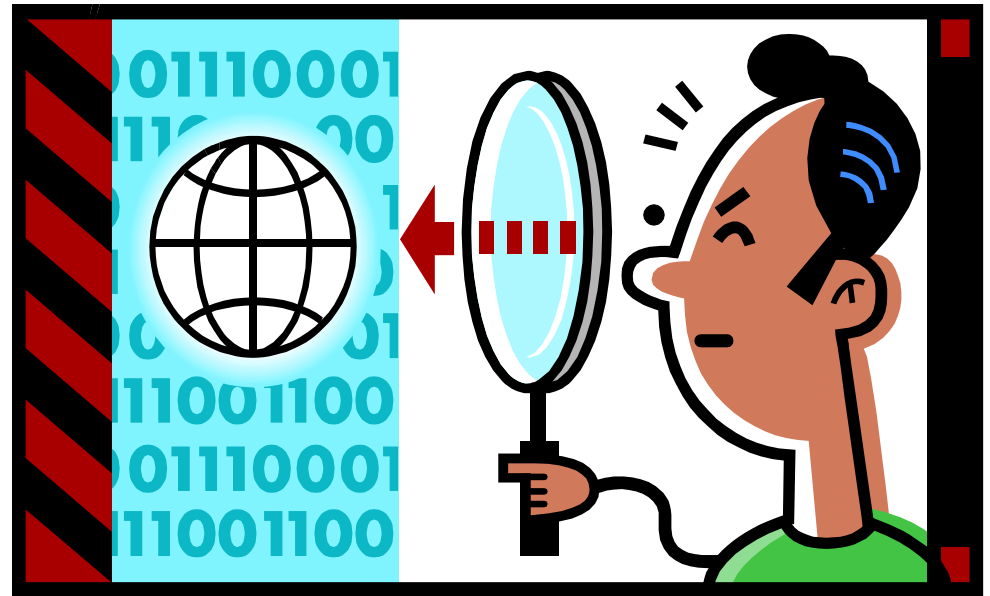


Wyszukiwanie informacji w internecie



Nguyen Hung Son

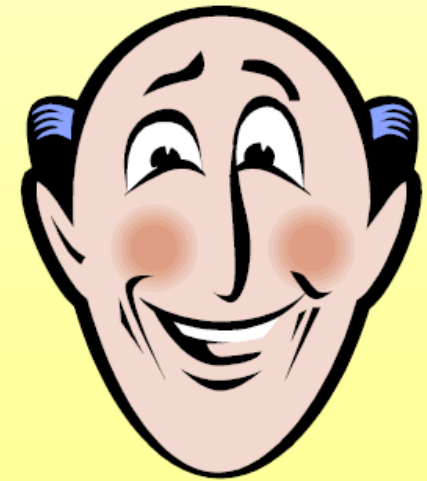
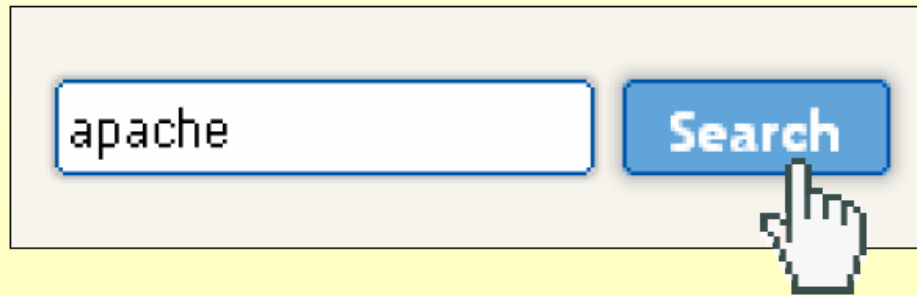
Jak znaleźć informację w internecie?

- Wyszukiwarki internetowe:
 - Potężne maszyny wykorzystujące najnowsze metody z różnych dziedzin
- Architektura: trzy główne moduły
 - Zarządzanie pająków;
 - Serwer indeksowania;
 - Interfejs użytkownika
- Wyniki wyszukiwania:
 - Lista rankingowa

