

# Efficient Algorithms for Shortest Partial Seeds in Words

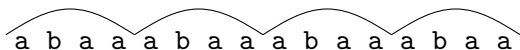
**Tomasz Kociumaka**<sup>1</sup>, Solon P. Pissis<sup>2</sup>,  
Jakub Radoszewski<sup>1</sup>, Wojciech Rytter<sup>1</sup>,  
Tomasz Waleń<sup>1</sup>

<sup>1</sup>University of Warsaw  
<sup>2</sup>King's College London

**CPM 2014**  
Moscow, June 16, 2014

# Periodicity and quasiperiodicity

Periodicity: occurrences are aligned

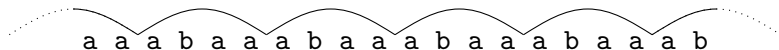


a b a a a b a a a b a a a b a a

The diagram shows the sequence "a b a a a b a a a b a a a b a a" with four wavy lines above it. Each wavy line spans the first four characters of one of the four "a b a a" blocks, demonstrating that the occurrences of this block are perfectly aligned and repeat every four characters.

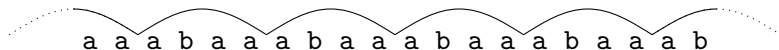
# Periodicity and quasiperiodicity

Periodicity: occurrences are aligned

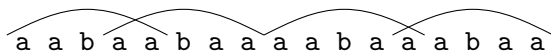


# Periodicity and quasiperiodicity

Periodicity: occurrences are aligned

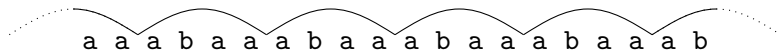


Quasiperiodicity: occurrences may overlap

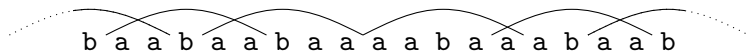


# Periodicity and quasiperiodicity

Periodicity: occurrences are aligned



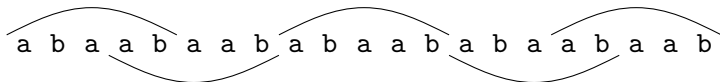
Quasiperiodicity: occurrences may overlap



# Covers and seeds

Definition (Apostolico, Farach, Iliopoulos; 1991)

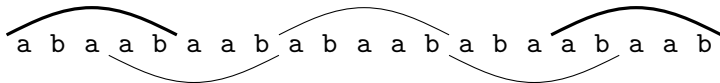
A factor  $u$  is a *cover* of  $w$  if each position (letter) in  $w$  lies within an occurrence of  $u$  in  $w$ .



# Covers and seeds

Definition (Apostolico, Farach, Iliopoulos; 1991)

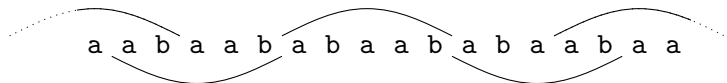
A factor  $u$  is a *cover* of  $w$  if each position (letter) in  $w$  lies within an occurrence of  $u$  in  $w$ .



# Covers and seeds

Definition (Iliopoulos, Moore, Park; 1993)

A factor  $u$  is a *seed* of  $w$  if  $u$  is a cover of a superstring of  $w$ .

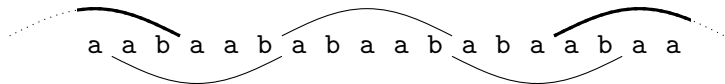




# Covers and seeds

Definition (Iliopoulos, Moore, Park; 1993)

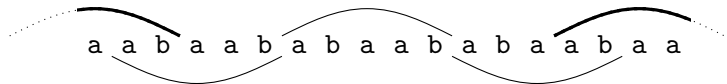
A factor  $u$  is a *seed* of  $w$  if  $u$  is a cover of a superstring of  $w$ .



# Covers and seeds

Definition (Iliopoulos, Moore, Park; 1993)

A factor  $u$  is a *seed* of  $w$  if  $u$  is a cover of a superstring of  $w$ .



Observation

A factor  $u$  is a seed of  $w$  iff each position (letter) in  $w$  lies within a possibly overhanging occurrence of  $u$  in  $w$ .

# Partial covers and partial seeds

## Definition (KPRRW; CPM'13)

The *cover index*  $\mathcal{C}(u)$  of  $u$  in  $w$  is the number of positions of  $w$  lying within an occurrence of  $u$  in  $w$ .

$\widehat{a} b \widehat{a a} b \widehat{a a a} b \widehat{a a a a} b \widehat{a a a} b \widehat{a} b \widehat{a a a} b$

$$\mathcal{C}(a) = 17$$

# Partial covers and partial seeds

## Definition (KPRRW; CPM'13)

The *cover index*  $\mathcal{C}(u)$  of  $u$  in  $w$  is the number of positions of  $w$  lying within an occurrence of  $u$  in  $w$ .

a b a a b a a a b a a a a b a a a b a b a a a b

$$\mathcal{C}(a) = 17$$

$$\mathcal{C}(abaa) = 19$$

# Partial covers and partial seeds

## Definition (KPRRW; CPM'13)

The *cover index*  $\mathcal{C}(u)$  of  $u$  in  $w$  is the number of positions of  $w$  lying within an occurrence of  $u$  in  $w$ .

a b a a b a a a b a a a b a a a b a b a a a b

$$\mathcal{C}(a) = 17$$

$$\mathcal{C}(abaa) = 19$$

$$\mathcal{C}(abaaaba) = 14$$

# Partial covers and partial seeds

## Definition (KPRRW; CPM'13)

The *cover index*  $\mathcal{C}(u)$  of  $u$  in  $w$  is the number of positions of  $w$  lying within an occurrence of  $u$  in  $w$ .

a b a a b a a a b a a a b a a a b a b a a a b

$$\mathcal{C}(a) = 17$$

$$\mathcal{C}(abaa) = 19$$

$$\mathcal{C}(abaaaba) = 14$$

## Definition

For a positive integer  $\alpha$  an  $\alpha$ -partial cover of  $w$  is a factor of  $w$  with cover index at least  $\alpha$ .

# Partial covers and partial seeds

## Definition

The *seed index*  $\mathcal{S}(u)$  of  $u$  in  $w$  is the number of positions of  $w$  lying within a possibly overhanging occurrence of  $u$  in  $w$ .

