as follows. We say that $a \in A$ is *reachable* from $b \in A$, denoted by $a \leq b$, if there is some $t \in \mathbf{F}A$ which uses b at least once, and which gives a under the multiplication operation of A . (Reachability can be seen as the forest algebra variant of the infix relation for semigroups.) Reachability is easily seen to be a pre-order, i.e. it is transitive and reflexive. We draw the

⁷ In the terminology of temporal logic, this logic is also called EF, which refers to the “exists finally” operator of (branching time) temporal logic.

reachability ordering in red when comparing forest-sorted elements. Thanks to the identities in the theorem, reachability is anti-symmetric when restricted to the forest sort, as explained in the following claim.

Claim 5.20. *If forest-sorted $a, b \in A$ are reachable from each other, then $a = b$.*

Proof For forest-sorted $a, b \in A$, reachability $a \geq b$ is equivalent to

$$\underbrace{a = c \cdot b}_{\text{for some context-sorted } c} \quad \text{or} \quad \underbrace{a = c + b}_{\text{for some forest-sorted } c} \quad \text{or} \quad a = b.$$

In the presence of the identities from the assumption of the theorem, all three conditions above imply $a = a + b$. For the same reason, $a \leq b$ implies $b = a + b$, and therefore $a = b$. \square

In every finite forest algebra there is a maximal forest-sorted element with respect to reachability, because every two forests can be combined using $+$, into a forest that is bigger than both of them. By the above, the maximal element is unique. Fix the maximal element a for the rest of the proof.

The following claim uses the induction assumption on algebra size to give a sufficient condition for definability.

Claim 5.21. *Assume that $c \in A$ is forest-sorted, and there is some non-maximal forest-sorted $b \in A$ from which c is not reachable. Then c is definable.*

Proof Let I be the set of elements in A that are reachable from b . This is an ideal, which means that if $t \in \text{FA}$ contains at least one letter from I , then its multiplication is in I . Furthermore, this ideal contains at least two forest-sorted elements, by assumption that b is non-maximal. Define \sim to be the equivalence relation on A which identifies two elements if they are equal, or both have the same sort and belong to I . Because I is an ideal, \sim is a congruence which means that the quotient function g which maps an element of A to its equivalence class is a homomorphism. Because b is non-maximal, the homomorphism makes the algebra smaller, and therefore the induction assumption on algebra size can be used to show that every language recognised by $g \circ h$ is definable in transitive modal logic. Because c is not reachable from b , it does not belong to the ideal I , and thus the equivalence class of c under \sim consists of c only. This means that h and $g \circ h$ have the same inverse images for c , and hence this inverse image is definable in transitive modal logic. \square

We will use the above claim to show that, with at most two exceptions, all forest-sorted elements of A are definable. Call an element b *sub-maximal* if it is not maximal and $c > b$ implies that c is maximal. If c is neither maximal nor

sub-maximal then it is definable by the above claim, because there is some sub-maximal b from which c is not reachable. For the same reason, if there are at least two sub-maximal elements, then all sub-maximal elements are definable. In particular, if there are at least two sub-maximal elements, then all elements are definable: the non-maximal elements are all definable, and the maximal element is definable as the complement of the remaining elements.

We are left with the case when there is exactly one sub-maximal element, call it b . We will show that the maximal element a is definable, and therefore b is definable (as the complement of the remaining elements). Define the *descendant forest* of a node x to be the forest that is obtained by keeping only the proper descendants of x . Since we do not allow empty forests, the descendant forest is defined only when x is not a leaf.

Claim 5.22. *A forest has type a if and only if it contains a node x , with label $\sigma \in \Sigma$, such that one of the following conditions hold:*

- (1) *the descendant forest of x has type $d < b$ and $h(\sigma) \cdot d = a$; or*
- (2) *the descendant forest of x has type $d \geq b$ and $h(\sigma) \cdot b = a$; or*
- (3) *the node x is a leaf and $h(\sigma) = a$.*

Proof To prove the bottom-up implication, we observe that each of the conditions (1,2,3) implies that the subtree of x has type a ; by maximality of a the entire forest must then also have type a . For conditions (1,3) the observation is immediate. For condition (2), there are two cases to consider: either the descendant forest has type a and the subtree of x has type a by maximality, or the descendant forest has type b and the subtree of x has type a by $h(\sigma) \cdot b = a$.

We now prove the top-down implication. We show that if every node in a forest violates all of the conditions (1,2,3), then the forest has type $\leq b$. If the forest has only one node, then we use condition (3). The second case is when the forest has at least two trees, i.e. it can be decomposed as $t = t_1 + t_2$. By induction assumption, both t_1 and t_2 have types $b_1, b_2 \leq b$. By the identity in the theorem, we get

$$b = b + b_1 + b_2,$$

which implies that $b_1 + b_2 \leq b$. The final case is when t is a tree, whose root x has label σ and descendant forest s . By induction assumption, s has a type $d \leq b$. If $d < b$ then we use condition (1) to infer that t has type b , otherwise we use condition (2). \square

To finish the proof of the theorem, it remains to show that the conditions in the above claim can be expressed using transitive modal logic. Condition (3) can easily be checked. In condition (1), the element d is definable because it is

neither maximal nor sub-maximal. Therefore, there is a formula of modal logic which is true in the world corresponding to a node x (in the descendant Kripke model) if and only if the descendant forest of x has type d . Therefore, there is a formula of modal logic which is true in the world corresponding to x if and only if it satisfies (1). For similar reasons, we can define condition (2), since the union of the languages for a and b is definable, by taking the complement of the remaining definable languages. \square

Exercises

Exercise 161. (2) Consider transitive modal logic for the monad F_∞ of infinite forests and contexts, as discussed in Section 5. Show that for infinite forests, the equations from Theorem 5.19 are sound (i.e. if a language is definable in transitive modal logic, then the syntactic algebra satisfies the equations) but not complete (i.e. the converse implication to soundness fails).

6

Hypergraphs of bounded treewidth

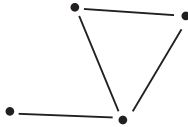
In this chapter, we study algebras for graphs and hypergraphs (which are like graphs, except that edges might connect more than two vertices). Although in principle the algebras can describe arbitrary graphs and hypergraphs, the more interesting results will assume bounded treewidth.

6.1 Graphs, logic, and treewidth

We begin by discussing graphs. Hypergraphs will appear later on, as they will provide the necessary structure used to define a monad.

Definition 6.1 (Graph). A *graph* consists of a set of vertices, together with a binary symmetric edge relation.

Here is a picture of a graph:



Unless otherwise stated, all graphs and hypergraphs in this chapter are finite. We use logic, mainly mso, to define properties of graphs.

Definition 6.2 (Graph languages definable in mso). Define the *incidence model* of a graph as follows. The universe is the disjoint union of the vertices and the edges, and there is a binary *incidence* relation, which is interpreted as

$$\{(v, e) : \text{vertex } v \text{ is incident with edge } e\}.$$

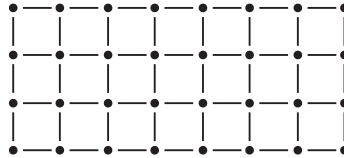
The two kinds of elements in the universe of the incidence model – vertices and edges – can be distinguished using first-order logic: an element of the universe is an edge if it is incident to some vertex, otherwise it is a vertex.

The incidence model is sometimes times called the mso_2 model. An alternative is the mso_1 model, where the universe is only the vertices, and there is a binary relation for the edges. For first-order logic, the two models are equivalent. This is because the edge relation of the mso_1 model can be defined in the incidence model by

$$E(v, w) \Leftrightarrow \exists e \text{ incident}(v, e) \wedge \text{incident}(w, e),$$

while first-order quantification over edges in the incidence model can be replaced by first-order quantification over pairs of vertices in the mso_1 model. For monadic second-order logic, the mso_2 model gives more expressive power than the mso_1 model, because of quantification over sets of edges. This is explained in the following example.

Example 37. A *clique* is a graph where every two vertices are connected by an edge. A *grid* is a graph that looks like this:



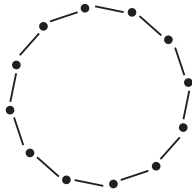
Both cliques and grids can be defined in mso (this is possible in both the incidence model and the mso_1 model). Consider now the set of graphs which are cliques of prime size. A clique has prime size if and only if it satisfies the following property that can be formalised in mso using the incidence model: one cannot remove edges so as to get a grid which has at least two rows and at least two columns. On the other hand, for graphs which are cliques, monadic second-order logic in the mso_1 model has the same expressive power as first-order logic. For cliques, every sentence of first-order logic will select finitely many or all but finitely many cliques. \square

In this chapter, we will not be studying first-order definable properties of graphs. Such properties are necessarily local¹, e.g. the existence of a cycle of length three:

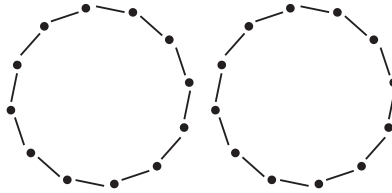
$$\exists u \exists v \exists w \quad E(u, v) \wedge E(v, w) \wedge E(w, u).$$

¹ The notion of locality is made precise by the Gaifman Theorem, see [21] Heinz-Dieter Ebbinghaus, *Finite Model Theory*, 2006, Theorem 2.5.1

A classical example of a property that is non-local, and therefore cannot be expressed in first-order logic, is graph connectivity. Using an Ehrenfeucht-Fraïssé argument, one can show that a sentence of first-order logic cannot distinguish between a large cycle and a disjoint union of two large cycles:



a connected graph



a disconnected graph

On the other hand, connectivity can be expressed in monadic second-order logic, already in the MSO_1 model, as witnessed by the following sentence

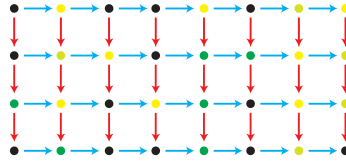
$$\underbrace{\exists X}_{\text{there is a set of vertices,}} \quad \underbrace{(\exists v \ v \in X) \wedge (\exists v \ v \notin X)}_{\text{which is neither empty nor full,}} \quad \wedge \quad \underbrace{(\forall v \ \forall w \ E(v, w) \wedge v \in X \Rightarrow w \in X)}_{\text{and which is closed under taking edges}}.$$

Example 38. [Hamiltonian cycles] Using MSO over the incidence model, we can say that a graph has a Hamiltonian cycle, i.e. a cycle that visits every vertex exactly once. A Hamiltonian cycle can be seen as a set of edges X such that: (a) every vertex in the graph has exactly one incoming and one outgoing edge from the set; (b) if only the edges from the set are used, then the graph is connected. The incidence model is important here. In the MSO_1 model, where quantification over sets of edges is not allowed, one cannot express the existence of a Hamiltonian path². \square

Already first-order logic is undecidable on graphs, in the following sense: it is undecidable whether or not a sentence of first-order logic is true in some graph. This undecidability is explained in the following example.

Example 39. Consider directed graphs with coloured vertices and edges. These extra features can be easily encoded, using first-order logic, in (undirected and unlabelled) graphs, see the exercises. A computation of a Turing machine can be visualised as a coloured grid, where: the rows represent configurations, and the colours of the vertices represent cell contents, as in the following picture:

² [13] Courcelle and Engelfriet, *Graph Structure and Monadic Second-Order Logic - A Language-Theoretic Approach*, 2012, Proposition 5.13



One can write a sentence of first-order logic which is true in exactly those grids that represent accepting computations of a given Turing machine. Therefore, the halting problem reduces to satisfiability for first-order logic over finite graphs. \square

Exercises

Exercise 162. (1) Show that for graphs without edges, mso has the same expressive power as first-order logic.

Exercise 163. (2) Unlike for the rest of this chapter, this exercise and the next one consider possibly infinite graphs. Consider two decision problems: (a) is a sentence of first-order logic true in some finite graph; (b) is a sentence of first-order logic true in some possibly infinite graph. Show that (a) is recursively enumerable (there is a Turing machine that accepts yes-instances in finite time, and does not halt on no-instances), while (b) is co-recursively enumerable (there is a Turing machine that does not halt on yes-instances, and rejects no-instances in finite time).

Exercise 164. (1) Show that grids, as described in Example 39, can be defined in first-order logic (with two extra unary predicates, which select the blue and red edges). We assume that the input graph has one connected component.

Exercise 165. (2) Show that the existence of a Hamiltonian cycle cannot be expressed in mso, using the mso_1 representation of graphs as a models.

Exercise 166. (2) Show that the existence of an Euler cycle (every edge is visited exactly once) cannot be expressed in mso, using the mso_1 representation of graphs as a models.

Exercise 167. (2) Show that for every $k \in \{1, 2, \dots\}$ the following property of

graphs is definable in mso using the incidence model: “the graph is connected, infinite, and has degree at most k ”. Show that “the graph is connected and infinite” is not definable.

Exercise 168. (2) Consider graphs which allow parallel edges (i.e. multiple edges connecting the same two vertices). The incidence model makes sense for such graphs as well. For $\ell \in \{1, 2, \dots\}$ define the ℓ -reduction of a graph to be the result the following operation: for each pair of vertices v, w we only keep the first ℓ edges that go from v to w . Show that for every mso sentence φ there is some ℓ such that φ is true in a graph (with parallel edges) if and only if it is true in its ℓ -reduction.

6.1.1 Treewidth

The undecidability problems described in Example 39 are avoided if we consider graphs that are similar to trees. In this chapter, the notion of similarity that we care about is treewidth, as defined below³.

Definition 6.3 (Tree decompositions and treewidth). A *tree decomposition* consists of:

- a graph, called the *underlying graph*;
- a set of *nodes*, equipped with a tree ordering (i.e. there is a least node called the root, and for every node x , the set of nodes $< x$ is totally ordered);
- for each node, an associated nonempty set of vertices called its *bag*.

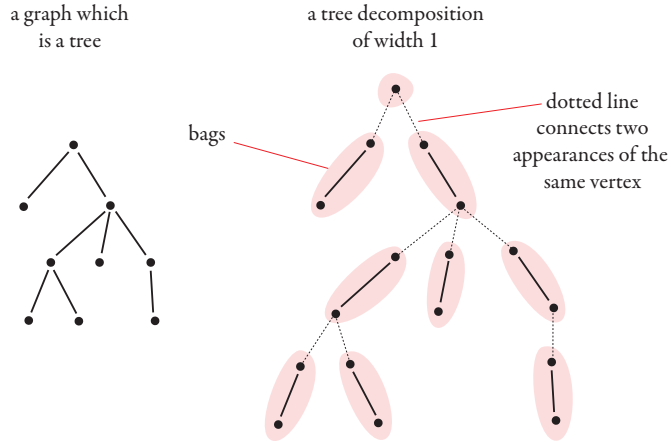
These should satisfy the following constraints:

- (1) for every edge in the underlying graph, there is some bag that contains both endpoints of the edge;
- (2) every vertex v of the underlying graph is *introduced* in exactly one node, which means there is exactly one node x such that v is in the bag of x and either x is the root or v is not in the bag of the parent of x .

The *width* of a decomposition is defined to be the maximal size of bags, minus one. The *treewidth* of a graph is the minimal width of a tree decomposition for the graph.

