

3.3 Index Access Scheduling

Given:

- index scans over m lists L_i ($i=1..m$), with current positions pos_i
- score predictors for $score(pos)$ and $pos(score)$ for each L_i
- selectivity predictors for document $d \in L_i$
- current top- k queue T with k documents
- candidate queue Q with c documents (usually $c \gg k$)
- min- k threshold = $\min\{\text{worstscore}(d) \mid d \in T\}$

Questions/Decisions:

- *Sorted-access (SA) scheduling:*

for the next batch of b scan steps, how many steps in which list?

(b_i steps in L_i with $\sum_i b_i = b$)

- *Random-access (RA) scheduling:*

when to initiate probes and for which documents?

- Possible constraints and extra considerations:

some dimensions i may support only sorted access or only random access, or have tremendous cost ratio C_{RA}/C_{SA}

Combined Algorithm (CA)

assume cost ratio $C_{RA}/C_{SA} = r$

perform NRA (TA-sorted)

with [worstscore, bestscore] bookkeeping in priority queue Q
and round-robin SA to m index lists

...

after every r rounds of SA (i.e. $m \cdot r$ scan steps)

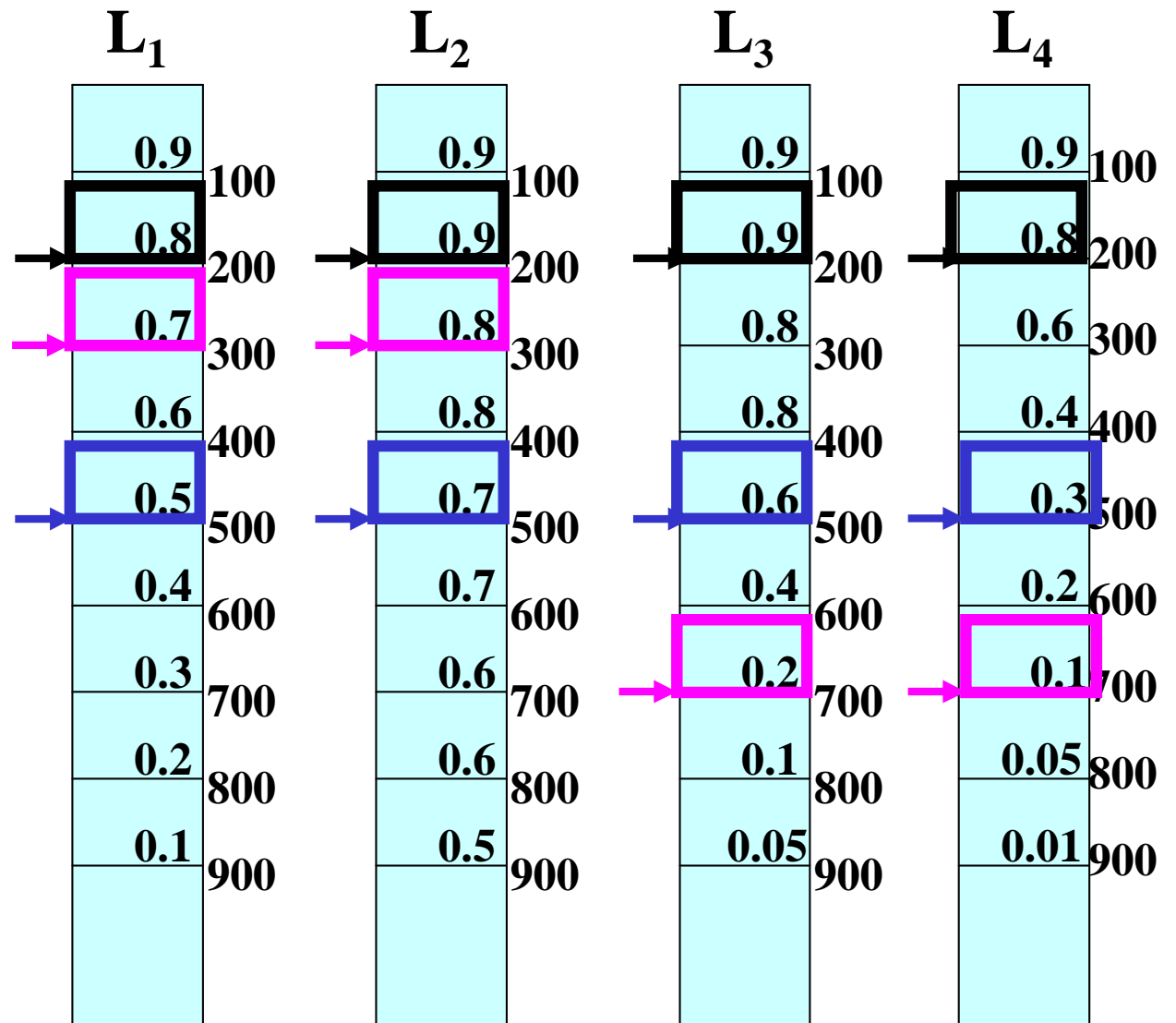
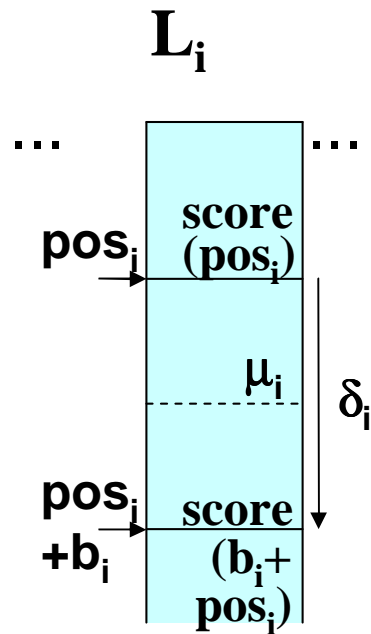
perform RA to look up all missing scores of „best candidate“ in Q
(where „best“ is in terms of
bestscore, worstscore, or $E[\text{score}]$, or $P[\text{score} > \text{min-k}]$)