$\widehat{a} b \widehat{a a} b \widehat{a a a} b \widehat{a a a a} b \widehat{a a a} b \widehat{a} b \widehat{a a a} b$

$$\mathcal{C}(a) = 17$$

$$\mathcal{S}(a) = 17$$

$$\mathcal{C}(abaa) = 19$$

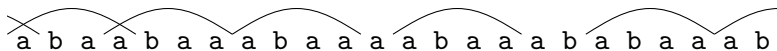
$$\mathcal{C}(abaaaba) = 14$$

# Partial covers and partial seeds

## Definition

The *seed index*  $\mathcal{S}(u)$  of  $u$  in  $w$  is the number of positions of  $w$  lying within a possibly overhanging occurrence of  $u$  in  $w$ .

a b a a b a a a b a a a a b a a a b a b a a a b



$$\mathcal{C}(a) = 17$$

$$\mathcal{C}(abaa) = 19$$

$$\mathcal{C}(abaaaba) = 14$$

$$\mathcal{S}(a) = 17$$

$$\mathcal{S}(abaa) = 21$$



# Partial covers and partial seeds

## Definition

The *seed index*  $\mathcal{S}(u)$  of  $u$  in  $w$  is the number of positions of  $w$  lying within a possibly overhanging occurrence of  $u$  in  $w$ .

a b a a b a a a b a a a a b a a a b a b a a a b

$$\mathcal{C}(a) = 17$$

$$\mathcal{S}(a) = 17$$

$$\mathcal{C}(abaa) = 19$$

$$\mathcal{S}(abaa) = 21$$

$$\mathcal{C}(abaaaba) = 14$$

$$\mathcal{S}(abaaaba) = 22$$

# Partial covers and partial seeds

## Definition

The *seed index*  $\mathcal{S}(u)$  of  $u$  in  $w$  is the number of positions of  $w$  lying within a possibly overhanging occurrence of  $u$  in  $w$ .

a b a a b a a a b a a a b a a a b a b a a a b

$$\mathcal{C}(a) = 17$$

$$\mathcal{S}(a) = 17$$

$$\mathcal{C}(abaa) = 19$$

$$\mathcal{S}(abaa) = 21$$

$$\mathcal{C}(abaaaba) = 14$$

$$\mathcal{S}(abaaaba) = 22$$

## Definition

For a positive integer  $\alpha$  an  $\alpha$ -partial seed of  $w$  is a factor of  $w$  with seed index at least  $\alpha$ .

# Other variants of covers and seeds

b a b a a a b a b a b a a a a b

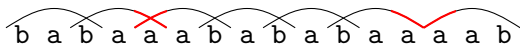
- *k*-covers and *k*-seeds (Iliopoulos, Smyth; 1998) – each position lies within a (possibly overhanging) occurrence of at least one of the few factors of length *k*, together forming a *k*-cover (*k*-seed).

# Other variants of covers and seeds

a a b a b a a a b a b a b a a a a b a

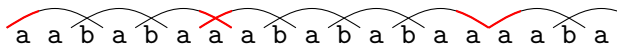
- *k*-covers and *k*-seeds (Iliopoulos, Smyth; 1998) – each position lies within a (possibly overhanging) occurrence of at least one of the few factors of length *k*, together forming a *k*-cover (*k*-seed).

# Other variants of covers and seeds



- *k*-covers and *k*-seeds (Iliopoulos, Smyth; 1998) – each position lies within a (possibly overhanging) occurrence of at least one of the few factors of length *k*, together forming a *k*-cover (*k*-seed).
- *approximate covers* (Sim, Park, Kim, Lee; 2002) and *approximate seeds* (Christodoulakis et al.; 2005) – each position is lies within a (possibly overhanging) occurrence of a factor *similar* to the approximate cover (or seed).

# Other variants of covers and seeds



- *k*-covers and *k*-seeds (Iliopoulos, Smyth; 1998) – each position lies within a (possibly overhanging) occurrence of at least one of the few factors of length *k*, together forming a *k*-cover (*k*-seed).
- *approximate covers* (Sim, Park, Kim, Lee; 2002) and *approximate seeds* (Christodoulakis et al.; 2005) – each position is lies within a (possibly overhanging) occurrence of a factor *similar* to the approximate cover (or seed).

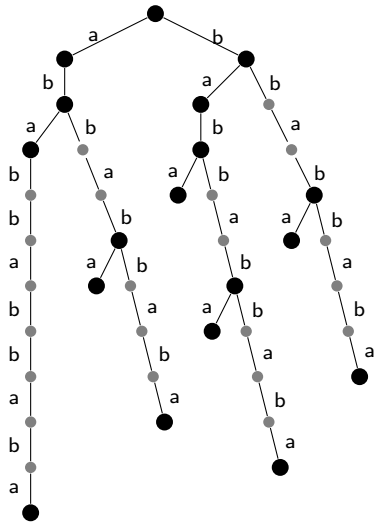
# Other variants of covers and seeds

b a b a a a b a b a b a a a a b

- *k*-covers and *k*-seeds (Iliopoulos, Smyth; 1998) – each position lies within a (possibly overhanging) occurrence of at least one of the few factors of length *k*, together forming a *k*-cover (*k*-seed).
- *approximate covers* (Sim, Park, Kim, Lee; 2002) and *approximate seeds* (Christodoulakis et al.; 2005) – each position is lies within a (possibly overhanging) occurrence of a factor *similar* to the approximate cover (or seed).

Main drawback:  $\Omega(n^2)$  algorithms.

# Cover Suffix Tree



$CST(ababbabbaba)$

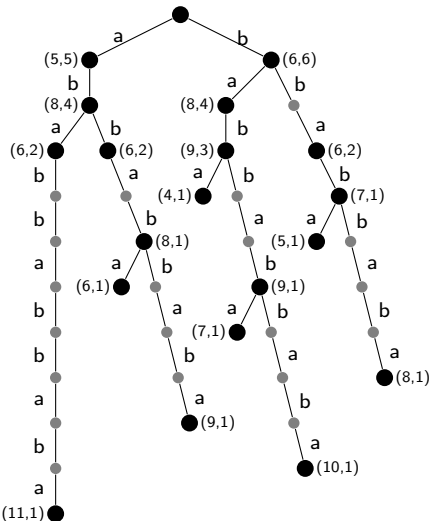
The cover suffix tree of  $w$  (denoted  $CST(w)$ ) is a suffix tree







# Cover Suffix Tree



$CST(ababbabbaba)$

The cover suffix tree of  $w$  (denoted  $CST(w)$ ) is a suffix tree

- augmented with  $\mathcal{O}(n)$  extra nodes,
- with each node *annotated* with a pair of integers  $(\mathcal{C}(v), \Delta(v))$ .