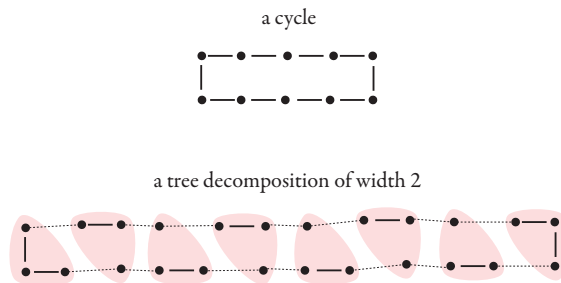
³ For an introduction to treewidth, including a brief history, see [15] Diestel, *Graph theory (electronic edition)*, 2006, Section 12.

If a graph is a tree (i.e. it is connected and has no cycle), then it has treewidth one. The nodes of the tree decompositions are edges of the graph, plus one extra root node. Here is a picture of a tree decomposition for a tree, where the underlying graph is a tree:



The bags in have size at most two, and therefore the treewidth is one (which explains why the width is defined to be the size of bags minus one). Forests (i.e. disjoint unions of trees) are the only graphs of treewidth one.

Cycles have treewidth 2, as illustrated in the following example:



The tree decomposition in the above picture is a *path decomposition*, i.e. every node in the tree decomposition has at most one child.

If a graph has $k + 1$ vertices, then it has treewidth at most k , since one can always use a trivial tree decomposition where all vertices of the graph are in the same bag. For cliques, the trivial tree decomposition is optimal, as explained in the following example.

Example 40. We show that for cliques, the trivial tree decomposition is op-

timal, i.e. in every tree decomposition, there must be some node whose bag contains all of the vertices in the clique. Consider a tree decomposition of a clique. The nodes of the tree decomposition where the vertices are introduced must be linearly ordered by the ancestor relation, since vertices introduced in nodes that are incomparable with respect to the ancestor relation cannot be connected by an edge. The maximal node in this linear order must have all vertices of the clique in its bag. \square

Another example of graphs with unbounded treewidth is grids, see the exercises. In fact, the Grid Theorem⁴, which is stated but not proved in the exercises, says that a class of graphs has unbounded treewidth if and only if it contains all grids as minors. We will show later in this chapter that for every $k \in \{1, 2, \dots\}$, the class of graphs of treewidth at most k has a decidable mso theory. In the exercises we also discuss a corollary of the Grid Theorem, which says that decidability of mso for bounded treewidth is optimal: if the mso theory of a class of graphs is decidable, then the class has bounded treewidth.

Exercises

Exercise 169. (2) We say that a graph G is a *minor* of a graph H if one can find a family of disjoint vertex sets

$$\{X_v \subseteq \text{vertices of } H\}_{v \in \text{vertices of } G}$$

such that every set of vertices U in G satisfies:

$$\underbrace{U \text{ is connected in } G}_{\substack{\text{a subset of vertices is connected if} \\ \text{the induced subgraph is connected}}} \quad \text{implies} \quad \left(\bigcup_{v \in U} X_v \right) \text{ is connected in } H.$$

(It is enough to check the implication for sets U with at most two vertices; and we assume that one vertex sets are connected). Show that if L is a property of graphs that is definable in mso using the incidence model (we use the incidence model for the remaining exercises), then the same is true for “some minor satisfies L ”.

Exercise 170. (1) The Grid Theorem says that if a class of graphs has unbounded treewidth, then for every $n \in \{1, 2, \dots\}$ there is some graph in the

⁴ For a recent paper about the Grid Theorem, see [11] Chuzhoy and Tan, “Towards Tight(er) Bounds for the Excluded Grid Theorem”, 2019

class which has an $n \times n$ grid as a minor. Using the Grid Theorem, prove that if a class of graphs has decidable mso theory, then it has bounded treewidth.

Exercise 171. (1) Building on the previous exercises, show that if a class of graphs has decidable mso theory, then it has bounded treewidth.

Exercise 172. (1) Show a class of graphs that has undecidable mso theory and bounded treewidth.

Exercise 173. (2) Show that the $n \times n$ grid has treewidth at least n .

Exercise 174. Show that for every language $L \subseteq \{a\}^*$ recognised in linear time by a (possibly nondeterministic) Turing machine, the language

$$\{G : G \text{ is an } n \times n \text{ grid such that } a^n \in L\}$$

is definable in mso.

6.2 The hypergraph monad

In this section, we describe a monad for hypergraphs⁵. Using algebras for this monad, we will study the connections between recognisability and definability in mso. The features of the hypergraph generalisation – labels, arities, ports – will be used to define the monad structure. Like any monad, the hypergraph monad will allow us to talk about algebras, homomorphisms, terms, recognisable languages, syntactic algebras, etc. The main result of this section is going to be Courcelle’s Theorem, which says that every graph property definable in mso is necessarily recognisable. Later, in Section 6.3, we also present a converse to Courcelle’s Theorem, which says that for bounded treewidth, recognisability implies definability in mso.

Definition 6.4. A *hypergraph* consists of:

- A set V of *vertices*.

⁵ This monad is based on the hyperedge replacement algebras of Courcelle. A discussion of these algebras can be found in

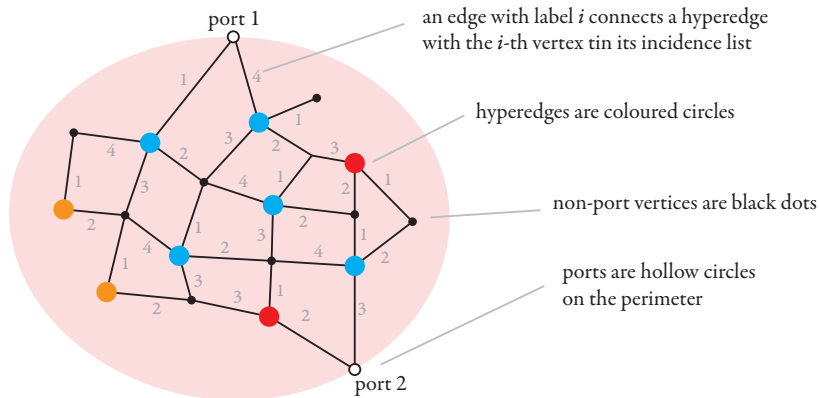
[13] Courcelle and Engelfriet, *Graph Structure and Monadic Second-Order Logic - A Language-Theoretic Approach*, 2012, Section 2.3.

The presentation of hyperedge replacement that uses monads is based on

[5] Bojańczyk, “Two Monads for Graphs”, 2018

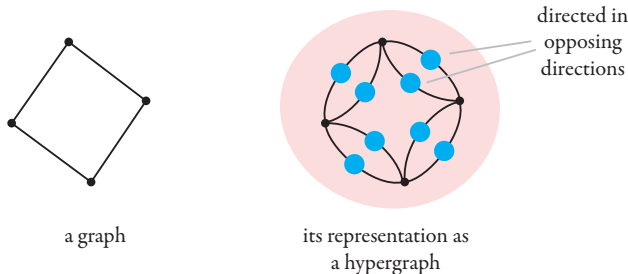
- A set E of *hyperedges*. Each hyperedge has an arity in $\{0, 1, \dots\}$.
- A set Σ of *labels*. Each label has an associated arity in $\{0, 1, \dots\}$.
- An arity $k \in \{0, 1, \dots\}$ and a non-repeating sequence of k distinguished vertices called *ports*;
- For each hyperedge e of arity n , an associated label in Σ of arity n , and a non-repeating sequence of *incident vertices* denoted by $e[1], \dots, e[n]$.

In this chapter, all hypergraphs are assumed to be finite, which means that there are finitely many vertices and hyperedges. We use the name *ports* for the vertices in the port sequence, and *non-port vertices* for the remaining vertices. We draw hypergraphs like this:



To avoid clutter in the pictures, we skip the numbers on the edges, if they are not important for the picture or implicit from the context.

A graph can be represented as a hypergraph. The representing hypergraph has zero ports, and the vertices are the same as in the graph. Each edge of the graph is represented by two binary hyperedges (with some fixed label), one in each direction. Here is a picture:



Directed graphs can be represented in the same way, but with the hyperedges not necessarily using both opposing directions.

The hypergraph monad. We now describe the monad structure of hypergraphs. The idea is that a hyperedge can be replaced by a hypergraph of matching arity. This replacement, which will be the free multiplication in the monad, is illustrated in the Figure 39.

Definition 6.5 (Hypergraph monad). The hypergraph monad, denoted by H , is defined as follows.

- The underlying category is the *category of ranked sets*

$$\text{Set}^{\{0,1,\dots\}}.$$

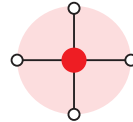
Objects in this category are *ranked sets*, i.e. sets where every element has an associated arity in $\{0, 1, \dots\}$. Morphisms are arity-preserving functions between ranked sets.

- For a ranked set Σ , the ranked set $H\Sigma$ consists of finite hypergraphs labelled by Σ , modulo isomorphism. This is indeed a ranked set; recall that the arity of a hypergraph is the number of ports.
- For a function $f : \Sigma \rightarrow \Gamma$, the function $Hf : H\Sigma \rightarrow H\Gamma$ applies f to the labels, without changing the rest of the hypergraph structure.
- The unit operation in the monad associated to every letter $a \in \Sigma$ of arity n is the hypergraph with ports $\{1, \dots, n\}$, no other vertices, and one hyperedge labelled by a with incidence list $1, \dots, n$. Here is a picture:

a letter

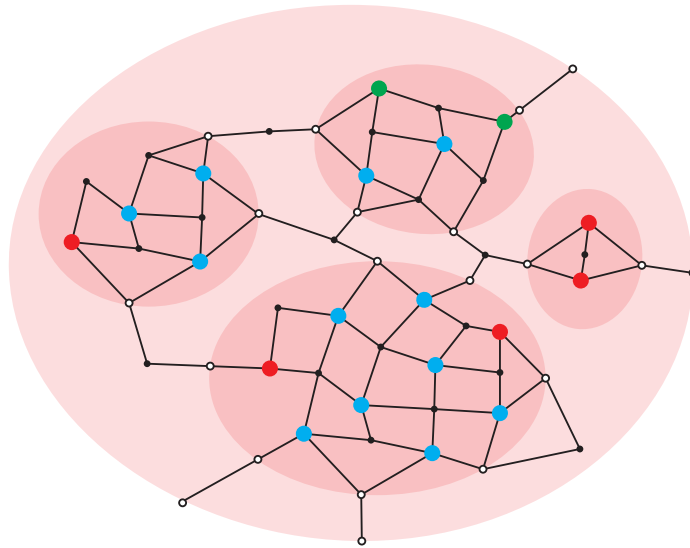


its unit hypergraph

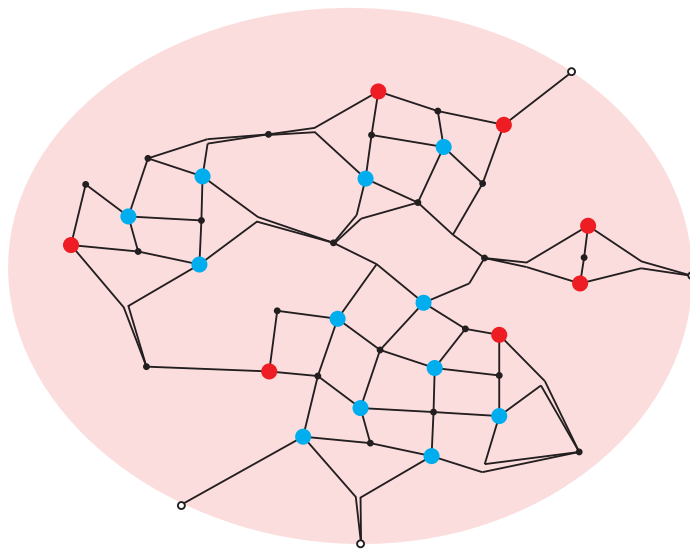


- Let $G \in HH\Sigma$ be a hypergraph labelled by hypergraphs. Its free multiplication is defined as follows. The vertices are vertices of G , plus pairs (e, v) such that e is a hyperedge of G and v is a vertex in the hypergraph G_e that is the label of the hyperedge e . The hyperedges are pairs (e, f) , where e is a hyperedge of G and f is a hyperedge in G_e . The arities and labels of hyperedges are inherited from the second coordinate, while the incidence lists are defined by

$$(e, f)[i] = \begin{cases} f[i] & \text{if } f[i] \text{ is a non-port vertex} \\ e[j] & \text{if } f[i] \text{ is the } j\text{-th port.} \end{cases}$$



a hypergraph labelled by hyperegraphs

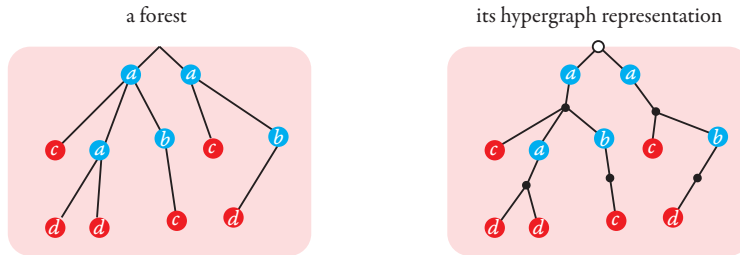


its free multiplication

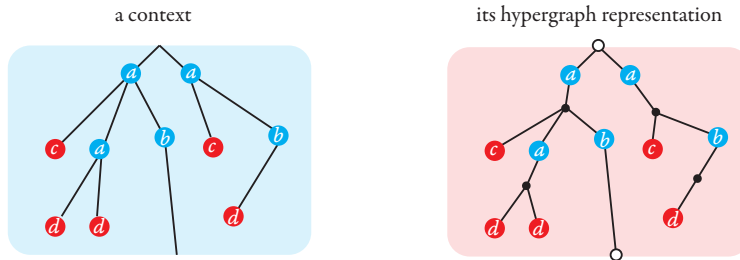
Figure 6.1 Free multiplication in the hypergraph monad.

We leave it as an exercise for the reader to check that the above definition satisfies the monad axioms. This completes the definition of the hypergraph monad.

Example 41. The hypergraph monad can be seen as a generalisation of the forest monad. A forest can be represented as a hypergraph of arity one, as explained in the following picture:



A context can be represented as a hypergraph of arity two:



This representation is consistent with the monad structures of the forest monad and the context monad. Therefore, we can think of the forest monad as being a sub-monad of the hypergraph monad (when we identify the forest sort with arity 1, and the context sort with arity 2). In particular, from every algebra of the hypergraph monad we can extract an algebra of the forest monad. \square

The rest of this section is devoted to discussing the algebraic notions that arise from this monad, such as terms and algebras.

Exercises

Exercise 175. (1) Show that \mathbf{H} satisfies the monad axioms.

Exercise 176. (1) Show that connected hypergraphs are also a monad.

6.2.1 Recognisable languages

We begin by discussing recognisable languages.

Example 42. Let M be a commutative monoid. Define an algebra in the hypergraph monad as follows. The underlying ranked set A has a copy of M on each arity. (More formally, an n -ary element of A is a pair (a, n) with $a \in M$.) The multiplication operation in the algebra maps a hypergraph $G \in \mathbf{HA}$ to the multiplication – in the monoid M – of all the labels of its hyperedges. (Because the monoid is commutative, the order of multiplication is not important.) The result of this multiplication is viewed as an element of the n -th copy of M , with n being the arity of G . It is not hard to see that this operation is associative, i.e. it satisfies the axioms of Eilenberg-Moore algebras.