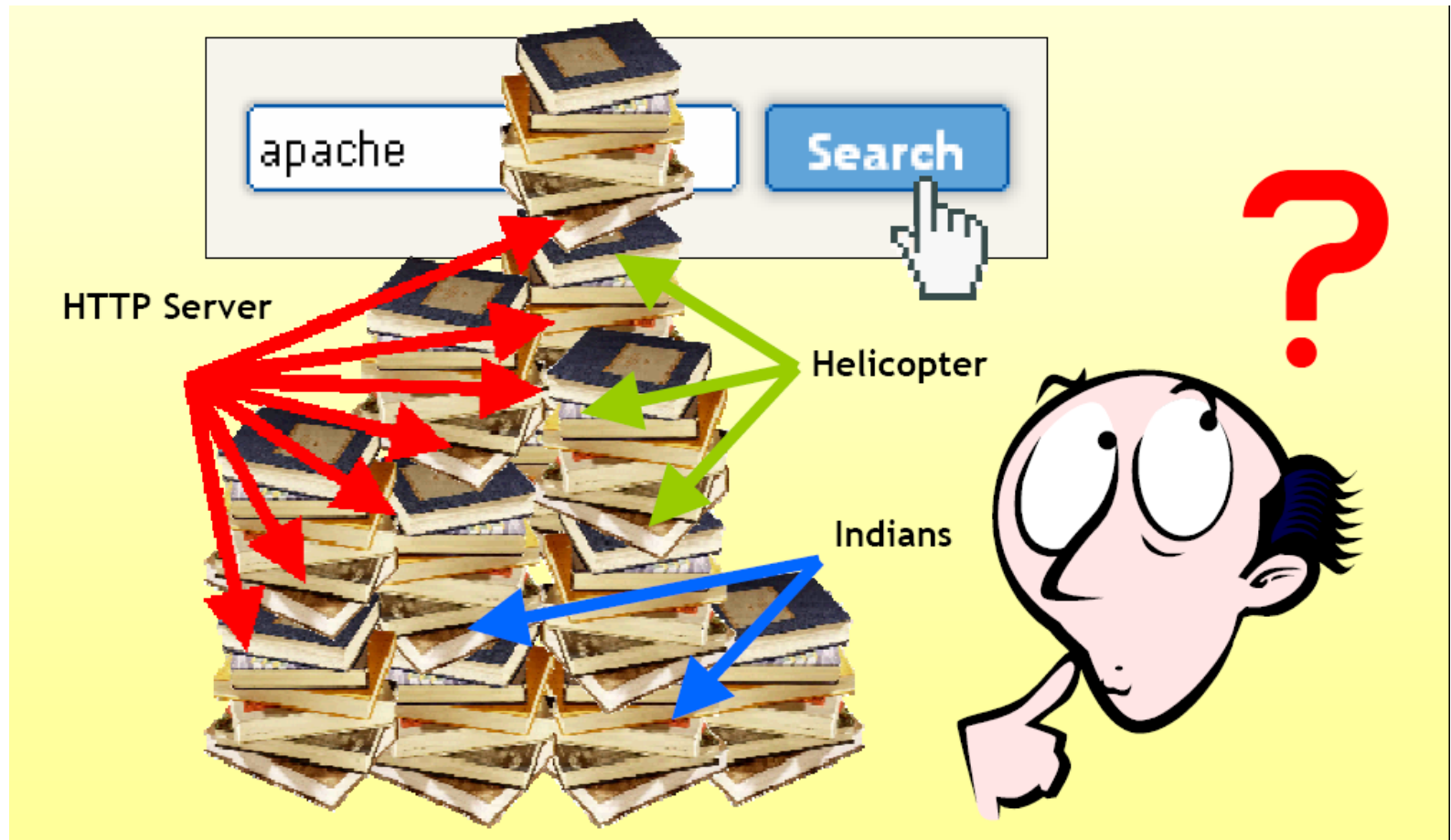
Architektura



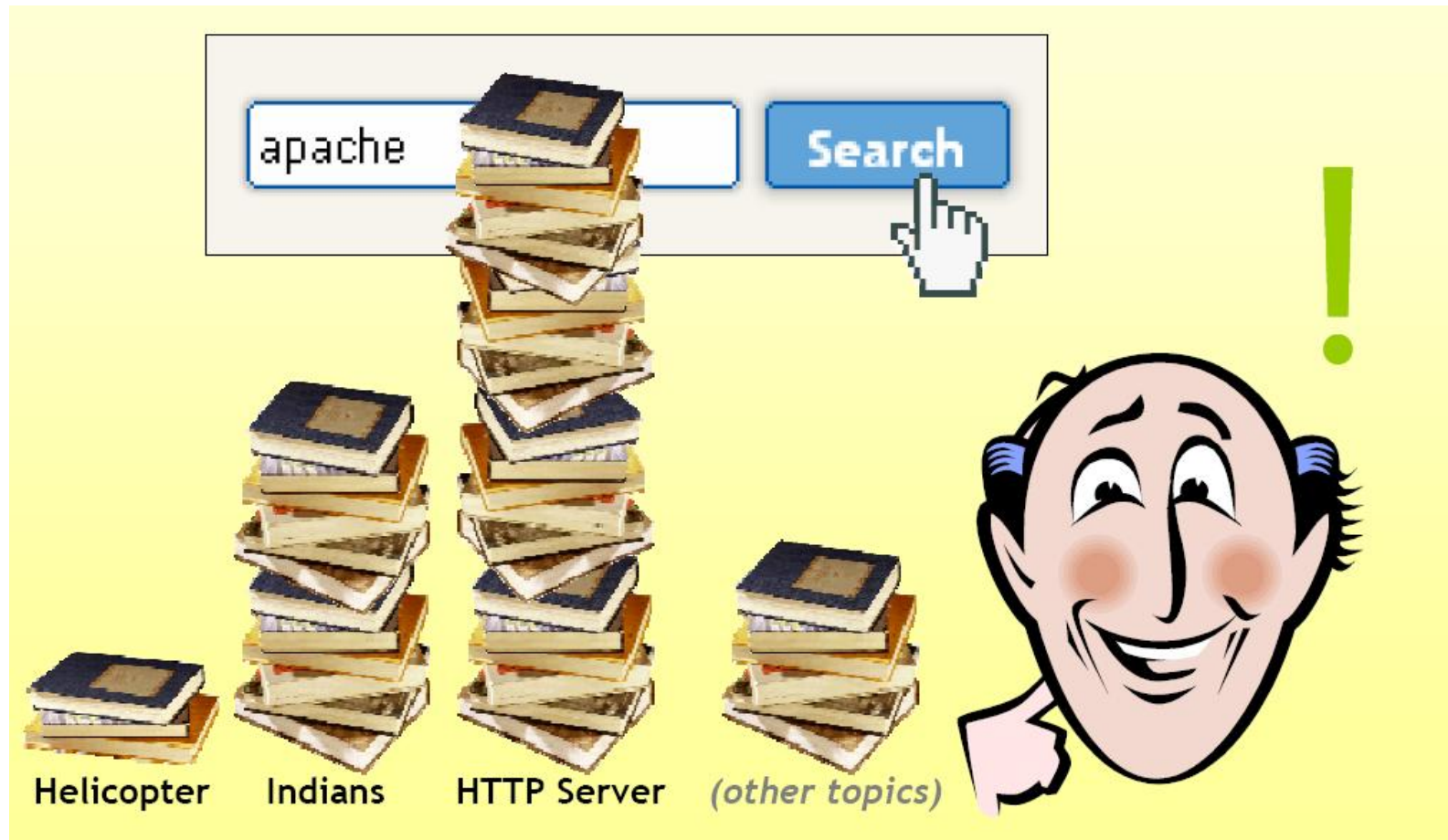
Lista rankingowa nie jest doskonała!



Lista rankingowa nie jest doskonała!



Grupowanie wyników wyszukiwania (ang. *SRC: Search Result Clustering*)



SRC korzysta z krótkich fragmentów tekstu (snippets)



[company](#) | [products](#) | [solutions](#) | [demos](#) | [partners](#) | [press](#)

Search the Web

Search

▶ [Advanced Search](#) ▶ [Help!](#) ▶ [Tell Us What You Think!](#)

Clustered Results

[jaguar](#) (194)

▶ [Jaguar Cars](#) (25)

▶ [Club](#) (15)

▶ [Parts, Auto](#) (16)

▶ [Cat](#) (14)

▶ [Mac](#) (12)

▶ [Type](#) (10)

▶ [Performance](#) (6)

▶ [Classic](#) (6)

▶ [Quote, Dealer](#) (17)

▶ [Panthera onca](#) (8)

▼ [More](#)

Find in clusters:



Top 194 results retrieved for the query **jaguar** ([Details](#))

New! Results now open in the full browser window by default. Click on the [\[frame\]](#) links next to the titles to get the old behavior and an updated toolbar with exciting new features.