**Theorem (KPRRW; CPM'13)**

*The tree  $CST(w)$  can be built in  $\mathcal{O}(n \log n)$  time for any word  $w$  of length  $n$ .*

## Problem (PARTIAL SEEDS)

*Given a word  $w$  of length  $n$  and a positive integer  $\alpha \leq n$  find all shortest factors  $u$  of  $w$  such that  $\mathcal{S}(u) \geq \alpha$ .*

## Problem (LIMITED LENGTH PARTIAL SEEDS)

*Given a word  $w$  of length  $n$  and an interval  $[\ell, r]$  find a factor  $u$  of  $w$  maximizing  $\mathcal{S}(u)$  among factors for which  $|u| \in [\ell, r]$ .*

# Our results

## Problem (PARTIAL SEEDS)

*Given a word  $w$  of length  $n$  and a positive integer  $\alpha \leq n$  find all shortest factors  $u$  of  $w$  such that  $\mathcal{S}(u) \geq \alpha$ .*

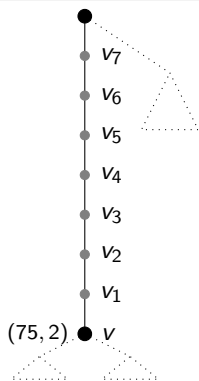
## Problem (LIMITED LENGTH PARTIAL SEEDS)

*Given a word  $w$  of length  $n$  and an interval  $[\ell, r]$  find a factor  $u$  of  $w$  maximizing  $\mathcal{S}(u)$  among factors for which  $|u| \in [\ell, r]$ .*

## Theorem

*Given  $CST(w)$  both PARTIAL SEEDS and LIMITED LENGTH PARTIAL SEEDS can be solved in linear time.*

# Determining the cover index

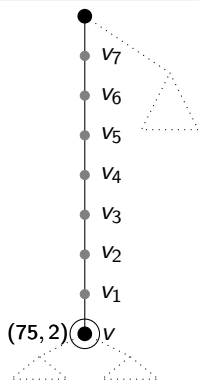


## Lemma (CPM'13)

Let  $v_0, v_1, \dots, v_k$  be the nodes of an edge of  $CST(w)$  with  $v = v_0$  being the lowest (explicit) node. Then

$$C(v_j) = C(v) - j\Delta(v).$$

# Determining the cover index



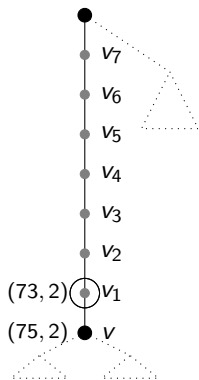
## Lemma (CPM'13)

Let  $v_0, v_1, \dots, v_k$  be the nodes of an edge of  $CST(w)$  with  $v = v_0$  being the lowest (explicit) node. Then

$$C(v_j) = C(v) - j\Delta(v).$$



# Determining the cover index



## Lemma (CPM'13)

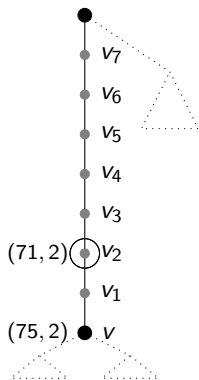
Let  $v_0, v_1, \dots, v_k$  be the nodes of an edge of  $CST(w)$  with  $v = v_0$  being the lowest (explicit) node. Then

$$\mathcal{C}(v_j) = \mathcal{C}(v) - j\Delta(v).$$





# Determining the cover index



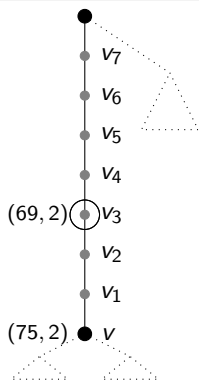
## Lemma (CPM'13)

Let  $v_0, v_1, \dots, v_k$  be the nodes of an edge of  $CST(w)$  with  $v = v_0$  being the lowest (explicit) node. Then

$$C(v_j) = C(v) - j\Delta(v).$$



# Determining the cover index



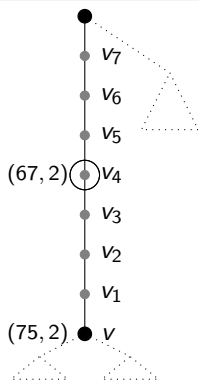
## Lemma (CPM'13)

Let  $v_0, v_1, \dots, v_k$  be the nodes of an edge of  $CST(w)$  with  $v = v_0$  being the lowest (explicit) node. Then

$$C(v_j) = C(v) - j\Delta(v).$$



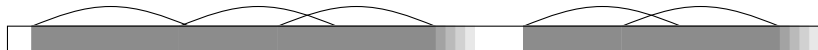
# Determining the cover index



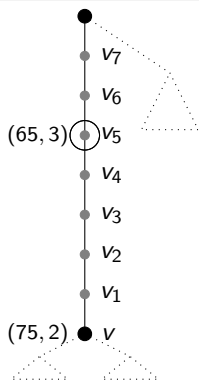
## Lemma (CPM'13)

Let  $v_0, v_1, \dots, v_k$  be the nodes of an edge of  $CST(w)$  with  $v = v_0$  being the lowest (explicit) node. Then

$$C(v_j) = C(v) - j\Delta(v).$$



# Determining the cover index



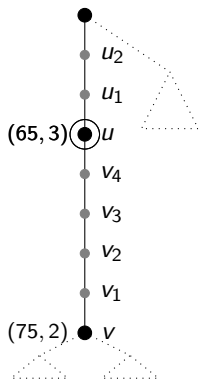
## Lemma (CPM'13)

Let  $v_0, v_1, \dots, v_k$  be the nodes of an edge of  $CST(w)$  with  $v = v_0$  being the lowest (explicit) node. Then

$$C(v_j) = C(v) - j\Delta(v).$$



# Determining the cover index



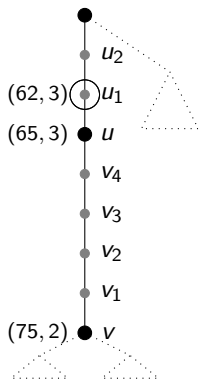
## Lemma (CPM'13)