The algebra constructed this way can be used to recognise some simple languages of hypergraphs. Consider the monoid

$$M = (\{0, \dots, n\}, \max).$$

If A is the algebra constructed for this monoid as above, then it can be used to recognise the hypergraph language

$$\{G \in \mathbf{H}\Sigma : \text{at least } n \text{ hyperedges have label in } \Gamma\} \quad \text{for ranked sets } \Gamma \subseteq \Sigma.$$

The homomorphism maps a hypergraph to the number of hyperedges with label in Γ , up to a threshold of n , with the number stored in the copy of the monoid that matches the arity of the input graph. Another application of this construction is recognising the language of hypergraphs with an even number of hyperedges; here the appropriate monoid is the two-element group. \square

The algebras in the above example are infinite, but finite on every arity. This is the best we can do in the hypergraph monad, because it is impossible for an algebra to have an underlying set that is finite altogether. This is because the multiplication operation $\mu : \mathbf{HA} \rightarrow A$ in an algebra is arity-preserving, and \mathbf{HA} is nonempty on every arity (as witnessed by hypergraphs without hyperedges). Therefore, the underlying set of an algebra must be nonempty on every arity. In the following definition, we assume that “finite algebras” are those which have finitely many elements for each arity.

Definition 6.6 (Recognisable language of hypergraphs). We say that a language $L \subseteq H\Sigma$ is *recognisable* if it is recognised by a homomorphism into an algebra which has finitely many elements on every arity.

This definition will turn out to be not restrictive enough, as far as general hypergraphs are concerned, see Example 45. In fact, no entirely satisfactory definition of “finite algebra” for general hypergraphs is known, and possibly does not exist. However, for hypergraphs of bounded treewidth, algebras that are finite on every sort will be behaved and equivalent to mso, as we will see in Section 6.3.

Example 43. Define a *path* in a hypergraph to be sequence of the form

$$v_0 \xrightarrow{e_1} v_1 \xrightarrow{e_2} \cdots \xrightarrow{e_{n-1}} v_{n-1} \xrightarrow{e_n} v_n,$$

where v_0, \dots, v_n are vertices or ports and e_1, \dots, e_n are hyperedges, such that each hyperedge e_i is incident with both v_{i-1} and v_i . Note that the notion of path does not depend on the order of the incidence lists for the hyperedges. The *source* of the path is v_0 , and v_n is its *target*. We say that the path *connects* v_0 with v_n . A hypergraph is called *connected* if every vertex or port can be connected to every other vertex or port via a path. Define h to be the function which maps a hypergraph to the following information: (a) its arity; (b) is there a pair of non-connected vertices such that at least one of them is not a port; and (c) which pairs of ports are connected. One can check that this is function is compositional, and therefore its image can be equipped with the structure of an algebra so that h is a homomorphism. The corresponding algebra is finite on every arity. Therefore, the language of connected hypergraphs is recognisable.

□

Example 44. In this example, we show that the language of k -colourable hypergraphs is recognisable. Define a k -*colouring* of a hypergraph to be a function from vertices to $\{1, \dots, k\}$ such that every hyperedge has an incidence list that uses each colour at most once. (In particular, all hyperedges must have arity at most k .) Define h to be the function which maps a hypergraph to the following information: (a) its arity; (b) which functions from the ports to $\{1, \dots, k\}$ can be extended to k -colourings. If the hypergraph has arity zero, then (b) is just one bit of information: is there a k -colouring or not. This function is compositional, and has finitely many values for each arity, and therefore the language of k -colourable hypergraphs is recognisable. □

The following example illustrates a problem with of our notion of recognisability, which is that it allows for too many algebras, at least as long as hypergraphs of unbounded treewidth are allowed.

Example 45. We say that a hypergraph is a clique if it has no ports, and for every two vertices are adjacent (i.e. connected by some hyperedge). Let P be any set of natural numbers, possibly undecidable. We will show that the language “cliques whose size is in P ” is recognisable. Let h be the function which maps a hypergraph to the following information: (a) its arity; (b) is there a pair of non-adjacent vertices such that at least one of them is not a port; (c) which ports are adjacent. If the arity is zero, then h also stores (d) is the number of vertices in P . This function is compositional, and has finitely many values for each arity, and therefore the language “cliques whose size is in P ” is recognisable. \square

As we will see later on, the problem from the above example will disappear once we restrict attention to hypergraphs of bounded treewidth.

Exercises

Exercise 177. (1) Show that every recognisable language in the hypergraph monad has a syntactic algebra.

6.2.2 Terms and tree decompositions.

In this section we discuss an alternative perspective on treewidth, which is defined using monad terminology. Tree decompositions and treewidth can be naturally extended to hypergraphs.

Definition 6.7 (Tree decompositions for hypergraphs). Tree decompositions are defined for hypergraphs in the same way as for graphs, with the following differences: (a) for every hyperedge there must be some bag which contains its entire incidence list (we say that such a bag *covers* the hyperedge); (b) every port of the hypergraph appears in the root bag.

As before, the width of a tree decomposition is the maximal bag size minus one, and the treewidth of a hypergraph is the minimal width of a tree decomposition. For hypergraphs which represent graphs (i.e. no ports, and every edge is represented by two binary hyperedges in opposing directions), the above definition coincides with Definition 6.3.

A bag in a tree decomposition can contain an unbounded number of hyperedges. This will not be a problem for our intended applications, since the properties of hypergraphs that we study will not depend in an important way on parallel hyperedges (i.e. hyperedges with the same incidence lists).

Tree decompositions as terms. The algebraic structure of the hypergraph monad can be used to give an alternative description of treewidth. Recall the notion of terms from Section 4.3.1: a term over variables X is any element of HX . As was the case for the forest monad, the terms in the hypergraph monad are sorted, which means that each variable used by the term has an arity, and the term itself has an arity. If A is an algebra, then a term $t \in HX$ induces a *term operation* t^A and defined by

$$\underbrace{\eta \in A^X}_{\substack{\text{an arity-preserving} \\ \text{valuation of the variables}}} \mapsto \underbrace{\text{multiplication in } A \text{ applied to } (H\eta)(t)}_{\substack{\text{an element of the algebra } A, \\ \text{whose arity is the same as the arity of } t}}.$$

As was the case for forest algebra, term operations are in general not arity-preserving.

Since a term is a hypergraph, it has some treewidth. The following lemma shows that hypergraphs of treewidth at most k are closed under applying (term operations induced by) terms of treewidth at most k . The algebra used in the lemma is the hypergraph free algebra.

Lemma 6.8. *Let $t \in HX$ be a term and let $\eta \in (H\Sigma)^X$ be a valuation of its variables. If the hypergraphs t and $\{\eta(x)\}_{x \in X}$ have treewidth at most k , then the same is true for the result of applying the term operation $t^{H\Sigma}$ to η .*

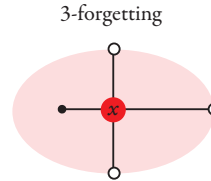
Proof Take a tree decomposition for the term t . For every hyperedge e which is labelled by a variable x , find a node whose bag contains the incidence list of the hyperedge, remove the hyperedge, and add a child to this node with a tree decomposition of $\eta(x)$. \square

A corollary of the above lemma is that there is a well-defined monad for hypergraphs of treewidth at most k . This monad, call it H_k , uses only hypergraphs with treewidth at most k , with all the monad structure inherited from H . The underlying category is ranked sets with arities at most $k + 1$; since hypergraphs with bigger arities will have treewidth at least $k + 1$. In particular, the alphabet Σ can have letters of arity at most $k + 1$.

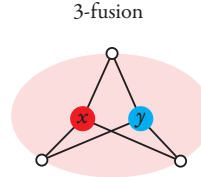
The treewidth terms. We now show a family of terms which can be used to generate all hypergraphs of given treewidth. Define the *treewidth terms* to be the terms⁶ from Figure 6.2. For an algebra in the hypergraph monad, define its *treewidth k operations* to be the term operations induced in the algebra by treewidth terms that have treewidth at most k . The following theorem

⁶ There is an inconsistency in our use of the words “introduce” and “forget”. When we say that a node in a tree decomposition introduces a vertex, we take a top-down perspective on tree decompositions. On the other hand, the name of the “forget” term in Figure 6.2 is based on a bottom-up perspective of the same phenomenon.

Forgetting. Let x be a variable of arity $k + 1$. The k -*forgetting term* is defined to be the hypergraph in $H\{x\}$ which has ports $\{1, \dots, k\}$, one non-port vertex v , and one hyperedge with label x and incidence list $(1, \dots, k, v)$.



Fusion. Let x, y be two variables of arity k . The k -*fusion term* is defined to be the hypergraph in $H\{x, y\}$ which has k ports, no vertices, and two hyperedges with labels x and y and incidence list $(1, \dots, k)$.



Rearrangement. Let $f : \{1, \dots, k\} \rightarrow \{1, \dots, \ell\}$ be an injective function. Let x be a variable of arity k . The f -*rearrangement term* is defined to be the hypergraph in $H\{x\}$ which has ℓ ports, no vertices, and one hyperedge with label x and incidence list $(f(1), \dots, f(k))$.

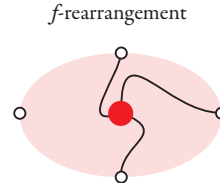


Figure 6.2 The treewidth terms. The parameters k, ℓ are from $\{0, 1, \dots\}$. In the pictures, the ports are ordered clockwise from top, and the same is true for the vertices incident to a hyperedge.

shows that the treewidth terms can be used to generate all hypergraphs of given treewidth.

Theorem 6.9. *Let $k \in \{1, 2, \dots\}$. A hypergraph has treewidth at most k if and only if it can be generated (in the hypergraph free algebra) from hypergraphs with no vertices and at most $k + 1$ ports, by applying treewidth k operations.*

Proof The right-to-left implication follows from Lemma 6.8.

Consider now the left-to-right implication. Every tree decomposition can easily be modified, without affecting its width, into a tree decomposition which satisfies: (*) the root bag contains only ports, and if a node has at least two children, then the bag of the node is equal to the bags of all of its children. To ensure condition (*), we can insert an extra node on every parent-child edge which has the same bag as the parent. By a simple induction on the size number of nodes, one shows that for every width k tree decomposition satisfying (*),

the underlying hypergraph can be generated using the treewidth k operations as in the statement of the lemma. \square

Using the above theorem, and the same argument as in Theorem 3.13, we get the following corollary.

Corollary 6.10. *Let $k \in \{1, 2, \dots\}$ and consider the monad H_k of hypergraphs with treewidth at most k . The multiplication operation in an algebra over this monad is uniquely determined by its treewidth k operations.*

6.2.3 Courcelle's Theorem

In this section, we prove Courcelle's Theorem, which says that all languages definable in mso are recognisable. To define properties of hypergraphs in mso, we use a hypergraph version of the incidence model, defined as follows.

Definition 6.11 (Incidence model). The incidence model of a hypergraph is defined as follows. The universe is vertices and hyperedges, and it is equipped with the following relations:

$$\underbrace{e[i] = v}_{\substack{v \text{ is the } i\text{-th} \\ \text{element of the} \\ \text{incidence list of } e}} \quad \underbrace{\text{port}_i(v)}_{\substack{v \text{ is the } i\text{-th port}}} \quad \underbrace{a(e)}_{\substack{\text{hyperedge } e \\ \text{has label } a \in \Sigma.}}$$

The arguments of the relations are e and v , while i and a are parameters. Each choice of parameters gives a different relation.

Recognisability holds also for slight strengthening of mso, called *counting* mso, which also allows the following form of modulo counting: for every $n \in \{2, 3, \dots\}$ there is a predicate

$$|X| \equiv 0 \pmod{n},$$

which inputs a set and says if the size of this set is divisible by n . The modulo counting predicate is a second-order predicate, since it inputs a subset of the universe, and not an element (or tuple of elements) in the universe.

Counting mso is more powerful than mso without counting, e.g. “the number of vertices is even” can be defined in counting mso but not in mso, see Exercise 162.

Theorem 6.12 (Courcelle's Theorem). *If a language $L \subseteq H\Sigma$ is definable in counting mso, over the incidence model, then it is recognisable, i.e. recognised by a homomorphism into an algebra that is finite on every arity.*

We use the same construction as in previous chapters. The main step of the proof, which deals with set quantification, is presented in the following lemma.

Lemma 6.13. *The recognisable languages images under functions of the form⁷*

$$\underbrace{Hf : H\Sigma \rightarrow H\Gamma \quad \text{for } f : \Sigma \rightarrow \Gamma}_{\text{such functions are called letter-to-letter homomorphisms.}}$$

Proof We use a powerset construction for algebras. Since we have already used powerset constructions before, we take this opportunity to discuss powerset constructions in more detail and generality, so that we can think about the kinds of monads that allow a powerset construction (hint: these are not all monads, e.g. the group monad does not have a powerset construction).

For a ranked set X , define its *powerset* to be the ranked set PX where elements of arity n are sets of elements from A that have arity n . For an arity-preserving function $f : X \rightarrow Y$ on ranked sets, define

$$Pf : PX \rightarrow PY$$

to be the arity-preserving function that maps a set to its image. In the language of category theory, P is the co-variant powerset functor (the contra-variant powerset functor uses inverse images instead of forward images). For a ranked set X , define *distribution on X* to be the function of type

$$HPX \rightarrow PHX$$

which inputs a hypergraph G , and outputs the set of hypergraphs that can be obtained from G by choosing for each edge an element of its label.

Claim 6.14. *Distribution is a natural transformation, which means that the following diagram commutes for every arity-preserving function $f : X \rightarrow Y$*

$$\begin{array}{ccc} HPX & \xrightarrow{HPf} & HPY \\ \text{distribute on } X \downarrow & & \downarrow \text{distribute on } Y \\ PHX & \xrightarrow{PHf} & PHY \end{array}$$

Proof The right-down path corresponds to the following procedure: for each hyperedge, choose an element of its label, and then apply f . The down-right path corresponds to the following procedure: for each hyperedge, take the image under f of its label, and then choose an element. The two procedures give the same result. This is true thanks to the following property of distribution for

⁷ In Exercise 180 we show that the assumption on letter-to-letter homomorphisms is important.

hypergraphs: if we apply distribution on X to some hypergraph, then every hypergraph in the resulting set will have the same vertices, ports and hyperedges as the original hypergraph. This property would not hold, for example, in the group monad; in fact the claim would be false in the group monad⁸. \square

We use the powerset and distribution to prove the lemma. Suppose that a language L is recognised by a homomorphism

$$h : H\Sigma \rightarrow A,$$

and consider a letter-to-letter homomorphism

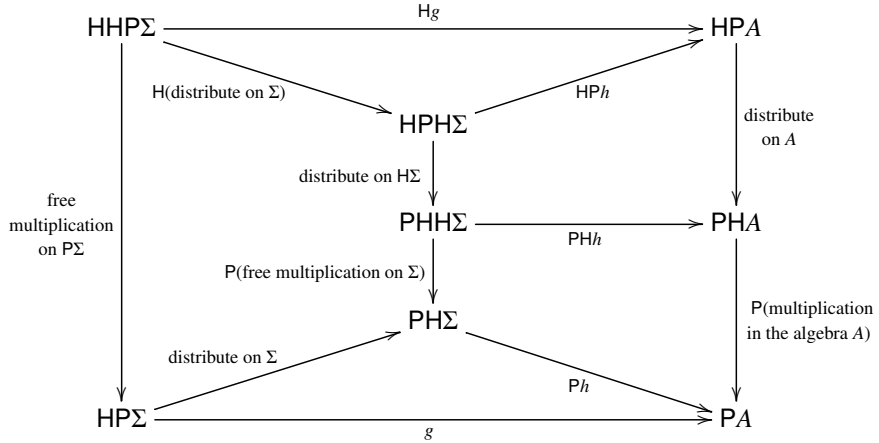
$$Hf : H\Sigma \rightarrow H\Gamma.$$

Define g to be the composition of the following two functions:

$$H\P\Sigma \xrightarrow{\text{distribute on } \Sigma} P H\Sigma \xrightarrow{P h} P A.$$

Claim 6.15. *The function g is compositional.*

Proof Consider the following diagram.