[Apple Mac OS X 10.2 Jaguar](#) [\[new window\]](#) [\[frame\]](#) [\[preview\]](#)

Sponsored Link

Find great prices on Apple Mac OS X 10.2 **Jaguar** at CNET Shopper.com, a comprehensive pricing guide that will help you find the latest tech products at great prices. - [shopper.cnet.com](#) - [show in clusters](#)

[Get a Free Jaguar Quote](#) [\[new window\]](#) [\[frame\]](#) [\[preview\]](#)

Sponsored Link

Get a free **Jaguar** quote from a local dealer with Yahoo! Autos. Choose a vehicle, enter your contact info and a local dealer will contact you with a great no-haggle price. - [autos.yahoo.com](#) - [show in clusters](#)

1. [Jaguar Cars](#) [\[new window\]](#) [\[frame\]](#) [\[preview\]](#)

URL: [www.jaguarcars.com](#) - [show in clusters](#)

Sources: [Lycos 1](#), [Lycos 4](#), [Looksmart 2](#), [MSN 1](#)

2. [www.jaguar-racing.com](#) [\[new window\]](#) [\[frame\]](#) [\[preview\]](#)

URL: [www.jaguar-racing.com](#) - [show in clusters](#)

Sources: [Lycos 2](#), [Lycos 11](#), [Looksmart 35](#), [MSN 3](#)

3. [Apple - Mac OS X](#) [\[new window\]](#) [\[frame\]](#) [\[preview\]](#)

Learn about the new OS X Server, designed for the Internet, digital media and workgroup management. Download a technical factsheet. ... Mac OS X version 10.2 **Jaguar** contains over 150 new features and provides significant enhancements to its modern, UNIX-based ...

URL: [www.apple.com/macosx](#) - [show in clusters](#)

Sources: [MSN 2](#), [Lycos 3](#)

SRC czy grupowanie dokumentów?

Grupowanie dokumentów:

- Miliardy stron;
- Ich treści ciągle się zmieniają;
- Skalowalność wzg. liczby dokumentów
- Są to niestrukturalne i różnorodne dane;
- dodatkowe informacje:
 - hiperłącze,
 - przejścia między stronami (click-through data), itp.

SRC

- Próba 100~400 wyników wyszukiwania
- Informacje są aktualne
- Działa na bieżąco
- Skalowalność wzg. potrzeby użytkownika
- Zbyt mała, zaszumiona informacja → gorsza jakość grup

Problemy w SRC

2. [Apache HTTP Server Project](#)

Effort to develop and maintain an open-source HTTP server for modern operating systems including UNIX and Windows NT.

<http://apache.org> - 32k - [Cached](#) - [More from this site](#) - [Save](#)

49. [pre-FAQ - The Apache Software Foundation](#)

... most of the common queries that we receive about our software and the **Apache** Software Foundation ... something similar indicating that **Apache** has been installed) on your screen ...

www.apache.org/foundation/preFAQ.html - 32k - [Cached](#) - [More from this site](#) - [Save](#)

587. [Apache C++ Standard Library](#)

Last Modified: \$Date: 2006-02-16 09:05:15 -0800 (Thu, 16 Feb 2006) \$ stdcxx. STDCXX - **Apache** C++ Standard Library. what is stdcxx? ... The goal of the **Apache** C++ Standard Library project is to provide a free implementation of the ISO ... C++ Standard Library to the **Apache** stdcxx project, a proven code base ...

incubator.apache.org/stdcxx - 41k - [Cached](#) - [More from this site](#) - [Save](#)

Wymagania

- Kryteria oceniania jakości metod SRC:
 - Semantyczność: dokumenty w jednej grupie powinny dotyczyć tego samego tematu
 - Znaczenie etykiet grup: powinny one dobrze opisać zawartość całej grupy.
 - Mała liczba grup: należy pokryć jak najwięcej dokumentów używając przy tym jak najmniej grup.
- Te kryteria są raczej subiektywne aniżeli obiektywne.

Model wektorowy dokumentów

- $T = \{t_1, \dots, t_n\}$ – zbiór wybranych wyrazów (słów, fraz)
- Dokument $d_i = [w_{i,1}, \dots, w_{i,n}]$ gdzie $w_{i,j}$ jest wagą wyrazu t_j w dokumencie d_i
- Schemat ważenia wyrazów – TFxIDF

$$w_{i,j} = f_{i,j} \times \log \frac{N}{df(t_j)}$$

$w_{i,j}$: częstość występowania wyrazu t_j w dokumencie d_i

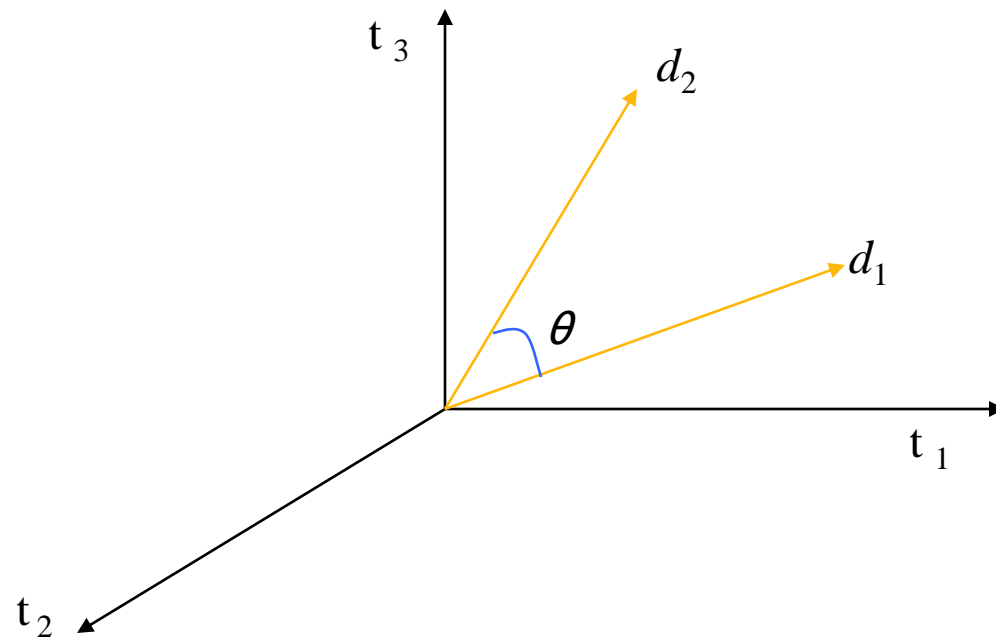
N : liczba dokumentów

$df(t_j)$: liczba dokumentów zawierających t_j

Podobieństwo dokumentów

- Miara cosinusa:

$$\text{sim}(d_j, d_k) = \frac{\vec{d}_j \cdot \vec{d}_k}{|\vec{d}_j| |\vec{d}_k|} = \frac{\sum_{i=1}^n w_{i,j} w_{i,k}}{\sqrt{\sum_{i=1}^n w_{i,j}^2} \sqrt{\sum_{i=1}^n w_{i,k}^2}}$$