Let  $v_0, v_1, \dots, v_k$  be the nodes of an edge of  $CST(w)$  with  $v = v_0$  being the lowest (explicit) node. Then

$$C(v_j) = C(v) - j\Delta(v).$$



# Determining the cover index



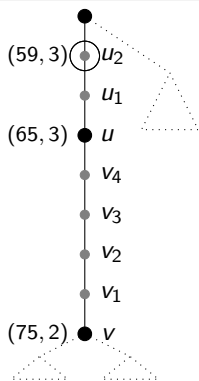
## Lemma (CPM'13)

Let  $v_0, v_1, \dots, v_k$  be the nodes of an edge of  $CST(w)$  with  $v = v_0$  being the lowest (explicit) node. Then

$$C(v_j) = C(v) - j\Delta(v).$$



# Determining the cover index



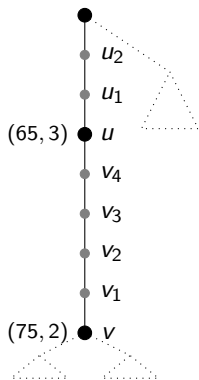
## Lemma (CPM'13)

Let  $v_0, v_1, \dots, v_k$  be the nodes of an edge of  $CST(w)$  with  $v = v_0$  being the lowest (explicit) node. Then

$$C(v_j) = C(v) - j\Delta(v).$$



# Determining the cover index



## Lemma (CPM'13)

Let  $v_0, v_1, \dots, v_k$  be the nodes of an edge of  $CST(w)$  with  $v = v_0$  being the lowest (explicit) node. Then

$$\mathcal{C}(v_j) = \mathcal{C}(v) - j\Delta(v).$$



## Corollary

Given a locus of  $v$  in  $CST(w)$ , the cover index  $\mathcal{C}(v)$  can be computed in  $\mathcal{O}(1)$  time.



# Seed index

$$\mathcal{S}(v) = \quad + \quad +$$



# Seed index

$$S(v) = \quad + \mathcal{C}(v) +$$

↑  
Full occs



# Seed index

$$S(v) = \underset{\substack{\nearrow \\ \text{Left-overhanging} \\ \text{occs only}}}{\text{Left}S(v)} + \underset{\substack{\uparrow \\ \text{Full occs}}}{C(v)} +$$



# Seed index

$$S(v) = \underset{\substack{\text{Left-overhanging} \\ \text{occs only}}}{\text{Left}S(v)} + \underset{\substack{\text{Full} \\ \text{occs}}}{C(v)} + \underset{\substack{\text{Right-overhanging} \\ \text{occs only}}}{\text{Right}S(v)}$$



# Seed index

$$S(v) = \underset{\substack{\text{Left-overhanging} \\ \text{occs only}}}{\text{LeftS}(v)} + \underset{\substack{\text{Full} \\ \text{occs}}}{C(v)} + \underset{\substack{\text{Right-overhanging} \\ \text{occs only}}}{\text{RightS}(v)}$$

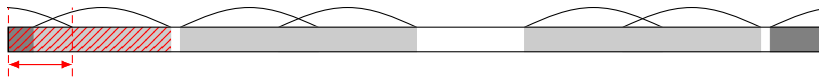


$$\text{LeftS}(v) = \min(B[\text{first}(v) + |v| - 1], \text{first}(v) - 1)$$

$\text{first}(v)$  start position of the first occurrence of  $v$ ,  
 $B[i]$  largest border of  $w[1..i]$ .

# Seed index

$$S(v) = \underset{\substack{\text{Left-overhanging} \\ \text{occs only}}}{\text{LeftS}(v)} + \underset{\substack{\text{Full occs}}}{C(v)} + \underset{\substack{\text{Right-overhanging} \\ \text{occs only}}}{\text{RightS}(v)}$$

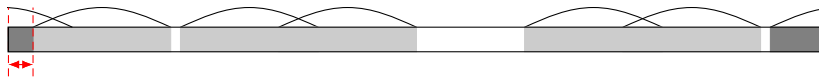


$$\text{LeftS}(v) = \min(B[\text{first}(v) + |v| - 1], \text{first}(v) - 1)$$

$\text{first}(v)$  start position of the first occurrence of  $v$ ,  
 $B[i]$  largest border of  $w[1..i]$ .

# Seed index

$$S(v) = \underset{\substack{\text{Left-overhanging} \\ \text{occs only}}}{\text{Left}S(v)} + \underset{\text{Full occs}}{C(v)} + \underset{\substack{\text{Right-overhanging} \\ \text{occs only}}}{\text{Right}S(v)}$$



$$\text{Left}S(v) = \min(B[\text{first}(v) + |v| - 1], \text{first}(v) - 1)$$

$\text{first}(v)$  start position of the first occurrence of  $v$ ,  
 $B[i]$  largest border of  $w[1..i]$ .

# Seed index

$$S(v) = \underset{\substack{\text{Left-overhanging} \\ \text{occs only}}}{\text{Left}S(v)} + \underset{\substack{\text{Full} \\ \text{occs}}}{C(v)} + \underset{\substack{\text{Right-overhanging} \\ \text{occs only}}}{\text{Right}S(v)}$$



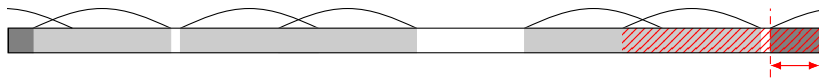
$$\text{Left}S(v) = \min(B[\text{first}(v) + |v| - 1], \text{first}(v) - 1)$$
$$\text{Right}S(v) = \min(B^R[\text{last}(v)], n - |v| + 1 - \text{last}(v))$$

$\text{last}(v)$  start position of the last occurrence of  $v$ ,  
 $B^R[i]$  largest border of  $w[i..n]$ .