If we prove that the perimeter of the diagram commutes, then we will prove that g is compositional. The upper and lower triangular faces commute by definition of g . The upper rectangular face commutes by naturality of distribution

⁸ In fact, there is no powerset construction for algebras in the group monad. Nevertheless, by a proof that does not use powerset algebras, one can show that in the group monad the recognisable languages are closed under images of letter-to-letter homomorphisms.

from Claim 6.14, and the lower rectangular face commutes because h is a homomorphism. It remains to show that the five-sided face on the left commutes. This again, is proved via simple check⁹, similarly to Claim 6.14. \square

Like for any compositional function, the image of g can be equipped with a multiplication operation which turns g into a homomorphism. (This multiplication operation is the composition of the two arrows on the right side from the diagram in the proof of the claim.) We will use the homomorphism g to recognise the image of L under Hf . Consider the following diagram:

$$\begin{array}{ccc}
 H\Gamma & \xrightarrow{\text{H(inverse image under } f\text{)}} & HP\Sigma \\
 \text{inverse image under } Hf \downarrow & \text{distribute on } \Sigma \nearrow & \downarrow g \\
 PH\Sigma & \xrightarrow{P_h} & PA
 \end{array}$$

The top-left triangular face in the diagram commutes by definition of distribution, and the bottom-right triangular face in the diagram commutes by definition of g . A hypergraph $G \in H\Gamma$ belongs to the image of L under Hf if and only if applying the function on the down-right path in the diagram gives a set that intersects the image $h(L)$. Since the diagram commutes, it follows that the right-down path in the diagram recognises the image $(Hf)(L)$. The right-down path is a homomorphism, as a composition of two homomorphisms. Also, PA is finite on every arity, because finiteness on every arity is preserved by powersets. \square

The above lemma implies that recognisable languages are closed under quantification of sets of hyperedges (since a subset of the hyperedges can be seen as a colouring of hyperedges with two colours “yes” and “no”). This motivates the following logic that only allows quantification over sets of hyperedges.

Definition 6.16. Define *hyperedge* mso to be the following variant of mso. There is no first-order quantification, and set quantifiers range over sets of hyperedges. The logic allows the following relations on sets of hyperedges:

| | | | | |
|---|--|--|--|---|
| $\underbrace{X \subseteq Y}$ set inclusion | $\underbrace{X \subseteq a}$ every hyperedge in X has label $a \in \Sigma$ | $\underbrace{i \in X[j]}$ there exists a hyperedge $e \in X$ such that $e[j]$ is the i -th port | $\underbrace{X[i] \cap Y[j] \neq \emptyset}$ there exist hyperedges $e \in X$ and $f \in Y$ such that $e[i] = f[j]$ | $\underbrace{ X \equiv 0 \pmod n}$ the number of hyperedges in X is divisible by n |
|---|--|--|--|---|

In the above relations, the arguments are the sets X, Y . The labels $a \in \Sigma$ and

⁹ In the language of category theory, the five-sided face is the main axiom of a distributive law of a monad over a functor.

numbers $i, j, n \in \{1, 2, \dots\}$ are parameters. Each choice of parameters gives a different relation.

Lemma 6.17. *If a language $L \subseteq H\Sigma$ is definable in hyperedge mso, then it is recognisable.*

Proof Same proof as for the monads for words and forests. Consider a formula of hyperedge mso

$$\varphi(\underbrace{X_1, \dots, X_n}_{\substack{\text{the free variables represent} \\ \text{sets of hyperedges}}}),$$

where Σ is the ranked set of labels used by the underlying hypergraphs. Define the *language* of this formula to be the set hypergraphs over an extended alphabet that consists of 2^n disjoint copies of Σ . This language is defined in the same way as for words and forests: for each hyperedge, the bits from 2^n in its label determine which of the sets X_1, \dots, X_n contain the hyperedge. By induction on formula size, we prove that every formula has a recognisable language. The induction step is proved in the same way as for words and forests: for Boolean combinations we use homomorphisms into product algebras, while for the quantifiers we use the powerset construction from Lemma 6.13.

We are left with the induction base. For the formulas $X \subseteq Y$ and $X \subseteq a$, the corresponding hypergraph language is of the form “every hyperedge has a label in $\Gamma \subseteq \Sigma$ ”. Such languages were shown to be recognisable in Example 42. In the same example, we showed how to count hyperedges modulo some number, thus showing recognisability of the modulo counting relation. Consider now the language which corresponds to the relation

$$i \in X[j].$$

Let h be the function h which maps a hypergraph to the following information: (a) its arity; (b) which ports belong to $X[j]$. This function is easily seen to be compositional, and it has finite image on every arity, and therefore it witnesses recognisability of the language corresponding to $i \in X[j]$. A similar argument works for

$$X[i] \cap Y[j] \neq \emptyset.$$

□

Using the above lemma, we finish the proof of Courcelle’s Theorem. There is a minor issue that needs to be resolved. Because the hyperedge mso can only

quantify over sets of hyperedges, it cannot express properties of isolated vertices, i.e. vertices which are not adjacent to any hyperedge. To finish the proof of Courcelle's Theorem, we need to take into account the isolated vertices.

For a hypergraph G , define $\alpha(G)$ to be the hypergraph of same arity obtained from G by removing all isolated vertices that are not ports, and define $\beta(G)$ to be the hypergraph obtained by removing all non-isolated vertices that are not ports and all hyperedges. In particular, G is the fusion of $\alpha(G)$ and $\beta(G)$. By induction on formula size, one can show that every sentence of counting mso is equivalent to Boolean combination of sentences of counting mso, each of which talks about only $\alpha(G)$ or $\beta(G)$. Therefore, to prove Courcelle's Theorem, it is enough to show that for every sentence φ of counting mso, both of the following languages are recognisable:

$$\underbrace{\{G \in \text{H}\Sigma : \alpha(G) \models \varphi\}}_{L_\alpha} \quad \underbrace{\{G \in \text{H}\Sigma : \beta(G) \models \varphi\}}_{L_\beta}$$

For the language L_α we use Lemma 6.17. Let n be the maximal arity of letters in the finite alphabet. Every set X of non-isolated vertices can be represented by n sets of hyperedges as

$$X = X_1[1] \cup \dots \cup X_n[n],$$

where X_i is the set of hyperedges whose i -th incident vertex is in X . Using this representation, we can quantify over sets of non-isolated vertices by using quantification over sets of hyperedges. It follows that the language L_α can be defined in hyperedge mso, and is therefore recognisable by Lemma 6.17.

For the language L_β , we observe that it can be defined by checking an ultimately periodic property of the numbers of ports and isolated vertices. For a hypergraph $G \in \text{H}\Sigma$, define $\gamma(G)$ to be the word $a^n b^m$ where n is the number of ports and m is the number of isolated vertices. This word is roughly the same thing as $\beta(G)$, except that it has an order. Therefore, for every sentence φ of counting mso on graphs, one can easily find a sentence ψ of counting mso over finite words which makes the following diagram commute:

$$\begin{array}{ccc} \text{H}\Sigma & \xrightarrow{\gamma} & a^* b^* \\ \beta \downarrow & & \downarrow \psi \\ \text{H}\Sigma & \xrightarrow{\varphi} & \{\text{yes, no}\} \end{array}$$

For finite words, counting mso has the same expressive power as mso, because modulo counting can be expressed using the order of the word. Since mso on finite words can only define regular languages, it follows that ψ defines a regular language contained in $a^* b^*$. Every such regular language is defined as the

intersection of a^*b^* with a finite Boolean combination of constraints of the form

$$\underbrace{\#_{\sigma} = n}_{\text{letter } \sigma \in \{a, b\} \text{ appears exactly } n \text{ times}} \quad \underbrace{\#_{\sigma} \equiv k \pmod n}_{\text{the number of appearances of } \sigma \in \{a, b\} \text{ is congruent to } k \text{ modulo } n}$$

It follows that L_{β} is a finite Boolean combination of languages of the form “exactly n ports”, “exactly n isolated vertices”, “the number of ports is congruent to k modulo n ”, “the number of isolated vertices is congruent to k modulo n ”. All of these languages are easily seen to be recognisable, using a construction similar to Example 42.

This completes the proof of Courcelle’s Theorem.

Exercises

Exercise 178. (1) Consider graphs (not hypergraphs). Show that the existence of an Eulerian cycle can be defined in counting mso, but not in mso.

Exercise 179. (2) Show that the existence of a Hamiltonian cycle cannot be defined in counting mso with set quantification restricted to sets of vertices (and not hyperedges).

Exercise 180. (1) Show Lemma 6.13 ceases to be true if we allow homomorphisms that are not necessarily letter-to-letter.

Exercise 181. (2) Show that for every mso formula $\varphi(X)$ with one free set variable, the following problem can be solved in linear time:

- *Input.* A tree decomposition T ;
- *Output.* The maximal size of a set of vertices X , such that $\varphi(X)$ is true in the underlying hypergraph.

6.2.4 Satisfiability for bounded treewidth

We finish this section with an algorithm for deciding satisfiability of counting mso. Recall that already first-order logic on graphs has undecidable satisfiability, and this undecidability carries over to the more general setting of hy-

pergraphs and counting mso. We recover decidability if we restrict attention to hypergraphs of bounded treewidth.

Theorem 6.18. *The following problem is decidable:*

- **Input.** *A sentence of counting mso and $k \in \{1, 2, \dots\}$.*
- **Question.** *Is the sentence true in some hypergraph of treewidth at most k ?*

Proof We use the proof of Courcelle’s Theorem, with an emphasis on computability. We say that ranked set is *computable* if its elements can be represented in a finite way, and there is an algorithm which inputs an arity k and either outputs the finite list of all elements of arity k (if there are finitely many), or starts enumerating these elements (if there are infinitely many). We say that an algebra A is *computable* if its underlying ranked set is computable, and its multiplication operation is also computable (the inputs to the multiplication are finite hypergraphs, which can be represented in a finite way).

Free algebras over computable alphabets are computable, all of the algebras that we used as recognisers for the atomic relations in the proof Courcelle’s Theorem are computable, and computability is preserved under the products and the powerset construction. Therefore, we get the following computable strengthening of Courcelle’s Theorem: given a sentence of counting mso, which defines a property of hypergraphs over a finite alphabet Σ , we can compute a recognising homomorphism

$$h : H\Sigma \rightarrow A$$

into a computable algebra. The algebra, homomorphism, and accepting set are represented by the corresponding algorithms.

Let $A_k \subseteq A$ be the image under h of all hypergraphs with treewidth at most k . By Theorem 6.9, A_k is equal to the smallest subset of A that contains the letters and which is closed under applying the treewidth k operations in the algebra A . Since A_k is contained in the finite part of A which has arity at most $k + 1$, and there are finitely many treewidth k operations, it follows that A_k can be computed. Finally, we check if A_k contains at least one element of the accepting set. \square

Exercises

Exercise 182. (2) Show that the following problem is decidable: given a first-order formula, decide if it is true in some grid. Here we are talking about unlabelled grids as in Example 37, and not labelled grids as in Example 39.

Exercise 183. (2) Show that if a hypergraph language L has bounded treewidth and is definable in counting mso, then its syntactic algebra is computable.

Exercise 184. (2) Show that the following problem is decidable: given $k \in \{1, 2, \dots\}$ and an mso sentence φ , decide if φ is true in infinitely many hypergraphs of treewidth at most k .

6.3 Definable tree decompositions

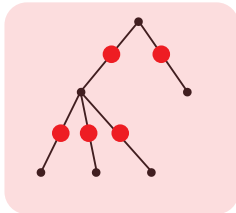
In this section, we show that for hypergraphs of bounded treewidth, tree decompositions can be defined in mso. The general idea is that an mso formula can guess a colouring of the hypergraph with a small number of colours, and use that colouring to recover a tree decomposition. As an application of this result, we will show in Section 6.3.1 that for hypergraphs of bounded treewidth, the converse of Courcelle's Theorem is also true: every recognisable property is definable in counting mso for hypergraphs of bounded treewidth.

Definable tree decompositions. We begin by explaining how a tree decomposition can be defined in mso. This is split into two ingredients: in Definition 6.19 we represent a tree decomposition using two binary relations on vertices in the underlying graph, and then in Definition 6.20 we explain how these binary relations can be defined in mso.

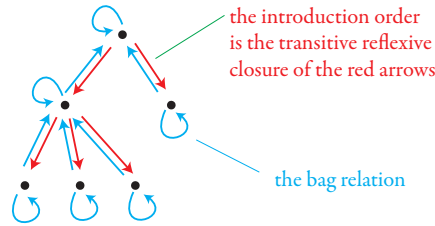
Definition 6.19 (Introduction and covering). For a tree decomposition, define two binary relations on vertices in the underlying hypergraph as follows:

- *Introduction order.* We say that vertex v is introduced before vertex w if v is introduced in an ancestor of the node that introduces w .
- *Bag relation.* We say that vertex v is in the bag of vertex w if w is in the bag of the node where v is introduced.

Here is a picture of the introduction order and bag relations for a width 1 tree decomposition:



a hypergraph which is a tree



its introduction order and bag relation

For the rest of this section, we only consider tree decompositions every two nodes have different bags; this assumption can be ensured without changing the width by merging nodes with the same bag. If all nodes have different bags, which means that every node introduces at least one vertex, then a tree decomposition can easily be recovered from its introduction ordering and bag relation, by defining the nodes to be vertices (modulo being introduced in the same node).

To represent an introduction and covering relations in the underlying graph, we will define them using an mso formula of constant size with set parameters, as described in the following definition.

Definition 6.20 (Relations defined by formulas with set parameters). A mso formula with set parameters is an mso formula of the form

$$\varphi(\underbrace{Y_1, \dots, Y_n}_{\substack{\text{set variables} \\ \text{called} \\ \text{set parameters}}}, \underbrace{x_1, \dots, x_m}_{\substack{\text{element variables} \\ \text{called} \\ \text{arguments}}})$$

We say that an m -ary relation R on elements in some model is definable by φ if there is some valuation of the set parameters, such that the resulting m -ary relation on arguments is exactly R .

Note that the logic used in the above definition is mso, and not counting mso. It will turn out that counting is not needed to define tree decompositions. We will be mainly interested in mso formulas with two arguments (and some number of set parameters), which will be used to define the introduction ordering and the bag relation in tree decompositions.

Definition 6.21 (Definable tree decomposition). We say that a set $L \subseteq \text{H}\Sigma$ of hypergraphs has *definable tree decompositions* if there is some $\ell \in \{0, 1, \dots\}$ and mso formulas φ, ψ with two arguments and set parameters, such that ev-

ery hypergraph in L has a tree decomposition of width at most ℓ where the introduction ordering is definable by φ , and the bag relation is definable by φ .