Istniejące metody

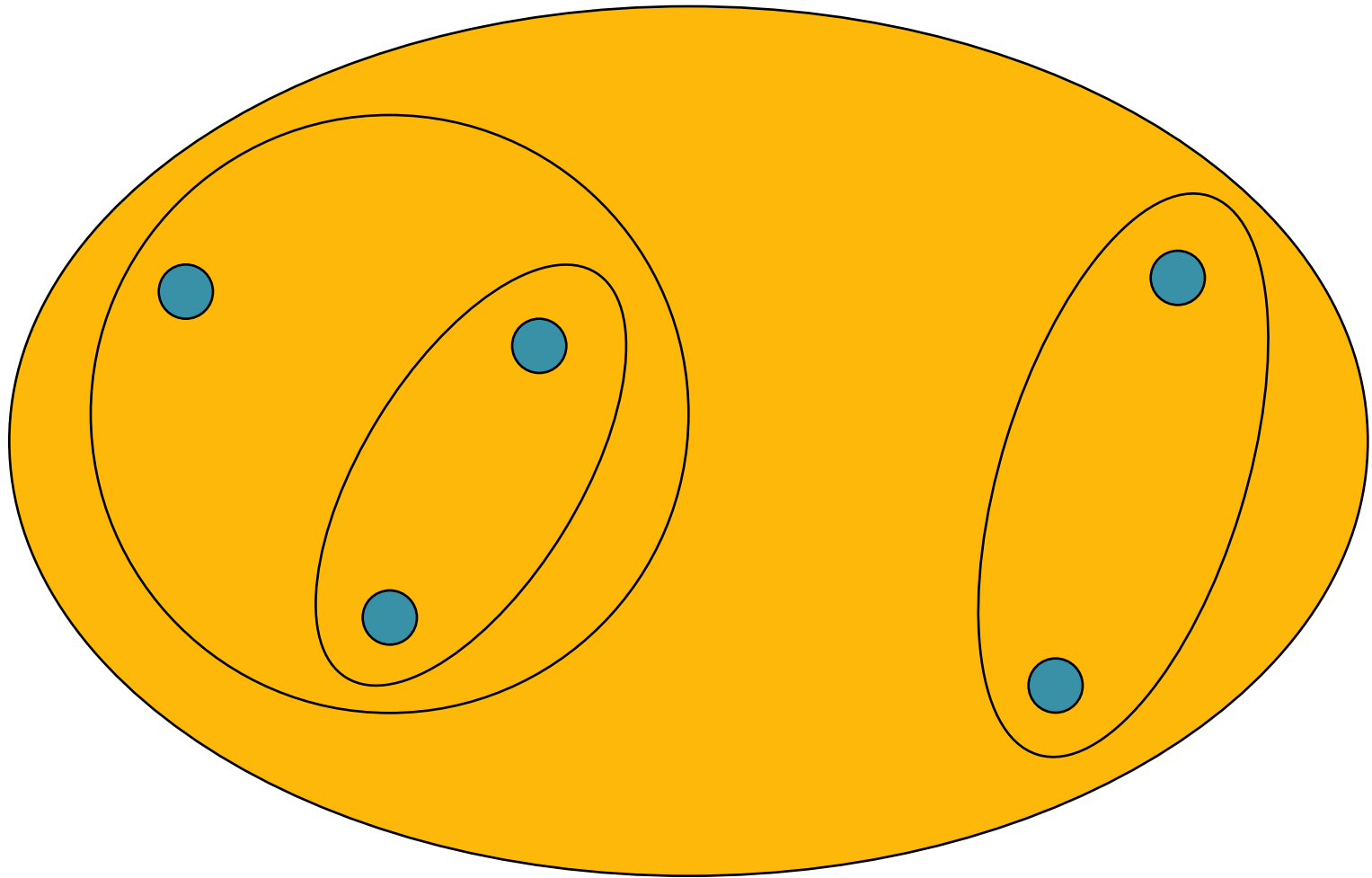
Name	Word Flat	Sentences Flat	Word Hier.	Sentences Hier.	Online/ Software
WebCat	+				+
Retriever	+				
Scatter/Gather	+				
Wang <i>et al.</i>	+				
Grouper		+			
Carrot		+			+
Lingo		+			+
Microsoft		+			
FICH			+		+
Credo			+		+
IBM			+		
SHOC				+	
CIIRarchies				+	+
LA				+	
Highlight				+	+
SnakeT				+	+
Mooter			+		+
Vivisimo				+	+

Klasyfikacja algorytmów grupowania

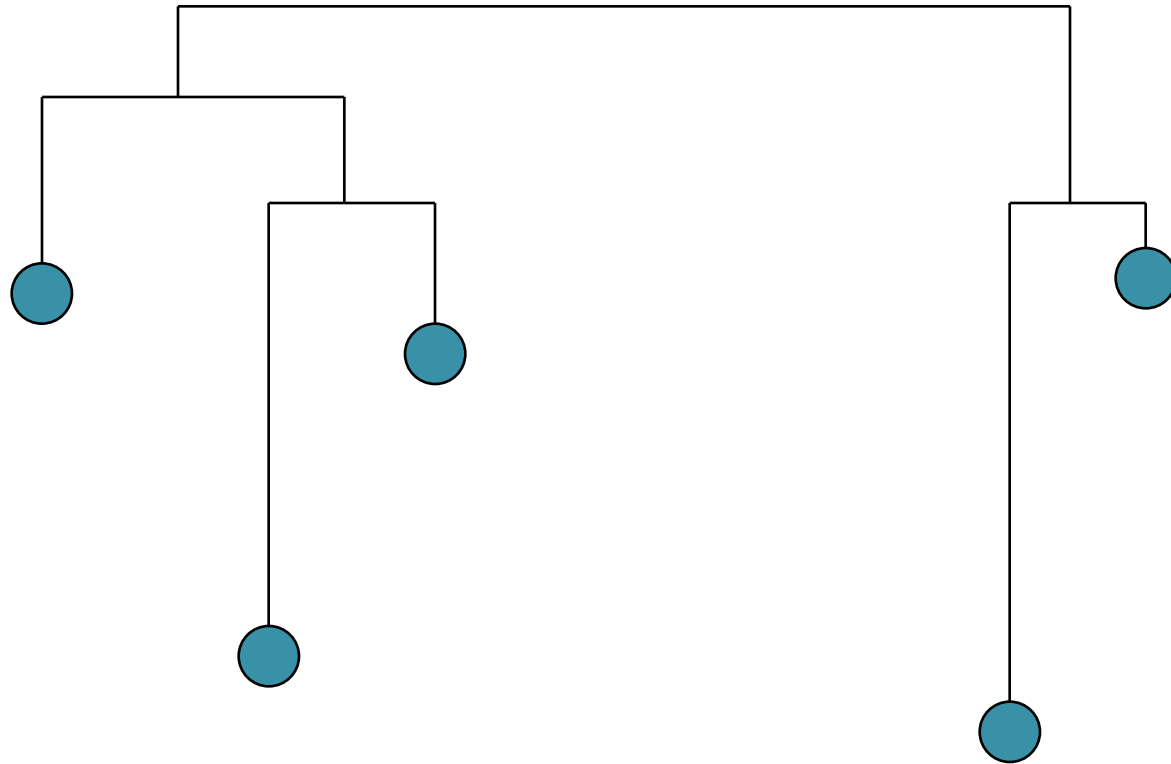
- Płaska struktura czy hierarchiczna?
- Czy grupy są rozłączne?
- Ostry czy miękki podział?
- Przyrostowa metoda?
- Czy liczba grup jest z góry zadana?
- Czy miary odległości lub podobieństwa muszą być zadane z góry?