# Seed index

$$S(v) = \underset{\substack{\text{Left-overhanging} \\ \text{occs only}}}{\text{LeftS}(v)} + \underset{\substack{\text{Full occs}}}{C(v)} + \underset{\substack{\text{Right-overhanging} \\ \text{occs only}}}{\text{RightS}(v)}$$

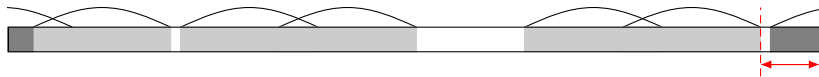


$$\text{LeftS}(v) = \min(B[\text{first}(v) + |v| - 1], \text{first}(v) - 1)$$
$$\text{RightS}(v) = \min(B^R[\text{last}(v)], n - |v| + 1 - \text{last}(v))$$

$\text{last}(v)$  start position of the last occurrence of  $v$ ,  
 $B^R[i]$  largest border of  $w[i..n]$ .

# Seed index

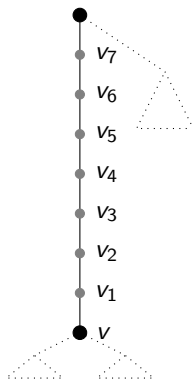
$$S(v) = \underset{\substack{\text{Left-overhanging} \\ \text{occs only}}}{\text{LeftS}(v)} + \underset{\substack{\text{Full occs}}}{C(v)} + \underset{\substack{\text{Right-overhanging} \\ \text{occs only}}}{\text{RightS}(v)}$$



$$\text{LeftS}(v) = \min(B[\text{first}(v) + |v| - 1], \text{first}(v) - 1)$$
$$\text{RightS}(v) = \min(B^R[\text{last}(v)], n - |v| + 1 - \text{last}(v))$$

$\text{last}(v)$  start position of the last occurrence of  $v$ ,  
 $B^R[i]$  largest border of  $w[i..n]$ .

# Seed index on an edge of $CST(w)$



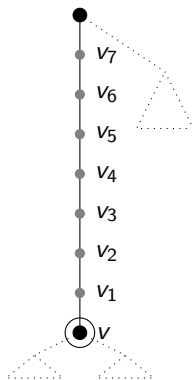
$$S(v_j) = LeftS(v_j) + C(v_j) + RightS(v_j)$$

$C(v_j)$  a linear function,

$RightS(v_j)$  becomes a linear upon creation of  $\leq 1$  extra node per edge,

$LeftS(v_j)$  less regular:  $\min(B[r_v - j], c_v)$ .

# Seed index on an edge of $CST(w)$

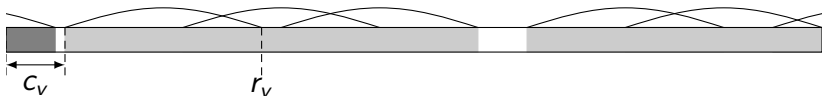


$$S(v_j) = LeftS(v_j) + C(v_j) + RightS(v_j)$$

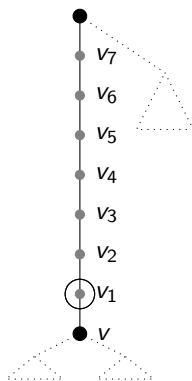
$C(v_j)$  a linear function,

$RightS(v_j)$  becomes a linear upon creation of  $\leq 1$  extra node per edge,

$LeftS(v_j)$  less regular:  $\min(B[r_v - j], c_v)$ .



# Seed index on an edge of $CST(w)$

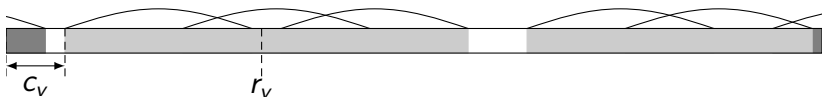


$$S(v_j) = LeftS(v_j) + C(v_j) + RightS(v_j)$$

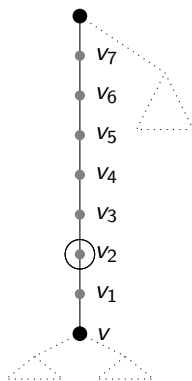
$C(v_j)$  a linear function,

$RightS(v_j)$  becomes a linear upon creation of  $\leq 1$  extra node per edge,

$LeftS(v_j)$  less regular:  $\min(B[r_v - j], c_v)$ .



# Seed index on an edge of $CST(w)$

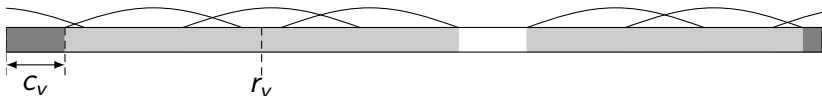


$$S(v_j) = LeftS(v_j) + C(v_j) + RightS(v_j)$$

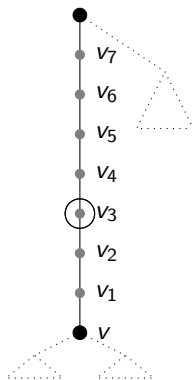
$C(v_j)$  a linear function,

$RightS(v_j)$  becomes a linear upon creation of  $\leq 1$  extra node per edge,

$LeftS(v_j)$  less regular:  $\min(B[r_v - j], c_v)$ .



# Seed index on an edge of $CST(w)$

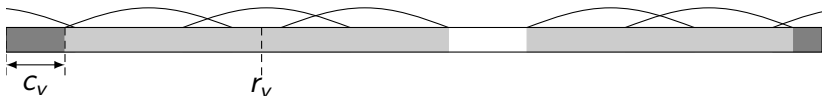


$$S(v_j) = LeftS(v_j) + C(v_j) + RightS(v_j)$$

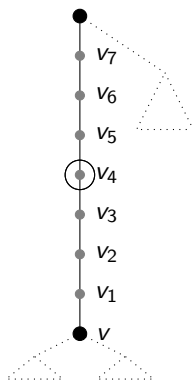
$C(v_j)$  a linear function,

$RightS(v_j)$  becomes a linear upon creation of  $\leq 1$  extra node per edge,

$LeftS(v_j)$  less regular:  $\min(B[r_v - j], c_v)$ .