In the above definition, we need two formulas, one for the introduction ordering and one for the bag relation. In fact, these formulas could be merged into a single formula, thanks to the following observation.

Lemma 6.22. *Suppose that φ_1, φ_2 are MSO formulas with set parameters, with the same number of arguments. There is an MSO formula φ with set parameters which defines every relation definable by either φ_1 or φ_2 .*

Proof Let the formulas from the assumption be

$$\{\varphi_i(Y_1, \dots, Y_{n_i}, x_1, \dots, x_m)\}_{i \in \{1,2\}}.$$

The formula from the conclusion adds an extra set parameter Y_0 , whose emptiness determines which of the two formulas φ_1 or φ_2 should be used:

$$(Y_0 = \emptyset \Rightarrow \varphi_1) \wedge (Y_0 \neq \emptyset \Rightarrow \varphi_2).$$

The number of parameters is one plus the maximal number of parameters in φ_1 and φ_2 . □

Thanks to the above lemma, we can assume without loss of generality that the introduction order and the bag relation in a tree decomposition are both defined by the same MSO formula φ with set parameters (using, of course, different values for the set parameters). In such a case, we say that the tree decomposition is defined by φ .

We are now ready to state the main result of this section.

Theorem 6.23. *If $L \subseteq H\Sigma$ has bounded treewidth, then it has definable tree decompositions.*

The definable tree decompositions in the conclusion of the theorem will have sub-optimal width, i.e. if the hypergraphs in L have treewidth $\leq k$, then the defined tree decompositions will have width ℓ which is doubly exponential in k . With more care in the proof, we could produce optimal width tree decompositions¹⁰, but the sub-optimal width will be enough for our purposes, namely proving Corollary 6.24, which says that recognisability implies definability for hypergraphs of bounded treewidth.

The MSO formulas defining the tree decompositions will not care about the

¹⁰ Definability of tree decompositions of optimal width is shown in [7] Bojańczyk and Pilipczuk, “Optimizing Tree Decompositions in MSO”, 2017, Theorem 2

labelling relation $a(e)$, but only about the incidence relation $v = e[i]$. For this reason, we will not specify the alphabet Σ for the rest of this section.

The rest of Section 6.3 is devoted to proving Theorem 6.23 and a corollary which shows the converse of Courcelle's Theorem.

Exercises

Exercise 185. (2) Define mso_1 to be the variant of mso where set quantification is restricted to sets of vertices (and not hyperedges). Show that for hypergraphs of bounded treewidth, mso_1 has the same expressive power as mso .

Exercise 186. (2) Let Σ be a ranked alphabet. Show that there is no mso formula φ such that for every hypergraph in $\text{H}\Sigma$ there is a linear order on the vertices that is definable by φ , for some choice of set parameters.

Exercise 187. (1) Show that for every $\ell \in \{0, 1, \dots\}$ there is no mso formula which defines every tree decomposition of width ℓ . Hint: use Exercise 186.

Exercise 188. (2) Show that a formula as in Exercise 186 can be found, if we want the linear order only for connected hypergraphs with of degree at most k (every vertex is adjacent to at most k hyperedges).

Exercise 189. (2) Show that a formula as in Exercise 186 can be found, if we want a spanning forest instead of a linear order.

Exercise 190. (2) We say that a set of hypergraphs L has *bounded treedepth* if there is some ℓ such that every hypergraph in L has a tree decomposition of width and height at most ℓ (the height is the maximal depth of nodes). Without using Theorem 6.23, show that if L is recognisable and has bounded treedepth, then it is definable counting mso .

6.3.1 Application to recognisability

Before proving Theorem 6.23 about the existence of definable tree decompositions, we apply the theorem to get a converse of Courcelle's Theorem for hypergraphs of bounded treewidth, which is Corollary 6.24 below. Apart from explaining the use of definable tree decompositions, the application will be

a motivation to introduce some terminology that will be used in the proof of Theorem 6.23.

Corollary 6.24. *If $L \subseteq \text{H}\Sigma$ is recognisable and has bounded treewidth, then it is definable in counting MSO.*

The idea is simple. By Theorem 6.23, there is some MSO formula φ and some ℓ such that every hypergraph in L has a tree decomposition of width at most ℓ which is definable by φ . The formula defining the language L guesses this tree decomposition, by existentially quantifying over the set parameters needed to define the introduction order and the bag relation, and then does a bottom-up pass through the decomposition to compute the value of the hypergraph with respect to the homomorphisms which recognises L .

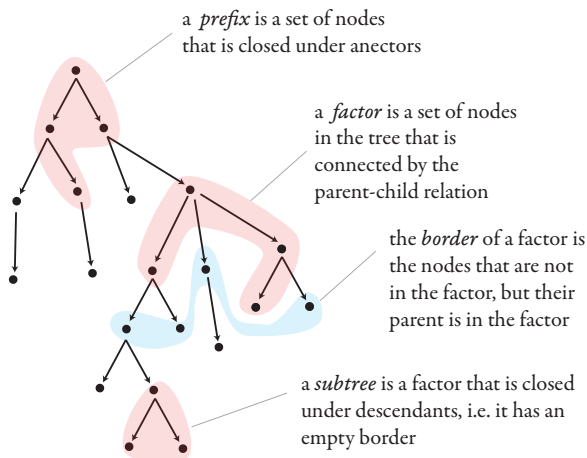
This idea is described in more detail below. To explain what we mean by a bottom-up pass through the tree decomposition, we begin by explaining how to assign a hypergraph to each subtree of a tree decomposition, which will be called the *torso* of the subtree.

Local colouring. In the torso of a subtree in the tree decomposition, there will be ports that describe the *adhesion* of the subtree, which is the vertices that appear in the subtree and in the rest of the tree decomposition. Since the ports in a hypergraph are numbered, we need a way of numbering the vertices in the adhesion. Furthermore, for different subtrees the numberings of the adhesions should be consistent with each other. To achieve this consistency, we use local colours, as defined below.

Definition 6.25. Define a *local colouring* of a tree decomposition of width k to be a colouring of vertices in the underlying hypergraph with colours $\{0, \dots, k\}$ such that for every bag, all vertices in the bag have different colours.

Every tree decomposition has at least one local colouring, which can be obtained in a greedy way by colouring the root bag, then colouring the bags of the children, and so on. If a tree decomposition is equipped with a local colouring, then every bag has an implicit total order, from the smallest colour to the biggest colour. For the rest of this section, we assume that every tree decomposition comes with a local colouring.

Factors and torsos. When working with tree decompositions, we use the following tree terminology:



Every factor has a root, which is the minimal node in the factor. A prefix is the special case of a factor where the root node is also the root of the entire tree.

In the following definition, to every factor in a tree decomposition, we assign its torso, which describes the part of the underlying hypergraph that is covered by the factor. For the purpose of Corollary 6.24, we will use torsos for subtrees, but later in the proof of Theorem 6.23 we will also use torsos for factors that are not subtrees, so we give the definition in full generality below.

Definition 6.26 (Torso). For a factor X of a tree decomposition T , define its torso T/X to be the hypergraph which is obtained as follows. The vertices of the torso are vertices of the underlying hypergraph which appear in X , i.e. they appear in at least one bag. The ports are the *adhesion* of X , i.e. vertices that appear in X but are not introduced in X (this is the same as the adhesion of $\{x\}$ where x is the root node of X). The ports of the torso are ordered according to the local colouring in the tree decomposition. There are two kinds of hyperedges in the torso:

- *non-border*: hyperedges of the underlying hypergraph which are introduced in nodes from T (the node of a tree decomposition which introduces a hyperedge is the least node of the tree decomposition which covers the hyperedge);
- *border*: for every node x of the tree decomposition which is in the border of X , there is a hyperedge whose incidence list is the adhesion of $\{x\}$, ordered according to the local colouring in the tree decomposition.

In the above definition, we do not specify the label of the border hyperedge.

This is because, as mentioned before, we do not use the labels in the hypergraph when defining tree decompositions.

Equipped with the terminology above, we are ready to prove Corollary 6.24, which says that recognisability implies definability in counting mso for languages of bounded treewidth.

Proof Suppose that L has bounded treewidth and is recognised by a homomorphism

$$h : H\Sigma \rightarrow A$$

into a hypergraph algebra that is finite on every arity. By Theorem 6.23 there exists $k \in \{0, 1, \dots\}$ and an mso formula φ with set parameters, such that every hypergraph in L has a tree decomposition T of width at most k that can be defined by φ .

The formula defining the language L works as follows. Let $G \in L$ and let T be the corresponding definable tree decomposition. For a node x of this tree decomposition, define G_x to be the torso of the subtree of x . The formula defining L first uses existential set quantification to guess: (a) set parameters for the formula φ which define the introduction order of T , (b) set parameters for the formula φ which define the bag relation of T ; (c) a local colouring of vertices which uses $\{0, \dots, k\}$; (d) for every node x of the tree decomposition T an element a_x of the algebra A such that

$$a_x = h(G_x). \tag{6.1}$$

The elements from item (d) are represented by storing the same value a_x in every vertex v that is introduced in node x . Note that the arity of a_x is at most $k + 1$, and therefore there are finitely many possibilities for a_x , which means that the elements from item (d) can be represented using finitely many set quantifiers. Next, the formula checks that the sets guessed in (a, b, c) indeed describe a tree decomposition of width at most k , that the equation (6.1) holds for every node x , and that if x is the root node of the tree decomposition then a_x belongs to the accepting set $h(L)$.

In the rest of the proof, we describe how (6.1) can be checked in counting mso. Recall the *fusion* operation described in Figure 6.2, which inputs two hypergraphs of the same arity, and fuses them along the ports. For fixed arity of input hypergraphs, the fusion operation is associative and commutative, and therefore it can be seen as an operation which inputs a multiset of hypergraphs of same arity. In the following, we use additive notation $G + H$ for fusion. The idea is that a_x is obtained by applying fusion to elements which correspond to

children of x in the tree decomposition, and fusion in a finite algebra can be computed in counting mso. A more detailed explanation follows below.

Let x be a node of the tree decomposition T . We define hypergraphs that correspond to hyperedges introduced in x and child node of x as follows:

- For a hyperedge e that is introduced in x , define E_e to be the hypergraph where the vertices are the bag of x , all vertices are ports, and which has a unique hyperedge, namely e .
- For a child node y of x in the tree decomposition T , define H_y to be the hypergraph that results from G_y by adding all vertices that are in the bag of x but not in the subtree of y , and setting the ports to be the bag of x .

In the hypergraph H_y from the second item, the hyperedges are the same as in G_y , it is only the vertices and ports that changed. All of the hypergraphs G_x , E_e and H_y have the same ports, namely the bag of x , ordered according to the local ordering in the tree decomposition. Therefore, we can use the fusion operation to get the following equality:

$$G_x = \underbrace{\sum_y H_y}_{\text{fusion ranging over children of } x} + \underbrace{\sum_e H_e}_{\text{fusion ranging over hyperedges introduced in } x}$$

By definition, the hypergraph H_y is obtained from G_y by forgetting some ports (as in the forget term operation from Figure 6.2) and adding new port vertices and renumbering the old ports (as in the rearrangement term operation from Figure 6.2). Therefore, the operation which transforms H_y into G_y is a term operation, call it f_y . Therefore, we have the following equation for every node x :

$$G_x = \sum_y f_y(G_y) + \sum_e H_e.$$

Since h is a homomorphism, and homomorphisms commute with term operations such as fusion and f_y , we have the following equation:

$$a_x = \sum_y f_y(a_y) + \sum_e h(H_e), \quad (6.2)$$

where both fusion and the term operations f_y are now interpreted in the finite algebra A . A straightforward bottom-up induction over nodes in the tree decomposition shows that (6.1) holds for every node x if and only if (6.2) holds for every node x . Therefore, it remains to show that a formula of counting mso can check (6.2).

Let A_x denote the elements of the algebra A whose arity is equal to the number of ports in G_x . In the equation (6.2), the term operation f_y has type $A_y \rightarrow A_x$, while fusion $+$ is applied to elements of A_x . Since fusion is associative and commutative, it follows that the fusion from (6.2) corresponds to multiplication in a finite commutative semigroup, and such multiplication can be computed in counting mso. The result follows, since the term operation f_y can be determined in mso based on the node y (given by a vertex introduced in it), and the same is true for $h(H_e)$. \square

Exercises

Exercise 191. (2) We say that a hypergraph algebra A is *aperiodic* if for every $n \in \{0, 1, \dots\}$ the semigroup

(elements of A with arity n , fusion)

is aperiodic. Show that if L is a set of hypergraphs of bounded treewidth, then L is definable in mso without counting if and only if it is recognised by a finite aperiodic hypergraph algebra.

Exercise 192. (2) Suppose that L has bounded treedepth, as described in Exercise 190. Show that if L is recognised by a aperiodic hypergraph algebra, then it is definable in first-order logic.

6.3.2 Nested tree decompositions

The rest of Section 6.3 is devoted to proving Theorem 6.23 about definable tree decompositions. The proof has three steps.

In this section, we prove the Merging Lemma, which shows how a definable tree decomposition of definable tree decompositions can be transformed into a single definable tree decomposition. Next, in Section 6.3.3 we prove the special case of Theorem 6.23, which says that if a set of hypergraphs has bounded pathwidth, then it has definable tree decompositions. Finally, in Section 6.3.4, we prove the theorem in its general form.

The Merging Lemma is based on the following simple idea. Suppose that we have an “external” tree decomposition, possibly of unbounded width, where every bag has an accompanying “internal” tree decomposition width at most k . We will show that if the external and internal tree decompositions can be

defined by bounded size formulas, then they can be merged into a tree decomposition which is also defined by bounded size formulas, and whose width is at most k . This lemma will be used several times in the proof, where the internal tree decompositions will typically be obtained by applying some kind of induction assumption.

Lemma 6.27 (Merging Lemma). *For every $k \in \{1, 2, \dots\}$ and every an MSO formula φ with set parameters, there is an MSO formula ψ with set parameters such that the following holds. Suppose that*

$$\underbrace{T}_{\text{external}} \quad \underbrace{\{T_x\}_{x \in \text{nodes of } T}}_{\text{internal}}$$

are tree decompositions such that each T_x is a width $\leq k$ tree decomposition of the torso $T/\{x\}$. Then the underlying hypergraph of T has a width $\leq k$ tree decomposition that is definable by ψ .

Proof Let G be the underlying hypergraph of the external tree decomposition T . We will merge the internal tree decompositions into a tree decomposition of G , which we call the *merged tree decomposition*. Nodes of the merged tree decomposition are pairs (x, y) such that x is a node of the external tree decompositions T and y is a node of the internal tree decomposition T_x . The bag of node (x, y) in the merged tree decomposition is the bag of node y in the internal tree decomposition T_x . In particular, the width of the merged tree decomposition is at most k . The parent of a node (x, y) in the merged tree decomposition is defined as follows:

- (1) If y has a defined parent y' in the internal tree decomposition T_x , then the parent of (x, y) in the merged tree decomposition is (x, y') .
- (2) If the previous case does not hold, then y is the root of the internal tree decomposition T_x . If also x is the root of the external tree decomposition, then (x, y) is the root of the merged tree decomposition and its parent is undefined. Otherwise, let x' be the parent of x in the external tree decomposition. Since x is in the border of $\{x'\}$, there must be a border hyperedge in the torso of $T/\{x'\}$ which corresponds to x . The parent of (x, y) in the merged tree decomposition is defined to be (x', y') , where y' is the node of the internal tree decomposition $T_{x'}$ that introduces the border hyperedge corresponding to x .