- Z użyciem odległości
 - Hierarchiczna struktura
 - Agglomerative Hierarchical Clustering (AHC)
 - Płaska struktura
 - K-centroidów (możliwe rozmycie)
 - Inkrementalna (Single-pass)
- Inne
 - Suffix Tree Clustering (Grouper)
 - SOM (Kohonen)
 - Latent Semantic Indexing (LSI) (zmniejsza wymiar)

Grupowanie hierarchiczne (AHC)

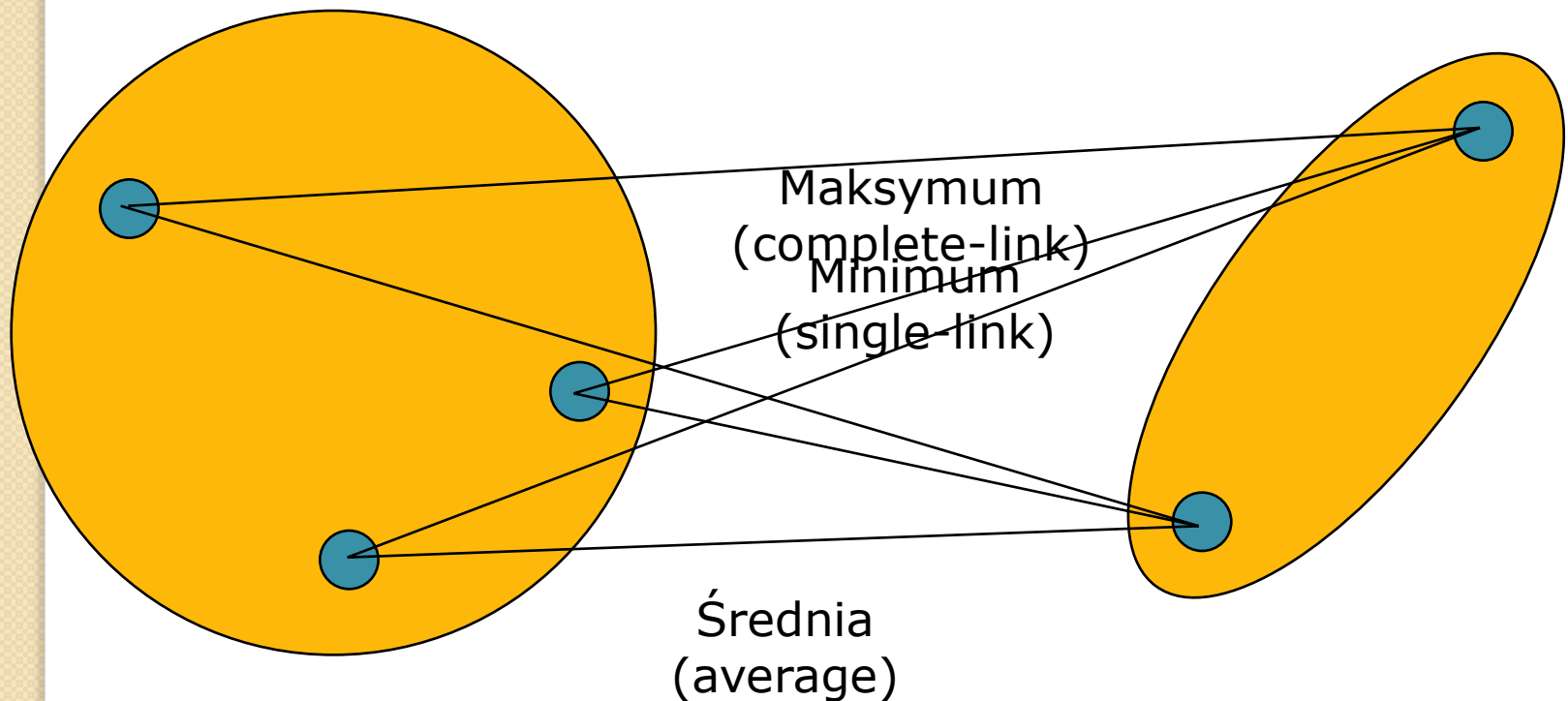


Wynik grupowania: hierarchia pojęć

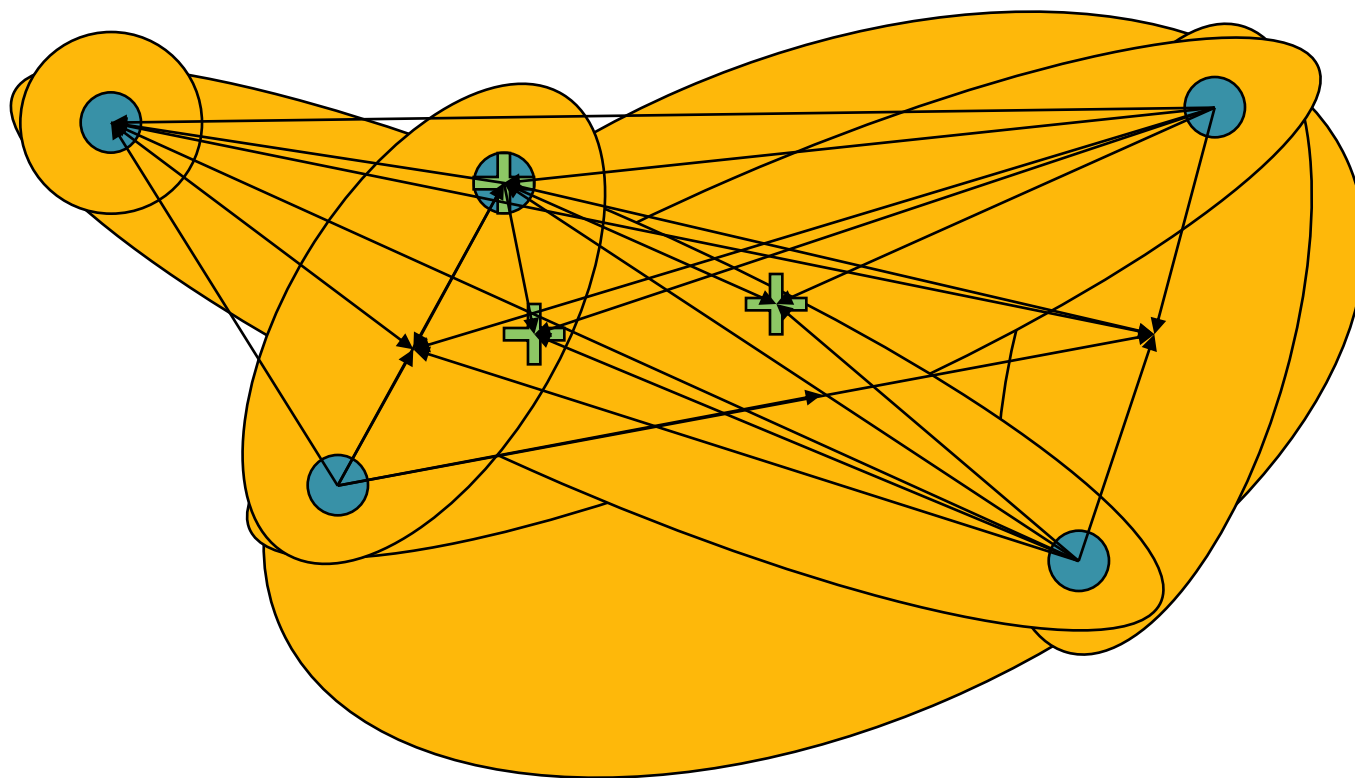


Różne wersje AHC

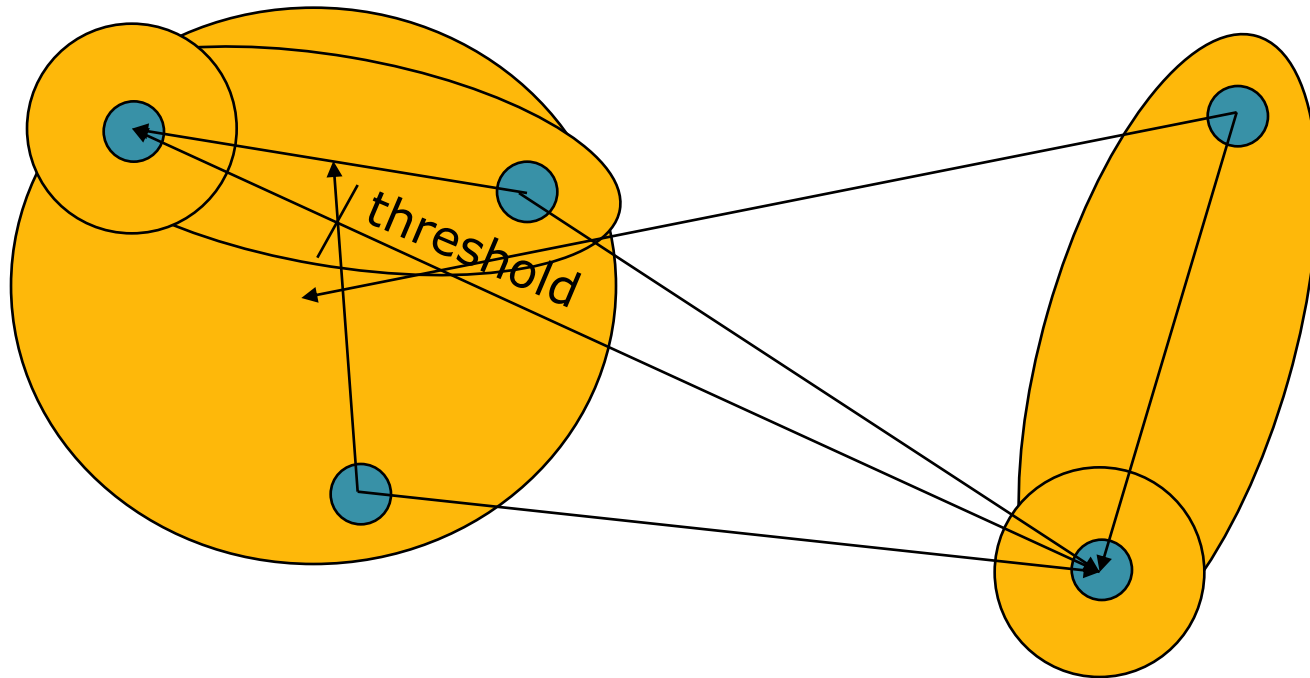
- Istnieją różne metody mierzenia podobieństwa grup



K-centroidów (k=3)



Metoda inkrementalna (single-pass)



„Grouper”

(Zamir and Etzioni 1997, 1999)