# Seed index on an edge of $CST(w)$

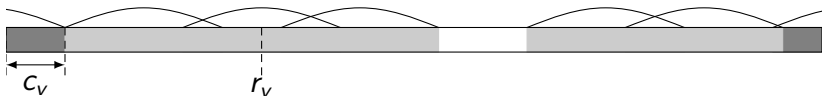


$$S(v_j) = LeftS(v_j) + C(v_j) + RightS(v_j)$$

$C(v_j)$  a linear function,

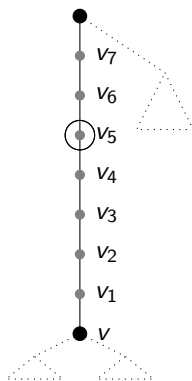
$RightS(v_j)$  becomes a linear upon creation of  $\leq 1$  extra node per edge,

$LeftS(v_j)$  less regular:  $\min(B[r_v - j], c_v)$ .





# Seed index on an edge of $CST(w)$

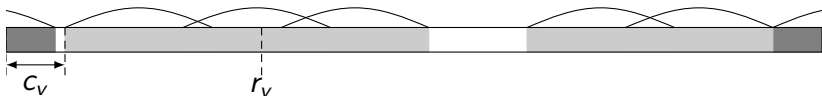


$$S(v_j) = LeftS(v_j) + C(v_j) + RightS(v_j)$$

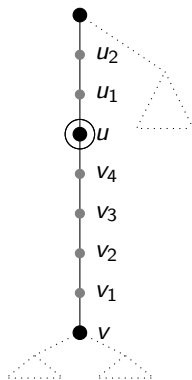
$C(v_j)$  a linear function,

$RightS(v_j)$  becomes a linear upon creation of  $\leq 1$  extra node per edge,

$LeftS(v_j)$  less regular:  $\min(B[r_v - j], c_v)$ .



# Seed index on an edge of $CST(w)$

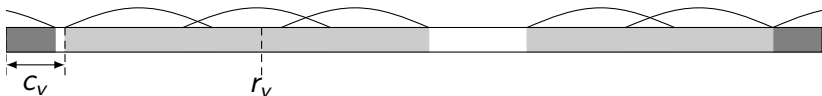


$$S(v_j) = LeftS(v_j) + C(v_j) + RightS(v_j)$$

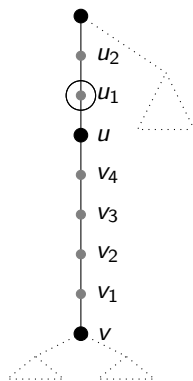
$C(v_j)$  a linear function,

$RightS(v_j)$  becomes a linear upon creation of  $\leq 1$  extra node per edge,

$LeftS(v_j)$  less regular:  $\min(B[r_v - j], c_v)$ .



# Seed index on an edge of $CST(w)$

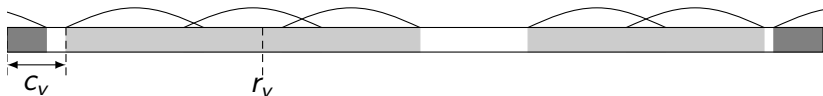


$$S(v_j) = LeftS(v_j) + C(v_j) + RightS(v_j)$$

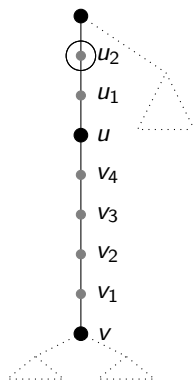
$C(v_j)$  a linear function,

$RightS(v_j)$  becomes a linear upon creation of  $\leq 1$  extra node per edge,

$LeftS(v_j)$  less regular:  $\min(B[r_v - j], c_v)$ .



# Seed index on an edge of $CST(w)$

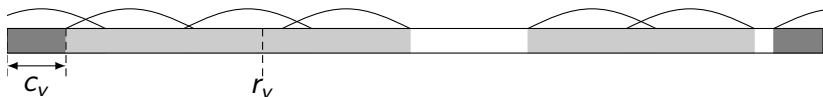


$$S(v_j) = LeftS(v_j) + C(v_j) + RightS(v_j)$$

$C(v_j)$  a linear function,

$RightS(v_j)$  becomes a linear upon creation of  $\leq 1$  extra node per edge,

$LeftS(v_j)$  less regular:  $\min(B[r_v - j], c_v)$ .



# Seed Suffix Tree

$CST(w)$  can be further augmented in  $\mathcal{O}(n)$  time to  $SST(w)$  (*Seed Suffix Tree*) such that

- for each node  $v$  there exists a function

$$\phi_v(x) = a_v x + b_v + \min(c_v, B[x])$$

and a range  $R_v = (\ell_v, r_v]$  such that  $\mathcal{S}(v_j) = \phi_v(r_v - j)$  for any  $v_j$  on the edge immediately above  $v$ ,

# Seed Suffix Tree

$CST(w)$  can be further augmented in  $\mathcal{O}(n)$  time to  $SST(w)$  (*Seed Suffix Tree*) such that

- for each node  $v$  there exists a function

$$\phi_v(x) = a_v x + b_v + \min(c_v, B[x])$$

and a range  $R_v = (\ell_v, r_v]$  such that  $\mathcal{S}(v_j) = \phi_v(r_v - j)$  for any  $v_j$  on the edge immediately above  $v$ ,

- $0 \leq a_v \leq \text{Occ}(v)$ , where  $\text{Occ}(v)$  is the number of occurrences of  $v$  in  $w$ .

# Seed Suffix Tree

$CST(w)$  can be further augmented in  $\mathcal{O}(n)$  time to  $SST(w)$  (*Seed Suffix Tree*) such that

- for each node  $v$  there exists a function

$$\phi_v(x) = a_v x + b_v + \min(c_v, B[x])$$

and a range  $R_v = (\ell_v, r_v]$  such that  $\mathcal{S}(v_j) = \phi_v(r_v - j)$  for any  $v_j$  on the edge immediately above  $v$ ,

- $0 \leq a_v \leq \text{Occ}(v)$ , where  $\text{Occ}(v)$  is the number of occurrences of  $v$  in  $w$ .

## Observation

*Given a locus of  $v$  in  $SST(w)$  and the border table  $B$ , the seed index  $\mathcal{S}(v)$  can be computed in  $\mathcal{O}(1)$  time.*

# Abstract problems

## Problem

**Input:** pairs  $(\phi_i, R_i)$ , where  $\phi_i(x) = a_i x + b_i + \min(c_i, B[x])$  is a function and  $R_i = (\ell_i, r_i] \subseteq [1, n]$  is a non-empty range

**Output:**

- (a)  $\operatorname{argmax}\{\phi_i(x) : x \in R_i\}$  for each pair,
- (b)  $\min\{x \in R_i : \phi_i(x) \geq \alpha\}$  for each pair.