Claim 6.28. *The merged tree decomposition is a tree decomposition.*

Proof Every hyperedge of G is a hyperedge in some torso $T/\{x\}$, and therefore it is covered by some bag of some internal tree decomposition, which is

also a bag of the merged tree decomposition. A short analysis shows that if v is a vertex of G , then it is introduced in exactly one node of the merged tree decomposition, namely the pair (x, y) where x is the unique node of external tree decomposition where v is introduced and y is the unique node of the internal tree decomposition T_x where v is introduced. \square

It remains to prove that introduction order and the bag relation of the merged tree decomposition can be defined by an mso formula that depends only on φ and k . For a vertex v in G , let us write $[v]$ for the node of the external tree decomposition where v is introduced. It is not hard to see that the introduction ordering and the bag relation for the merged tree decomposition can be defined in mso using the following ternary relations on vertices of G :

- $A(u, v, w)$: u is introduced before v in the internal tree decomposition $T_{[w]}$;
- $B(u, v, w)$: u is in the bag of v in the internal tree decomposition $T_{[w]}$;

Therefore, to finish the proof of the Merging Lemma, it remains to show that the ternary relations A and B can be defined by mso formulas that depend only on φ and k . We only do the proof for the relation A , the proof for B is the same.

Suppose that the formula φ from the assumption of the lemma has n set parameters:

$$\varphi(Y_1, \dots, Y_n, u, v)$$

By assumption of the lemma, for every node x in the external tree decomposition there is a choice of set parameters

$$Y_1^x, \dots, Y_n^x \subseteq \text{vertices and hyperedges in the torso } T//\{x\}$$

such that fixing this choice defines the introduction ordering of the internal tree decomposition T_x . To define the relation A , we need to aggregate all of these choices of set parameters into a single choice of set parameters in G . This is done in the following claim.

Claim 6.29. *For every $i \in \{1, \dots, n\}$, the following relations on G can be defined by mso formulas that depend only on i, k and φ .*

- $C_i(\alpha, v)$: α is a vertex or hyperedge, and v is a vertex, such that $\alpha \in Y_i^{[v]}$;
- $D_i(v)$: v is a vertex such that $[v]$ has a parent node x in the external tree decomposition, and Y_i^x contains the border hyperedge corresponding to $[v]$.

Proof The unary relation D_i is simply a set of vertices, so it can be defined using one set parameter. We focus on the binary relation C_i . Each set parameter Y_i^x uses three kinds of elements:

- border hyperedges of $T/\{x\}$;
- ports of $T/\{x\}$, which are the same as the adhesion of $\{x\}$;
- hyperedges and vertices of G that are introduced in x ;

The first kind does not participate in the relation C_i , because the border hyperedges are not hyperedges of G . Therefore, we need to represent only elements of the second and third kinds. For the second kind, there are at most k elements for every node x , and therefore elements of the second kind can be represented by the sets $Z_{i,1}, \dots, Z_{i,k}$ that are defined by

$$Z_{i,j} = \{v : Y_i^{[v]} \text{ contains the } j\text{-th port}\}. \quad (6.3)$$

Finally, elements of the third kind are disjoint for different choices of x , and therefore they can be aggregated using set union:

$$\bigcup_x \text{vertices and hyperedges that are introduced in } x \text{ and belong to } X_i^x. \quad (6.4)$$

Since the vertices, hyperedges and ports of the torso $T/\{[v]\}$ can be recovered in mso using the formulas defining the external tree decompositions, it follows that the relation C_i can be recovered from the set parameters (6.3) and 6.4. \square

Using the relations $\{C_i, D_i\}_i$ from the above claim, we can define the ternary relation A in mso, by running the formulas defining the internal tree decompositions on the set parameters recovered from $\{C_i, D_i\}_i$. The same argument works for B , and completes the proof of the Merging Lemma. \square

Exercises

Exercise 193. (1) A cut hyperedge in a hypergraph is a hyperedge e such that for some two vertices, every path connecting them must pass through e . Let L, K be sets of hypergraphs, such that for every $G \in K$, if all cut hyperedges are removed from G , then the resulting hypergraph is in L . Show that if L has definable tree decompositions, then the same is true for K .

6.3.3 Bounded pathwidth

In this section, we present the second step in the proof of Theorem 6.23. We show Theorem 6.30 below, which says that bounded pathwidth implies definable tree decompositions. Recall that the pathwidth of a hypergraph is

defined to be the smallest width of an associated path decomposition, and a path decomposition is defined to be the special case of a tree decomposition where all nodes are totally ordered by the ancestor relation.

Theorem 6.30. *If a set of hypergraphs has bounded pathwidth, then it has definable tree decompositions of bounded width.*

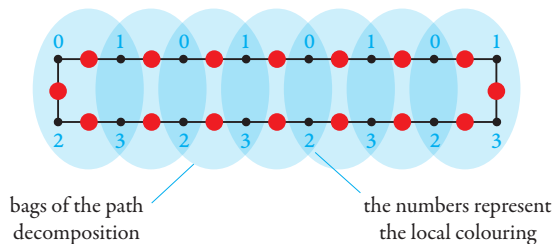
There is an asymmetry in the above theorem: the assumption has path decompositions but the conclusion has tree decompositions. To see why this is necessary, consider a hypergraph with no ports or hyperedges, but with many isolated vertices, like this one:



This hypergraph has a path decomposition of width 0, where every vertex is in its own bag. However, the introduction ordering for such a path decomposition is the same thing as a linear order on the vertices of the hypergraph. To define linear orders on hypergraphs without hyperedges, one needs formulas of mso of unbounded size, see Exercise 186.

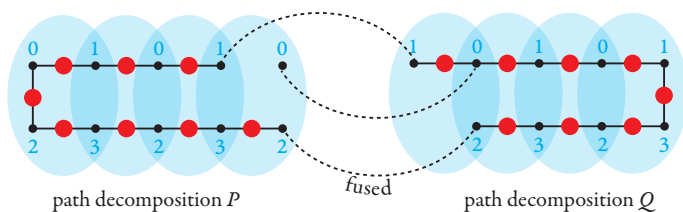
The rest of Section 6.3.3 is devoted to proving Theorem 6.30. To find definable tree decompositions, we will view path decompositions as semigroup, and use the Factorisation Forest Theorem.

Path decompositions as a semigroup. Fix $k \in \{1, 2, \dots\}$ for the rest of this section. We will view path decompositions of width at most k as a semigroup, as described below. The elements of this semigroup are path decompositions of width at most k . We assume that a path decomposition comes together with its underlying hypergraph, and a local colouring with colours $\{0, \dots, k\}$. Here is a picture of a path decomposition:



We also assume that the ports are exactly the vertices in the first bag, and that they are numbered so that the local colouring is increasing.

The multiplication $P \cdot Q$ of two path decompositions P, Q is defined as follows. The underlying hypergraph is obtained by taking the disjoint union of the two underlying hypergraphs, and fusing vertices using the local colouring at follows: for every $i \in \{0, \dots, k\}$, if a vertex with local colour i appears both in the last bag of P and in the first bag of Q , then these two vertices are fused into a single vertex. The path decomposition is then simply all bags of P followed by all bags of Q , with the last bag of P and first bag of Q merged into a single bag. For example, the path decomposition in the previous picture can be obtained as the multiplication $P \cdot Q$ described in the following picture:



It is easy to see that the operation $P \cdot Q$ is associative, and hence path decompositions of width at most k are a semigroup. This semigroup is generated by the finite set of path decompositions which have width at most k and at most two bags.

The reachability homomorphism. To apply the Factorisation Forest Theorem, we define a semigroup homomorphism from the semigroup of path decompositions of width at most k , which stores a finite amount of information about paths in the input path decomposition. Recall that a path in a hypergraph is a sequence of vertices and hyperedges, which alternates between vertices and hyperedges, such that every hyperedge in the sequence is incident to the vertices on the next and previous positions in the sequence. An *inner path* in hypergraph is the special case of a path which is not allowed to use port vertices (i.e. vertices from the first bag), with the possible exception of the source and target of the path. Define the *reachability homomorphism* to be the function which maps a path decomposition of width at most k to the answers to the following questions, for all $i, j \in \{0, \dots, k\}$ and $\sigma, \tau \in \{\text{first}, \text{last}\}$:

- does the σ bag contain a vertex with local colour i ?
- is there an inner path from a vertex with local colour i in the σ bag to a vertex with local colour j in the τ bag?
- is there are vertex with local colour i that is both in the first and last bag?

It is not hard to see that the reachability homomorphism is compositional in

the semigroup sense, and therefore it can be viewed as a semigroup homomorphism from the semigroup of path decompositions of width at most k into a finite semigroup (which consists of all possible sets of answers to the finitely many questions described above).

Using the Factorisation Forest Theorem. As we have already mentioned, the semigroup of path decompositions of width at most k is finitely generated, namely by the set Δ of path decompositions that have at most two bags. Therefore, we can view the reachability homomorphisms as a homomorphism

$$h : \Delta^+ \rightarrow S$$

to which we can apply the Factorisation Forest Theorem. For $\ell \in \{0, 1, \dots\}$ define L_ℓ to be the underlying hypergraphs of path decompositions in Δ^+ which have an h -factorisation tree of height at most ℓ . By the Factorisation Forest Theorem, there is some ℓ such that all hypergraphs of pathwidth at most k belong to L_ℓ . Therefore, to prove Theorem 6.30 about definable tree decompositions for hypergraphs of bounded pathwidth, it remains to show the following lemma.

Lemma 6.31. *For every $\ell \in \{0, 1, \dots\}$, L_ℓ has definable tree decompositions.*

Proof For the induction step, it will be useful to prove a slightly stronger result, which involves inner components. For a vertex or hyperedge x in a hypergraph, define its *inner component* to be the set of vertices and hyperedges that can be reached from x via some inner path, plus all the ports. A set X of vertices and hyperedges is called an inner component if it is the inner component of some vertex or hyperedge. If X is an inner component in a hypergraph G , then we write $G|X$ for the hypergraph that is obtained from G by keeping only the vertices and hyperedges from X (the ports are the same, since inner components contain all ports). Every hypergraph is the fusion of its inner components X :

$$G = \sum_X G|X.$$

By induction on $\ell \in \{0, 1, \dots\}$, we will prove that the language

$$K_\ell = \{G|X : G \in L_\ell \text{ and } X \text{ is an inner component}\}$$

has definable tree decompositions. This will imply the lemma, as explained in the following claim.

Claim 6.32. *Let L be a set of hypergraphs. If*

$$\{G|X : G \in L \text{ and } X \text{ is an inner component}\}$$

has definable tree decompositions, then so does L .

Proof We will find a definable tree decomposition for every hypergraph $G \in L$ using the Merging Lemma. Define an external tree decomposition which has a root plus one node for every inner component. The bag of the root is the ports, while for every inner component the bag of the corresponding node is the vertices from the inner component. Since the ports and the inner components can be defined in mso, the external tree decomposition can be defined using formulas of constant size. The root bag has constant size, while the internal tree decompositions for the remaining nodes are definable thanks to the assumption on L . Therefore, we can apply the Merging Lemma to get a definable tree decomposition for G . \square

It remains to prove that K_ℓ has definable tree decompositions for every $\ell \in \{0, 1, \dots\}$. The proof is by induction on ℓ . For the induction base of $\ell = 0$, hypergraphs in K_0 have at most $k + 1$ vertices, and therefore we can use trivial definable tree decompositions where all vertices are in the same bag.

We are left with the induction step. Consider a hypergraph in $K_{\ell+1}$, which by definition of this language is of the form $G|X$ where $G \in L_{\ell+1}$ and X is an inner component. By definition of $L_{\ell+1}$, there is a path decomposition P of G which can be factorised as

$$P = P_1 \cdots P_n$$

in the semigroup of path decompositions, so that for every $i \in \{1, \dots, n\}$, P_i is a path decomposition whose underlying hypergraph G_i belongs to L_ℓ . Also, either $n = 2$, or all of the path decompositions P_1, \dots, P_n have the value under the reachability homomorphism which is furthermore idempotent. For $i \in \{1, \dots, n\}$, define the *i -th embedding* to be the function which maps vertices and hyperedges of G_i to the corresponding vertices and hyperedges in G . Define X_i to be the vertices and hyperedges in G_i which are in the inverse image of X under the i -th embedding. Let Q be the path decomposition defined as follows. The nodes are $\{1, \dots, n\}$ ordered in the natural way, and the bag of node i is the vertices from X_i . It is easy to see that this is a path decomposition of $G|X$. We will apply the Merging Lemma, with the external tree decomposition being the path decomposition Q and the internal tree decompositions coming from the following claim.

Claim 6.33. *For every $i \in \{1, \dots, n\}$, the torso $Q/\{i\}$ is taken from a set of hypergraphs with definable tree decompositions.*

Proof Since the i -th embedding maps inner paths in G_i to inner paths in G , it follows that X_i is a union of inner components of X_i . Therefore, $G_i|X_i$ can be obtained by fusing a family of graphs from K_ℓ , and as such it has a definable tree decomposition. Finally, in order to get the torso of $\{i\}$, we need some bag to cover the border hyperedge, which represents node $i + 1$. Since there is only one border hyperedge, we can simply modify the tree decomposition of $G_i|X_i$ so that all that vertices incident to this border hyperedge are in all bags; the width of the resulting tree decomposition changes by a constant amount, as do the mso formulas needed to define it. \square

It remains to show that the external path decomposition Q is definable. If $n = 2$ then there is not much to do: the external path decomposition has two nodes, and therefore it can be defined by formulas using two set parameters that say which vertices are in which bags. We are left with the case where $n > 2$ and all path decompositions P_1, \dots, P_n have the idempotent image under the reachability homomorphism. We begin with the following observation about how vertices are distributed in the bags of the external path decompositions T .

Claim 6.34. *Every vertex $v \in X$ appears either in all bags of the external tree decomposition, or at most two consecutive bags.*

Proof If a vertex appears in at most two bags, then they must be consecutive by definition of path decompositions. Suppose that v is introduced in node i of the external tree decomposition Q , and it also appears in bag $j > i + 1$. This means that v must appear in both the first and last bag of P_{i+1} . Since all of the path decompositions P_1, \dots, P_n have the same idempotent image under the reachability homomorphisms, it follows that v must appear in both the first and last bag of $P_1 \cdots P_{i+1}$. Applying this reasoning again, it follows that v must appear in both the first and last bag of each of the path decompositions P_1, \dots, P_n . \square

In the remainder of the proof, we will show that the introduction ordering of the external path decomposition Q can be defined by constant size formulas with set parameters. In light of the above claim, the bag relation will also be easily seen to be definable, since vertices can be partitioned into three kinds: those which appear in all bags, those which appear only in the node where they are introduced, and those which appear only in the node where they are introduced and the next one.

Define the *index* of a vertex v to be the node $i \in \{1, \dots, n\}$ of the external path decomposition Q where x is introduced. The key to defining the introduction ordering of Q is the following claim, which is where we use the assumption that X is an inner component.

Claim 6.35. *Let $v, w \in X$ be vertices with indices $1 < i \leq j$. Then*

- (1) *every inner path in $G|X$ that connects them uses all indices in $\{i, \dots, j\}$;*
- (2) *some connecting inner path in $G|X$ uses only indices in $\{i - 1, \dots, j + 1\}$.*

Proof Neither v nor w can be a port of G , since ports must belong to the root bag of the path decomposition, and therefore they have index 1.