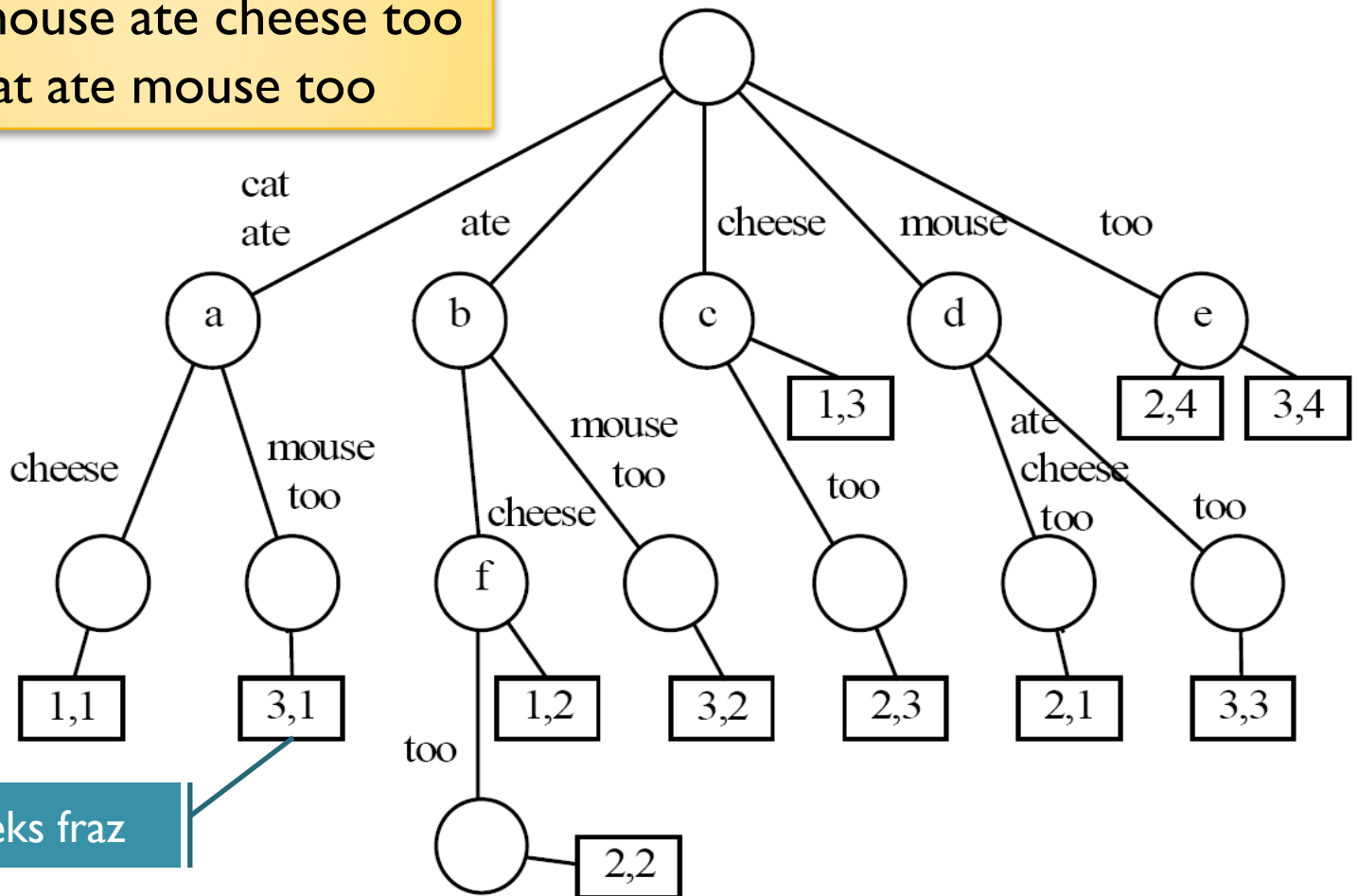
- Działa na bieżąco (online)
- Grupuje wyniki wyszukiwania (snippets)
- Grupuje dokumenty, które mają wiele wspólnych fraz
- Grupowanie drzewem sufiksowym (STC - Suffix Tree Clustering)
 - Czas liniowy
 - Metoda inkrementalna
 - Grupy nie są rozłączne
 - Może być hierarchiczna.

Algorytm STC (Suffix Tree Clustering)

- Krok 1: Czyszczenie danych:
 - Normalizacja (stemming, stop-words elimination)
 - Identyfikacja fraz i zdań.
 - Eliminacja znaków interpunkcyjnych.
- Krok 2: Budowa drzewa sufiksowego:
 - Stworzenie grup bazowych
 - Ocena grup bazowych za pomocą ich rozmiaru i ocen fraz
- Krok 3: Łączenie grup bazowych:
 - Grupy mające „dużą część wspólną” są połączone.

Drzewo sufiksowe = minimalne drzewo zawierające sufiksy wszystkich napisów

1. cat ate cheese
2. mouse ate cheese too
3. cat ate mouse too



Odwrotny indeks fraz

Krok 2 – Identyfikacja grup bazowych

a	cat ate	1,3
b	ate	1,2,3
c	cheese	1,2
d	mouse	2,3
e	too	2,3
f	ate cheese	1,2

- Wierzchołki reprezentują grupy dokumentów mających wspólną frazę
- Każda grupa B definiowana przez frazę P jest oceniona przez $S(B) = |B|f(|P|)$

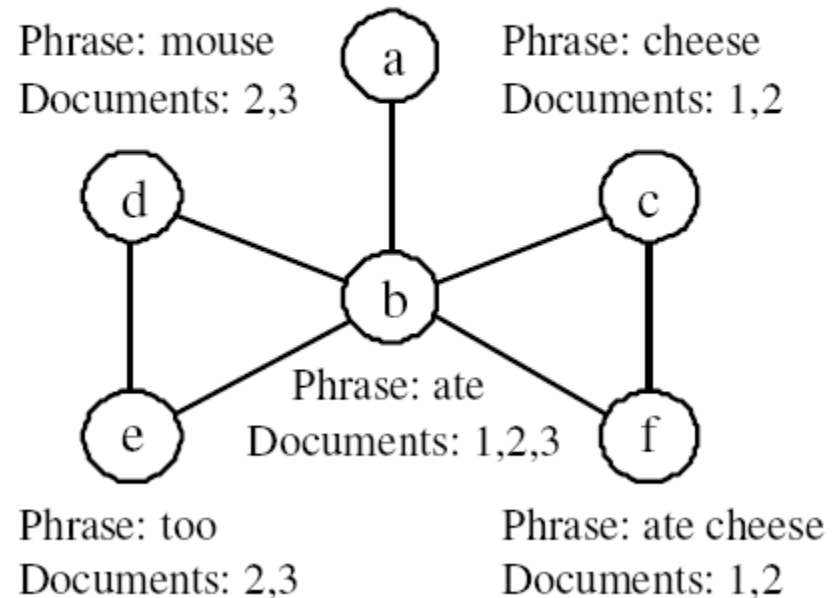
Krok 3 – Łączenie grup bazowych

- Podobieństwo między grupami bazowymi:

$$sim = \begin{cases} 1 & \frac{|B_n \cap B_m|}{|B_n|} > 0.5 \text{ oraz } \frac{|B_n \cap B_m|}{|B_m|} > 0.5 \\ 0 & \text{wpp.} \end{cases}$$

Phrase: cat ate
Documents: 1,3

- Łączymy grupy algorytmem przyrostowym



„Lingo”

(S.Osiński, D.Weiss)

- Korzysta z rozkładu macierzy wzg. wartości osobliwych (SVD)
- Reprezentacja zbioru dokumentów (snippets) w przestrzeni rzutowej o małym wymiarze
- Wektory osobliwe wyznaczają etykiety grup
- Dokumenty są dopisane do grup według miary cosinusa.
- Implementacja:

Carrot2: Search Results Clustering Framework

<http://www.carrot2.org>

<http://carrot.cs.put.poznan.pl>

Rozkład wzg. wartości osobliwych

(ang. SVD - Singular Value Decomposition)

$$\begin{array}{l} \text{term 1} \\ \text{term 2} \\ \text{term 3} \\ \text{term 4} \\ \text{term 5} \\ \text{term 6} \end{array} \begin{bmatrix} \text{doc 1} & \text{doc 2} & \text{doc 3} & \text{doc 4} \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} = A$$