# Abstract problems

## Problem

**Input:** pairs  $(\phi_i, R_i)$ , where  $\phi_i(x) = a_i x + b_i + \min(c_i, B[x])$  is a function and  $R_i = (\ell_i, r_i] \subseteq [1, n]$  is a non-empty range

**Output:**

- (a)  $\operatorname{argmax}\{\phi_i(x) : x \in R_i\}$  for each pair,
- (b)  $\min\{x \in R_i : \phi_i(x) \geq \alpha\}$  for each pair.

## Lemma

Values (a) and (b) can be computed (offline) in linear time.

# Abstract problems

## Problem

**Input:** pairs  $(\phi_i, R_i)$ , where  $\phi_i(x) = a_i x + b_i + \min(c_i, B[x])$  is a function and  $R_i = (\ell_i, r_i] \subseteq [1, n]$  is a non-empty range

**Output:**

- (a)  $\operatorname{argmax}\{\phi_i(x) : x \in R_i\}$  for each pair,
- (b)  $\min\{x \in R_i : \phi_i(x) \geq \alpha\}$  for each pair.

## Lemma

Values (a) and (b) can be computed (offline) in linear time.  
Additional assumption required for (b):  $\sum a_i = \mathcal{O}(n)$ .

# Abstract problems

## Problem

**Input:** pairs  $(\phi_i, R_i)$ , where  $\phi_i(x) = a_i x + b_i + \min(c_i, B[x])$  is a function and  $R_i = (\ell_i, r_i] \subseteq [1, n]$  is a non-empty range

**Output:**

- (a)  $\operatorname{argmax}\{\phi_i(x) : x \in R_i\}$  for each pair,
- (b)  $\min\{x \in R_i : \phi_i(x) \geq \alpha\}$  for each pair.

## Lemma

Values (a) and (b) can be computed (offline) in linear time.  
Additional assumption required for (b):  $\sum a_i = \mathcal{O}(n)$ .

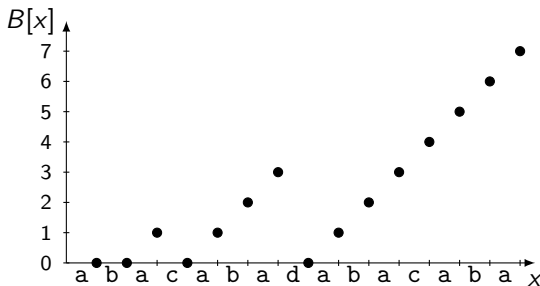
Workaround for  $\sum a_i = \mathcal{O}(n)$ :

- use (a) queries to restrict the set of edges queried for (b).

# Toy problem

## Problem

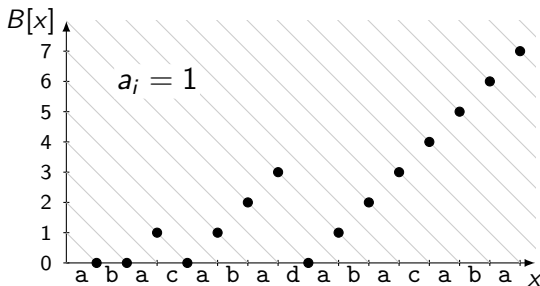
For the border array  $B$  to answer (off-line) the following queries: given a non-negative coefficient  $a_i$  and a range  $R_i = (\ell_i, r_i]$  compute  $x_i = \operatorname{argmax}\{a_i x + B[x] : x \in R_i\}$ .



# Toy problem

## Problem

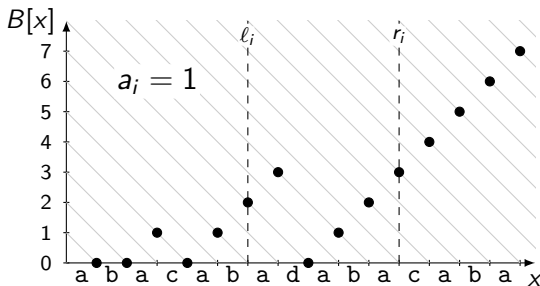
For the border array  $B$  to answer (off-line) the following queries: given a non-negative coefficient  $a_i$  and a range  $R_i = (\ell_i, r_i]$  compute  $x_i = \operatorname{argmax}\{a_i x + B[x] : x \in R_i\}$ .



# Toy problem

## Problem

For the border array  $B$  to answer (off-line) the following queries: given a non-negative coefficient  $a_i$  and a range  $R_i = (\ell_i, r_i]$  compute  $x_i = \operatorname{argmax}\{a_i x + B[x] : x \in R_i\}$ .



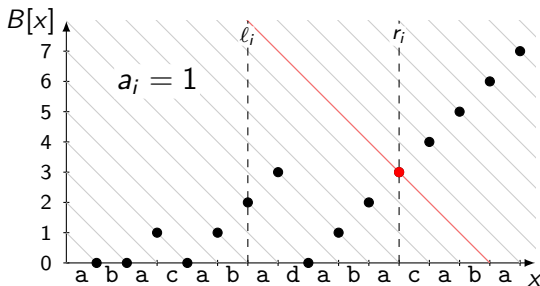
## Observation

For each query we have  $x_i = r_i$  or  $B[x_i + 1] < B[x_i] - a_i$ .

# Toy problem

## Problem

For the border array  $B$  to answer (off-line) the following queries: given a non-negative coefficient  $a_i$  and a range  $R_i = (\ell_i, r_i]$  compute  $x_i = \operatorname{argmax}\{a_i x + B[x] : x \in R_i\}$ .



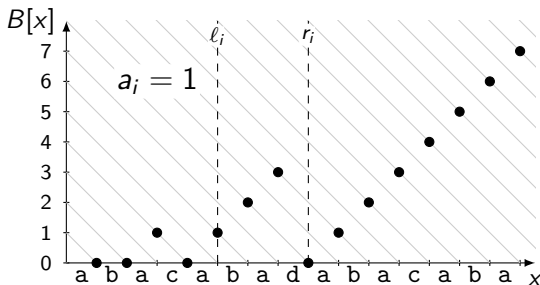
## Observation

For each query we have  $x_i = r_i$  or  $B[x_i + 1] < B[x_i] - a_i$ .