- (1) Using the same proof as in Claim 6.34, we show that if a non-port vertex v is incident to a hyperedge e introduced in node i of Q , then its index is i or $i - 1$. This implies that indices of vertices in an inner path are consecutive, which then implies (1).
- (2) Consider an inner path that connects v and w , which must exist because these vertices are in the same inner component X . By the first item, this path can be split into segments of three kinds:
 - a path in the underlying hypergraph of $P_i \cdots P_j$; or
 - a path in the underlying hypergraph of $P_1 \cdots P_{i-1}$, which begins and ends in vertices from the last bag of P_{i-1} ; or
 - a path in the underlying hypergraph of $P_{j+1} \cdots P_n$, which begins and ends in vertices from the first bag of P_{j+1} .

To prove the claim, we show that each of these paths can be modified, without changing its source or target, so that it visits only vertices with indices in $\{i - 1, \dots, j + 1\}$. For paths of the first kind, there is nothing to do. For the second case, we use the fact that P_{i-1} has the same image under the reachability homomorphisms as $P_1 \cdots P_{i-1}$, and a similar argument is used for the third kind.

□

To finish the proof of the lemma, we use the above claim and a modulo counting argument to show that the introduction ordering of the external path decomposition Q can be defined in $G|X$ by a constant size formula with set parameters.

For $m \in \{0, 1\}$, let $R_m(v, w)$ be the binary relation on vertices of $G|X$ which says that

$$\text{index of } w = m + \text{index of } v.$$

We will prove that R_0 and R_1 can be defined by a constant size formulas with set parameters. This will imply the claim, since the relation in the claim is the transitive closure of the union $R_0 \cup R_1$, and transitive closures of binary relations can be defined in mso.

Define the *modulus* of a vertex or hyperedge to be its index modulo 100. The

moduli can be described by a constant number of unary relations, and therefore they can also be defined by an mso formula with set parameters of constant size. From Claim 6.35 it follows that $R_m(v, w)$ holds if and only if there is an inner path from v to w which sees at most $m + 3$ moduli and

$$\text{modulus of } w = m + \text{modulus of } v \quad \text{modulo } 100.$$

Since the latter condition can clearly be defined by a constant size formula with set parameters, it follows that the introduction ordering of the external tree decomposition Q can be defined by a constant size formula of mso with set parameters. This in turn completes the proof of the lemma. \square

Exercises

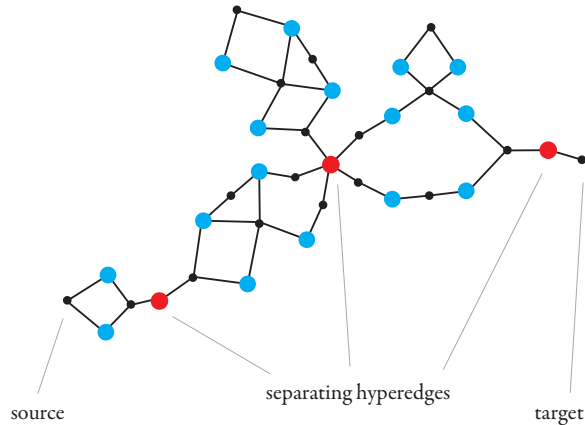
Exercise 194. (1) We say that a hypergraph is a tree if removing any hyperedge increases the number of connected components. Show that trees have unbounded pathwidth.

Exercise 195. (1) Give an algorithm which inputs a tree hypergraph, and computes its pathwidth.

6.3.4 Unbounded pathwidth

We now prove the general case of Theorem 6.23, i.e. we show that bounded tree width implies definable tree decompositions. The idea is to show that every tree decomposition can be partitioned into factors so that: (1) the torsos of the factors have bounded pathwidth; (2) the partition into factors can be defined by a bounded size formula of mso. To define tree decompositions for the torsos in item (1), we use Theorem 6.30 from the previous step in the proof. To define the partition into factors from item (2), we will use separators and paths with small overlap, as described below.

Separators. Consider two vertices in a hypergraph, called the source and target. We say that a vertex or hyperedge separates these two vertices if it must appear on every path that from the source to the target. Here is a picture:

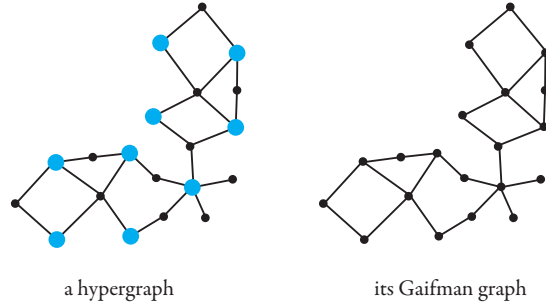


The following lemma shows that separating vertices and hyperedges are the only obstacle against finding two disjoint paths from the source to the target which are hyperedge disjoint

Lemma 6.36. *Consider a connected hypergraph with distinguished source and target vertices. There exist two paths from the source to the target such that if a vertex or hyperedge is used by both paths, then it separates the source and target.*

Proof We use Menger's Theorem about separators and disjoint paths in undirected graphs. Menger's Theorem¹¹ says that if, in an undirected graph with designated source and target vertices, one cannot find a set of k separating vertices, then there are k paths from the source to target which are disjoint (apart from the source and target). The lemma follows immediately by applying Menger's Theorem, in the case of $k = 1$, to the Gaifman graph of a hypergraph, which is an undirected graph that is defined as in the following example:

¹¹ Reinhard Diestel. *Graph theory (electronic edition)*. Vol. 173. Graduate texts in mathematics. Springer-Verlag, 2006, Theorem 3.3.3.



□

Paths with small overlap. Another ingredient that will be used in the proof of Theorem 6.23 is families of paths with small overlap. These will be used to define in mso define binary relations, such as the introduction order or bag relation of a tree decomposition.

Lemma 6.37. *Let \mathcal{P} be a family of paths in a hypergraph G such that every vertex is used by at most ℓ paths. Then the binary relation*

$$\{(s, t) : s, t \text{ are vertices such that some path in } \mathcal{P} \text{ has source } s \text{ and target } t\}$$

is definable by an mso formula with set parameters that has size $O(\ell)$.

Proof For every path $P \in \mathcal{P}$ choose a colouring

$$\lambda_P : (\text{vertices used by } P) \rightarrow \{1, \dots, \ell\}$$

so that if a vertex appears in two different paths $P, Q \in \mathcal{P}$, then it has different colours under λ_P and λ_Q . Such a family of colourings $\{\lambda_P\}_{P \in \mathcal{P}}$ can easily be defined using a greedy algorithm, thanks to the assumption that every vertex is used in at most ℓ paths. Define a directed graph H where the vertices are pairs (vertex of G , number in $\{1, \dots, \ell\}$) and which has a directed edge

$$(v, i) \rightarrow (w, j) \tag{6.5}$$

if there is some path $P \in \mathcal{P}$ where v, w are connected by a single hyperedge and $i = \lambda_P(v)$ and $j = \lambda_P(w)$. The following claim shows that the graph H can be defined in G using constant size formulas of mso.

Claim 6.38. *For every $i, j \in \{1, \dots, \ell\}$, the following binary relation on vertices of G can be defined by a constant size mso formula with set parameters:*

$$R_{ij}(u, v) \stackrel{\text{def}}{=} \text{there is a path from } (u, i) \text{ to } (v, j) \text{ in } H$$

Proof Define E_{ij} in the same way as the relation R_{ij} in the statement of the claim, except that the paths used in E_{ij} are required to have exactly one edge. For every choice of i, j , the binary relation E_{ij} contains only pairs which are connected by some hyperedge in G , and therefore it can be defined in G by a constant size formula of mso with set parameters, which uses the connecting hyperedges as set parameters. In presence of the binary relation E_{ij} , the relations R_{ij} can be defined using transitive closure.

Since the variant of transitive closure used here is a more sophisticated vectorial transitive closure, we explain the defining formula in more detail. Suppose that we want to define the relation R_{ij} . Fix a vertex u . Find the coordinate-wise least tuple of sets

$$V_1, \dots, V_\ell \subseteq \text{vertices of } G$$

which satisfies the following condition:

$$u \in V_i \wedge \bigwedge_{n,m} \forall x \forall y E_{nm}(x,y) \wedge x \in V_n \Rightarrow y \in V_m.$$

The definition of this tuple can be formalised in mso, and V_j consists of vertices v such that there is a path in H from (u, i) to (v, j) . This implies that each of the relations $\{R_{ij}\}_{ij}$ is definable in mso, using the relations $\{E_{nm}\}_{n,m}$. \square

It is not hard to see that a pair (s, t) belongs to the binary relation from the conclusion of the lemma if and only if there exist some $i, j \in \{1, \dots, \ell\}$ such that all of the following properties hold in the directed graph H : (a) there is no path with at least one edge that ends in (s, i) (b) there is a path from (s, i) to (t, j) ; and (c) there is no path with at least one edge that begins in (t, j) . All of these conditions can be defined using constant size formulas of mso, thanks to the claim above. \square

The Two Path Lemma. We are now state the Two Path Lemma, which is the main step in this section. This lemma says that for every tree decomposition of width at most k , and every choice of source and target vertices from the same bag, one can find a factor of the tree decomposition whose torso has bounded pathwidth and contains two roughly disjoint paths connecting the source and target.

The Path Lemma uses a mild assumption on tree decompositions, defined as follows. We say that a tree decomposition T is *sane* if for every factor X , the torso T/X is inner connected, which means that all vertices are in the same inner component (in other words, there is some non-port vertex or hyperedge from which all other vertices can be reached via inner paths). If a hypergraph is inner connected, then it has a sane tree decompositions of optimal width.

Indeed, if we take any tree decomposition, and we find a factor X with root x such that the torso T/X is not inner connected, then we can distribute the subtree of node x into separate subtrees, one for each inner connected component. Finally, if a hypergraph is not inner connected, then it is the fusion of all of its inner components. For these reasons, we can consider without loss of generality tree decompositions which are sane.

Lemma 6.39 (Two Path Lemma). *Let T be a sane tree decomposition of width k , and let s, t be two vertices from the bag of the some node x . There exists a factor X with root x such that:*

- (1) *the torso T/X has pathwidth at most $2k$;*
- (2) *there are two inner paths from s to t in the torso T/X , such that every border hyperedge is used by at most one of the paths.*

Proof In the proof, we will use path decompositions which satisfy a certain invariant, as described in the following claim. Maintaining the invariant is the reason why we have $2k$ instead of k in the statement of the lemma.

Claim 6.40. *For every hypergraph of pathwidth at most k , and distinguished vertices s, t there is a path decomposition of width at most $2k$ and satisfies the following property:*

- (#) *if e is a hyperedge that separates s from t , then there is some node z in P which covers e , and such that no vertex appears both in nodes $< z$ and $> z$.*

Proof For every separating hyperedge, find the node which introduces it. Next, move all connected components separated by this hyperedge to one side or the other of the introducing node. \square

To prove the lemma, we iterate the following claim starting with $X = \{x\}$.

Claim 6.41. *Let X be a factor of the tree decomposition T such that:*

- (*) *the root of X is x and there is a path decomposition P of T/X of width at most $2k$ which satisfies (#) with respect to the vertices s, t from the assumption of the lemma.*

Then either T/X has no border hyperedges which separate v and w , or otherwise one can add a new node to X so that it still satisfies ().*

Proof Suppose that there is a border hyperedge e in T/X which separates s and t . Let y be the node from the border of X which corresponds to the border hyperedge e . We will add y to the factor X and preserve the invariant (*). By the invariant, there is a node z in the path decomposition P that covers e , and

such that no vertex of the torso T/X appears both in nodes $< z$ and in nodes $> z$. Since the torso $T/\{y\}$ has at most k vertices, we can replace the contents of bag z with a path decomposition of $T/\{y\}$ that satisfies (#) and preserve the invariant. \square

Start with $\{x\}$ and keep iterating the above claim, until a factor X is reached such that T/X has no border hyperedge that separates s and t . Because the tree decomposition is sane, the torso T/X is inner connected. Therefore, thanks to Lemma 6.39 there are two paths from s to t in T/X which are disjoint on border hyperedges. \square

Partitioning a tree decomposition into factors of bounded pathwidth. We now repeatedly apply the Two Path Lemma to find a partition of a sane tree decomposition into factors, so that each factor has a torso with bounded pathwidth, and an important piece of information about the partition into factors can be defined by an mso formula of bounded size.

Lemma 6.42. *Let T be a sane tree decomposition of width k . There is a family of factors \mathcal{X} which partitions the nodes of T such that:*

- (1) *for every factor $X \in \mathcal{X}$, the torso T/X has pathwidth at most $2k$;*
- (2) *the following relation $R(v_0, \dots, v_k, w_0, \dots, w_k)$ is definable by an mso formula with set parameters whose size is bounded by a function of k :*
 - *there exists a node x of the tree decomposition T which is a root of some factor in \mathcal{X} , such that the vertices introduced by x are exactly $\{v_0, \dots, v_k\}$ and the bag of x is exactly $\{w_0, \dots, w_k\}$.*

Proof We repeatedly use the Two Path Lemma to get the family of factors, together with a family of paths with bounded overlap which describes the root bags of the factors. The induction step in the construction is described in the following claim.

Claim 6.43. *Suppose that \mathcal{X} is a family factors in T which satisfies:*

- (*) *the factors in \mathcal{X} partition some prefix X of the tree decomposition T and there exists a family of paths \mathcal{P} in the torso T/X such that:*
 - (1) *for every factor $Y \in \mathcal{X}$, the torso T/Y has pathwidth at most $2k$;*
 - (2) *for every factor $Y \in \mathcal{X}$ there is some vertex u introduced in Y such that for every vertex v in the root bag of X there is a path from u to v in \mathcal{P} ;*
 - (3) *every vertex of T/X appears in at most $2k^3 + k$ paths from \mathcal{P} ;*
 - (4) *every border hyperedge of T/X appears in at most $2k^3$ paths from \mathcal{P} .*

Then either X is all nodes of the tree decomposition, or otherwise one can add a new factor Y to \mathcal{X} so that $\mathcal{X} \cup \{Y\}$ still satisfies (*).

Proof Suppose that the prefix X is not all nodes in the tree decomposition. Choose some border hyperedge e in the torso T/X which corresponds to a border node y of X . The vertices incident to this border hyperedge are the adhesion of y , i.e. the vertices that appear in both y and its parent. The adhesion has size at most k , since bags have size at most $k+1$ and the adhesion is a proper subset of the bag (otherwise a node would have the same bag as its parent). For two vertices v, w we say that a path $P \in \mathcal{P}$ has a profile (v, w) if it contains an infix of the form

$$v \xrightarrow{e} w.$$

There are less than k^2 choices for v, w . Every path from \mathcal{P} has at most one profile, because we can assume without loss of generality that the hyperedge e is used at most once on each path (we can always eliminate loops from \mathcal{P} without affecting the invariant (*)). Choose a pair of vertices (s, t) in the adhesion of y so that the number of paths with this profile is maximal; let ℓ be the number of such paths. Apply the Two Path Lemma to the vertex pair (s, t) in the tree decomposition T , yielding a factor Y in T with root node y . We will prove that the invariant (*) is still satisfied after adding Y to the family \mathcal{X} . Clearly item (1) of (*) is satisfied, because T/Y has pathwidth at most $2k$ thanks to the Two Path Lemma. It remains to find a family of paths, call it \mathcal{R} , which will witness items (2)–(4) of (*). This family is constructed using the two paths from the Two Path Lemma as follows:

- (a) If a path from \mathcal{P} does not use the border hyperedge e , then add this path to \mathcal{R} without any changes. If a path in \mathcal{P} does use the border hyperedge e , but with a profile (v, w) which is different from (s, t) , then replace the segment

$$v \xrightarrow{e} w$$

with some inner path in T/Y that goes from v to w (which exists because the tree decomposition is sane), and add the resulting path into \mathcal{R} .