- A – macierz $m \times n$

$$A = U \Sigma V^T$$

- Kolumny U – wektory własne AA^T
- Kolumny V – wektory własne $A^T A$
- $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$:
 $\sigma_1, \dots, \sigma_n$ wartości osobliwe A
 $\sigma_1 > \dots > \sigma_k \dots > \sigma_n$

- **Aproksymacja:**

$$A \approx U_k \Sigma_k V_k^T = U_k C_k$$

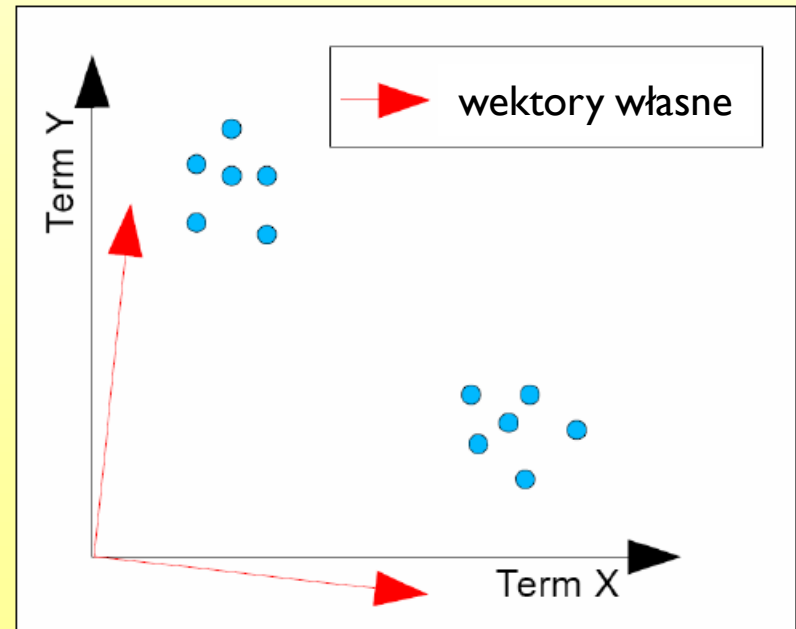
SVD

$$A = \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} \approx \begin{bmatrix} * & * \\ * & * \\ * & * \\ * & * \\ * & * \end{bmatrix} \times \begin{bmatrix} * & * & * & * \\ * & * & * & * \end{bmatrix}$$

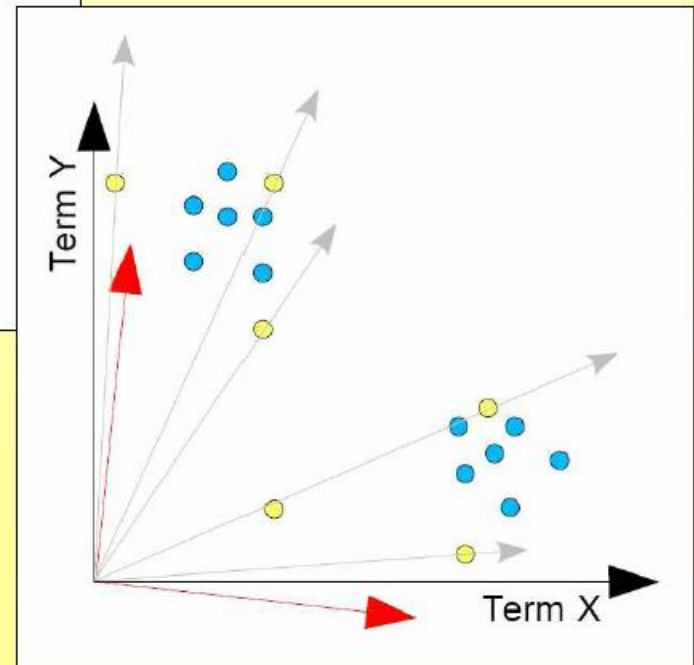
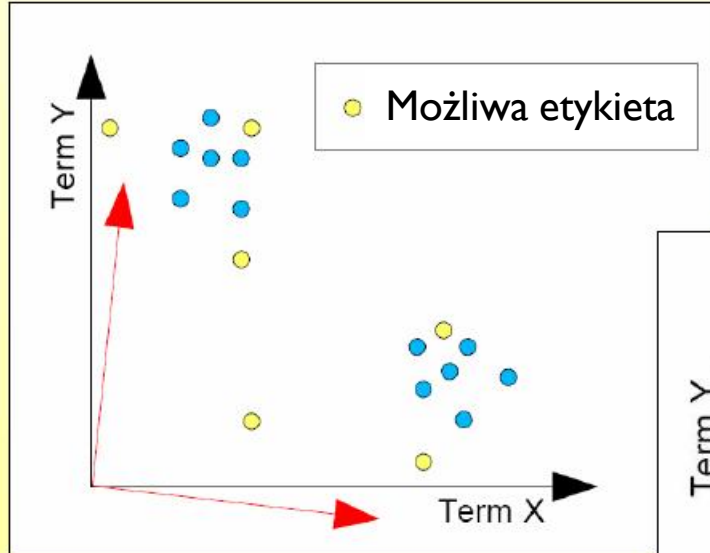
Współrzędne dokumentów
w przestrzeni rzutowej

Wektory
własne

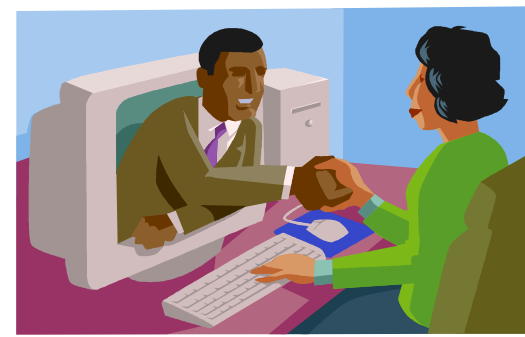
$$A \approx U_k C_k$$



SVD wyznacza etykiety grup



Konkluzje



- SRC – próba przyśpieszania procesu wyszukiwania informacji w internecie i w bibliotekach elektronicznych.
- Temat atrakcyjny również dla dużych graczy
- Problemy:
 - Brak obiektywnego kryterium oceny
 - Brak personalizacji
- Źródła informacji:
 - Historie procesów wyszukiwania w przeszłości
 - Publiczne katalogi internetowe
 - Leksykon semantycznych powiązań, np. Wordnet
 - Profil użytkownika