# Toy problem

## Problem

For the border array  $B$  to answer (off-line) the following queries: given a non-negative coefficient  $a_i$  and a range  $R_i = (\ell_i, r_i]$  compute  $x_i = \operatorname{argmax}\{a_i x + B[x] : x \in R_i\}$ .



## Observation

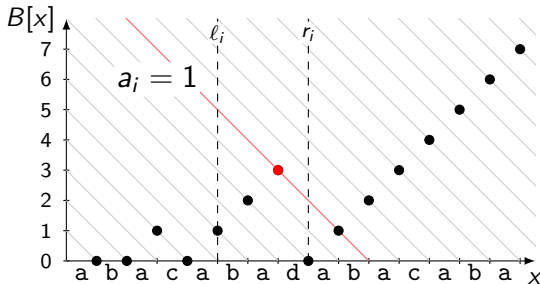
For each query we have  $x_i = r_i$  or  $B[x_i + 1] < B[x_i] - a_i$ .



# Toy problem

## Problem

For the border array  $B$  to answer (off-line) the following queries: given a non-negative coefficient  $a_i$  and a range  $R_i = (\ell_i, r_i]$  compute  $x_i = \operatorname{argmax}\{a_i x + B[x] : x \in R_i\}$ .



## Observation

For each query we have  $x_i = r_i$  or  $B[x_i + 1] < B[x_i] - a_i$ .

# Toy problem: solution

## Observation

Let  $F_a = \{x : B[x + 1] < B[x] - a\}$ . Then  $\sum_{a \geq 0} |F_a| = \mathcal{O}(n)$ .

# Toy problem: solution

## Observation

Let  $F_a = \{x : B[x+1] < B[x] - a\}$ . Then  $\sum_{a \geq 0} |F_a| = \mathcal{O}(n)$ .

## Proof.

$B[x+1] \leq B[x] + 1$ , i.e. the total increase in  $B$  is at most  $n$ .  
 $\sum_{a \geq 0} |F_a|$  is bounded by the total decrease of in  $B$ .  $\square$

# Toy problem: solution

## Observation

Let  $F_a = \{x : B[x + 1] < B[x] - a\}$ . Then  $\sum_{a \geq 0} |F_a| = \mathcal{O}(n)$ .

## Proof.

$B[x + 1] \leq B[x] + 1$ , i.e. the total increase in  $B$  is at most  $n$ .  
 $\sum_{a \geq 0} |F_a|$  is bounded by the total decrease of in  $B$ .  $\square$

- 1 Apply (offline) predecessor queries to translate the range  $R_i$  into the range of positions in  $F_{a_i}$ .

# Toy problem: solution

## Observation

Let  $F_a = \{x : B[x+1] < B[x] - a\}$ . Then  $\sum_{a \geq 0} |F_a| = \mathcal{O}(n)$ .

## Proof.

$B[x+1] \leq B[x] + 1$ , i.e. the total increase in  $B$  is at most  $n$ .  
 $\sum_{a \geq 0} |F_a|$  is bounded by the total decrease of in  $B$ .  $\square$

- 1 Apply (offline) predecessor queries to translate the range  $R_i$  into the range of positions in  $F_{a_i}$ .
- 2 Use range maximum queries (RMQ) for  $a_i x + B[x]$  and  $x \in F_{a_i}$  to compute  $\operatorname{argmax}\{a_i x + B[x] : x \in R_i \cap F_{a_i}\}$ .

# Toy problem: solution

## Observation

Let  $F_a = \{x : B[x + 1] < B[x] - a\}$ . Then  $\sum_{a \geq 0} |F_a| = \mathcal{O}(n)$ .

## Proof.

$B[x + 1] \leq B[x] + 1$ , i.e. the total increase in  $B$  is at most  $n$ .  
 $\sum_{a \geq 0} |F_a|$  is bounded by the total decrease of in  $B$ .  $\square$

- 1 Apply (offline) predecessor queries to translate the range  $R_i$  into the range of positions in  $F_{a_i}$ .
- 2 Use range maximum queries (RMQ) for  $a_i x + B[x]$  and  $x \in F_{a_i}$  to compute  $\operatorname{argmax}\{a_i x + B[x] : x \in R_i \cap F_{a_i}\}$ .
- 3 For each query check the possibility of  $x_i = r_i$ .

# Conclusions and open problems

Two problems regarding partial seeds can be solved in  $\mathcal{O}(n)$  time provided that  $CST(w)$  is already computed:

# Conclusions and open problems

Two problems regarding partial seeds can be solved in  $\mathcal{O}(n)$  time provided that  $CST(w)$  is already computed:

- find the shortest factor  $u$  with  $\mathcal{S}(u)$  exceeding a given threshold  $\alpha$ ,



# Conclusions and open problems

Two problems regarding partial seeds can be solved in  $\mathcal{O}(n)$  time provided that  $CST(w)$  is already computed:

- find the shortest factor  $u$  with  $\mathcal{S}(u)$  exceeding a given threshold  $\alpha$ ,
- find the factor  $u$  maximizing  $\mathcal{S}(u)$  and satisfying length restrictions,
  - other kinds of restrictions also possible.

# Conclusions and open problems

Two problems regarding partial seeds can be solved in  $\mathcal{O}(n)$  time provided that  $CST(w)$  is already computed:

- find the shortest factor  $u$  with  $\mathcal{S}(u)$  exceeding a given threshold  $\alpha$ ,
- find the factor  $u$  maximizing  $\mathcal{S}(u)$  and satisfying length restrictions,
  - other kinds of restrictions also possible.

Open problems:

- improve the construction algorithm for  $CST(w)$  (currently  $\mathcal{O}(w \log n)$ ),

# Conclusions and open problems

Two problems regarding partial seeds can be solved in  $\mathcal{O}(n)$  time provided that  $CST(w)$  is already computed:

- find the shortest factor  $u$  with  $\mathcal{S}(u)$  exceeding a given threshold  $\alpha$ ,
- find the factor  $u$  maximizing  $\mathcal{S}(u)$  and satisfying length restrictions,
  - other kinds of restrictions also possible.

Open problems:

- improve the construction algorithm for  $CST(w)$  (currently  $\mathcal{O}(w \log n)$ ),
- for each length find the factor  $u$  maximizing  $\mathcal{S}(u)$ 
  - for partial covers  $\mathcal{O}(n \log n)$  time.

Thank you for your attention!