- (b) For paths with profile (s, t) we do a similar procedure as in the previous item, except that we use two inner paths instead of one. Let P_1 and P_2 be the two inner paths in the torso T/Y from the conclusion of the Two Path Lemma. Split the ℓ paths in \mathcal{P} which have profile (s, t) into two groups, of sizes $\lceil \ell/2 \rceil$ and $\lfloor \ell/2 \rfloor$. For every path in the first group, replace the segment

$$s \xrightarrow{e} t$$

with P_1 , and for every path in the second group replace this segment with

- P_2 . Add the resulting paths to \mathcal{R} . By construction, every border hyperedge of T/Y , which is also a border hyperedge of $T/(X \cup Y)$, is visited by paths from at most one of the two groups, and therefore by at most $\lceil \ell/2 \rceil$ paths.
- (c) Choose some vertex u that is introduced in y . For every vertex in the bag of y , add to \mathcal{R} a path in the torso T/Y which goes from that vertex to u . Such a path must exist by the assumption that tree decomposition T is nice.

We now argue that the set of paths \mathcal{R} defined above satisfies items (2)–(4) of the invariant, with respect to the family obtained from \mathcal{X} by adding Y .

Item (2) is satisfied thanks to the paths from \mathcal{R} that are described in item (c) in the definition of \mathcal{R} .

Item (3) says that every vertex of T/X appears in at most $2k^3 + k$ paths from \mathcal{R} . If a vertex is introduced in X , then the invariant is satisfied because the number of paths that use this vertex is the same for \mathcal{P} and \mathcal{R} . This is because we used inner paths of T/Y when defining \mathcal{R} . If a vertex is introduced in Y , then it is visited by at most

$$\underbrace{2k^3}_{\text{paths from } \mathcal{P} \text{ that used } e} + \underbrace{k}_{\text{paths added in item (c)}}$$

paths and therefore the invariant is also satisfied.

We are left with item (4), which says that every border hyperedge of T/X appears in at most $2k^3$ paths from \mathcal{R} . Consider a border hyperedge of $T/(X \cup Y)$. If this hyperedge is not a border hyperedge of T/Y , then we use the same argument as for vertices. Otherwise, the number of paths in \mathcal{R} that use this border hyperedge is at most

$$\underbrace{\text{number of paths in } \mathcal{P} \text{ that do not have profile } (s, t)}_{\text{paths added in item (a)}} + \underbrace{\lceil \ell/2 \rceil}_{\text{paths added in item (b)}} + \underbrace{k}_{\text{paths added in item (c)}}$$

We claim that the above value is at most $2k^3$. Indeed, if $\ell < 2k$, then the summands for items (a) and (b) are at most $k^2(2k - 1)$, and adding k will not exceed the threshold $2k^3$. Otherwise, if $\ell \geq 2k$, then

$$\lceil \ell/2 \rceil + k \leq \ell$$

and therefore the sum of all items at most the number of paths in \mathcal{P} which uses the border hyperedge e . \square

Iterate the above claim, starting with the empty family, until reaching a family of factors \mathcal{X} which satisfies (*) and partitions all nodes of the tree decomposition. Let \mathcal{P} be the family of paths from (*). By item (c) of (*), for every factor $Y \in \mathcal{X}$ there is a subset \mathcal{P}_Y such that: (1) all paths in \mathcal{P}_Y have the same

source, which is introduced in Y ; (2) the targets of the paths in \mathcal{P}_Y are exactly the vertices in the root bag of Y .

Define P to be the set of (source, target) pairs for the family of paths

$$\bigcup_{Y \in \mathcal{X}} \mathcal{P}_Y. \quad (6.6)$$

Thanks to the bounded overlap condition from item (3) and Lemma 6.37, the relation P can be defined by a constant size mso formula with set parameters.

Since the source node of all paths from \mathcal{P}_Y is introduced in the factor Y , and every vertex is introduced in exactly one factor, it follows that for different choices of Y , the sources of the paths \mathcal{P}_Y will be different. Therefore, a set of vertices U is equal to equal to the root bag of some factor in \mathcal{X} if and only if

$$\exists s \forall u \ u \in U \Leftrightarrow P(s, u).$$

Therefore, the relation R in the statement of the lemma can be defined by a constant size formula of mso with set parameters. \square

We are now ready to complete the proof of Theorem 6.23, about definable tree decompositions for hypergraphs of bounded treewidth. By Claim 6.32, it is enough to find definable tree decompositions for hypergraphs that are inner connected, i.e. have one inner component. Such hypergraphs have sane tree decompositions.

Let then T be a sane tree decomposition of width at most k . To prove Theorem 6.23, we will show that the underlying hypergraph has a tree decomposition that can be defined using a constant size formula of mso with set parameters. Apply Lemma 6.42 to T of width at most k , yielding a partition \mathcal{X} of its nodes into factors. Define S to be the tree decomposition where the nodes are the factors from \mathcal{X} , and the ancestor order on nodes is inherited from T . The bag of a factor $X \in \mathcal{X}$ is the union of bags of the nodes in the factor X . This is easily seen to be a tree decomposition. We will apply the Merging Lemma, with S being the external tree decomposition. A node of S is the same thing as a factor X in T , and the torso $S/\{X\}$ (which is a torso of one node) is the same thing as the torso T/X (which is a torso of a set of nodes). These torsos have pathwidth at most $2k$ thanks to the conclusion of Lemma 6.42, and therefore they have definable tree decompositions thanks to Theorem 6.30 from the previous section. We are left with the external tree decomposition S . The following lemma shows that S can be defined using constant size mso formulas with set parameters; and therefore the Merging Lemma can be applied to complete the proof Theorem 6.30.

Lemma 6.44. *The introduction order and bag relation of S can be defined by constant size MSO formulas with set parameters.*

Proof We will show that the introduction order and the bag relation of S can be defined in MSO (without additional set parameters), using the relation R from the conclusion of Lemma 6.42. The set parameters will then be used to define the relation R .

For a node x of the tree decomposition T , define V_x to be the vertices which are introduced by the tree decomposition in node x or its descendants. The key observation is in the following claim, which describes the set V_x in terms of information that is contained in the relation R .

Claim 6.45. *A vertex v belongs to V_x if and only if there is a path from v to some vertex that is introduced in x , such that the path does not visit any vertex in the adhesion of x .*

Proof From the definition of sane tree decompositions. □

Define Y to be the set of root nodes of the factors from \mathcal{X} . If we have the relation R from the conclusion of Lemma 6.42, then we can use quantification over elements of Y , since a node $y \in Y$ can be represented by any vertex introduced in the node. The introduction ordering is defined as follows: u is introduced before v in S if there exists some node $y \in Y$ such that u is introduced by T in x and $v \in V_x$. Thanks to Claim 6.45, this condition can be defined in MSO using the relation R from the conclusion of Lemma 6.42. Similarly, we can define the bag relation of S : u is in the bag of the node where v is introduced if either □

Bibliography

- [1] H Appelgate et al. *Seminar on triples and categorical homology theory*. Springer, 1969.
- [2] Mustapha Arfi. “Polynomial Operations on Rational Languages”. In: *Symposium on Theoretical Aspects of Computer Science, STACS, Passau, Germany*. 1987, pp. 198–206.
- [3] Stephen L. Bloom and Zoltán Ésik. “The equational theory of regular words”. In: *Information and Computation* 197.1 (2005), pp. 55–89.
- [4] Achim Blumensath. “Regular Tree Algebras”. In: *CoRR* abs/1808.03559 (2018).
- [5] Mikołaj Bojańczyk. “Two Monads for Graphs”. In: *CoRR* abs/1804.09408 (2018).
- [6] Mikołaj Bojańczyk and Bartek Klin. “A non-regular language of infinite trees that is recognizable by a sort-wise finite algebra”. In: *Logical Methods in Computer Science* 15.4 (2019).
- [7] Mikołaj Bojańczyk and Michal Pilipczuk. “Optimizing Tree Decompositions in MSO”. In: *Symposium on Theoretical Aspects of Computer Science, STACS, Hannover, Germany*. Ed. by Heribert Vollmer and Brigitte Vallée. Vol. 66. LIPIcs. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2017, 15:1–15:13.
- [8] J. Richard Büchi. “On a decision method in restricted second order arithmetic”. In: *Logic, Methodology and Philosophy of Science (Proc. 1960 Internat. Congr.)* Stanford, Calif.: Stanford Univ. Press, 1962, pp. 1–11.
- [9] J. Richard Büchi. “Weak second-order arithmetic and finite automata”. In: *Z. Math. Logik und Grundl. Math.* 6 (1960), pp. 66–92.
- [10] Olivier Carton, Thomas Colcombet, and Gabriele Puppis. “An algebraic approach to MSO-definability on countable linear orders”. In: *The Journal of Symbolic Logic* 83.3 (2018), pp. 1147–1189.

- [11] Julia Chuzhoy and Zihan Tan. “Towards Tight(er) Bounds for the Excluded Grid Theorem”. In: *Symposium on Discrete Algorithms, (SODA), San Diego, USA*. Ed. by Timothy M. Chan. SIAM, 2019, pp. 1445–1464.
- [12] Thomas Colcombet and A. V. Sreejith. “Limited Set Quantifiers over Countable Linear Orderings”. In: *International Colloquium on Automata, Languages and Programming, ICALP, Kyoto, Japan*. Ed. by Magnús M. Halldórsson et al. Vol. 9135. Lecture Notes in Computer Science. Springer, 2015, pp. 146–158.
- [13] Bruno Courcelle and Joost Engelfriet. *Graph Structure and Monadic Second-Order Logic - A Language-Theoretic Approach*. Vol. 138. Encyclopedia of Mathematics and Its Applications. Cambridge University Press, 2012.
- [14] Wojciech Czerwinski et al. “A Characterization for Decidable Separability by Piecewise Testable Languages”. In: *Discret. Math. Theor. Comput. Sci.* 19.4 (2017).
- [15] Reinhard Diestel. *Graph theory (electronic edition)*. Vol. 173. Graduate texts in mathematics. Springer-Verlag, 2006.
- [16] Samuel Eilenberg and Marcel-Paul Schützenberger. *On pseudovarieties*. IRIA. Laboratoire de Recherche en Informatique et Automatique, 1975.
- [17] Calvin C. Elgot. “Decision problems of finite automata design and related arithmetics”. In: *Trans. Amer. Math. Soc.* 98 (1961), pp. 21–51.
- [18] Ronald Fagin. “Generalized first-order spectra and polynomial-time recognizable sets”. In: *Complexity of computation (Proc. SIAM-AMS Sympos. Appl. Math., New York, 1973)*. Providence, R.I.: Amer. Math. Soc., 1974, 43–73. SIAM–AMS Proc., Vol. VII.
- [19] Gudmund Skovbjerg Frandsen, Peter Bro Miltersen, and Sven Skyum. “Dynamic Word Problems”. In: *J. ACM* 44.2 (Mar. 1997), pp. 257–271.
- [20] J. A. Green and D. Rees. “On semi-groups in which $x^r = x$ ”. In: *Mathematical Proceedings of the Cambridge Philosophical Society* 48.1 (1952), pp. 35–40.
- [21] Jörg Flum Heinz-Dieter Ebbinghaus. *Finite Model Theory*. 2nd. Springer Monographs in Mathematics. Springer, 2006.
- [22] J.A. Kamp. “Tense Logic and the Theory of Linear Order”. PhD thesis. Univ. of California, Los Angeles, 1968.
- [23] Manfred Kufleitner. “The Height of Factorization Forests”. In: *Mathematical Foundations of Computer Science 2008, 33rd International Symposium, MFCS 2008, Torun, Poland, August 25-29, 2008, Proceedings*. Ed. by Edward Ochmanski and Jerzy Tyszkiewicz. Vol. 5162. Lecture Notes in Computer Science. Springer, 2008, pp. 443–454.

- [24] H Läuchli and J Leonard. “On the elementary theory of linear order”. In: *Fundamenta Mathematicae* 59.1 (1966), pp. 109–116.
- [25] Robert McNaughton. “Testing and generating infinite sequences by a finite automaton”. In: *Information and Control* 9 (1966), pp. 521–530.
- [26] Robert McNaughton and Seymour Papert. *Counter-free automata*. The M.I.T. Press, Cambridge, Mass.-London, 1971.
- [27] J.-E. Pin and P. Weil. “Polynomial closure and unambiguous product”. In: *Theory Comput. Syst.* 30.4 (1997), pp. 383–422.
- [28] Thomas Place and Marc Zeitoun. “Going Higher in First-Order Quantifier Alternation Hierarchies on Words”. In: *J. ACM* 66.2 (2019), 12:1–12:65.
- [29] Frank D. Ramsey. “On a problem of formal logic”. In: *Proc. of the London Math. Soc.* 30 (1929), pp. 338–384.
- [30] Jan Reiterman. “The Birkhoff theorem for finite algebras”. In: *Algebra Universalis* 14.1 (1982), pp. 1–10.
- [31] Chloé Rispal and Olivier Carton. “Complementation of Rational Sets on Countable Scattered Linear Orderings”. In: *International Journal of Foundations of Computer Science* 16.04 (Aug. 2005), pp. 767–786.
- [32] Hanamantagouda P Sankappanavar and Stanley Burris. “A course in universal algebra”. In: *Graduate Texts Math* 78 (1981).
- [33] Marcel-Paul Schützenberger. “On finite monoids having only trivial subgroups”. In: *Information and Control* 8 (1965), pp. 190–194.
- [34] Marcel-Paul Schützenberger. “Sur Le Produit De Concatenation Non Ambigu”. In: *Semigroup Forum* 13 (1976), pp. 47–75.
- [35] Saharon Shelah. “The Monadic Theory of Order”. In: *Annals of Mathematics* (1975), pp. 379–419.
- [36] Imre Simon. “Factorization Forests of Finite Height”. In: *Theoretical Computer Science* 72.1 (1990), pp. 65–94.
- [37] Imre Simon. “Piecewise testable events”. In: *Automata Theory and Formal Languages*. Ed. by H. Brakhage. Berlin, Heidelberg: Springer Berlin Heidelberg, 1975, pp. 214–222. ISBN: 978-3-540-37923-2.
- [38] J. W. Thatcher and J. B. Wright. “Generalized Finite Automata Theory with an Application to a Decision Problem of Second-Order Logic”. In: *Mathematical systems theory* 2.1 (Mar. 1968), pp. 57–81.
- [39] Boris A. Trakhtenbrot. “The synthesis of logical nets whose operators are described in terms of one-place predicate calculus (Russian)”. In: *Dokl. Akad. Nauk SSSR* 118.4 (1958), pp. 646–649.
- [40] Thomas Wilke. “Classifying Discrete Temporal Properties”. In: *STACS 99*. Ed. by Christoph Meinel and Sophie Tison. Berlin, Heidelberg: Springer Berlin Heidelberg, 1999, pp. 32–46.

Author index

Subject index