cost competitiveness w.r.t. „optimal schedule“

(scan until $\sum_i \text{high}_i \leq \min\{\text{bestscore}(d) \mid d \in \text{final top-k}\}$,
then perform RAs for all d' with $\text{bestscore}(d') > \text{min-k}$): $4m + k$

Sorted-Access Scheduling

available info:



goal:
eliminate candidates quickly
aim for quick drop in high_i bounds

SA Scheduling: Objective and Heuristics

plan next b_1, \dots, b_m index scan steps

for batch of b steps overall s.t. $\sum_{i=1..m} b_i = b$

and $\text{benefit}(b_1, \dots, b_m)$ is max!

possible benefit definitions:

$$\text{benefit}(b_1..b_m) = \sum_{i=1..m} \Delta_i \quad \text{with} \quad \Delta_i = (\text{high}_i - \text{score}_i(\text{pos}_i + b_i)) / b_i$$

score gradient

$$\text{benefit}(b_1..b_m) = \sum_{i=1..m} \delta_i \quad \text{with} \quad \delta_i = \text{score}_i(\text{pos}_i) - \text{score}_i(\text{pos}_i + b_i)$$

score reduction

Solve knapsack-style NP-hard optimization problem
(e.g. for batched scans) or use greedy heuristics:

$$b_i := b * \text{benefit}(b_i=b) / \sum_{v=1..m} \text{benefit}(b_v=b)$$

SA Scheduling: Benefit Aggregation Heuristics

Consider current top-k T and candidate queue Q ;
 for each $d \in T \cup Q$ we know $E(d) \subseteq 1..m$, $R(d) = 1..m - E(d)$,
 $bestscore(d)$, $worstscore(d)$, $p(d) = P[\text{score}(d) > \text{min-k}]$

$$benefit(d, b_1..b_m) =$$

$$surplus(d)^{-1} \cdot \sum_{i \notin E(d)} (high_i - score(pos_i + b_i))$$

$$+ gap(d)^{-1} \cdot \sum_{i \notin E(d)} \mu_i$$

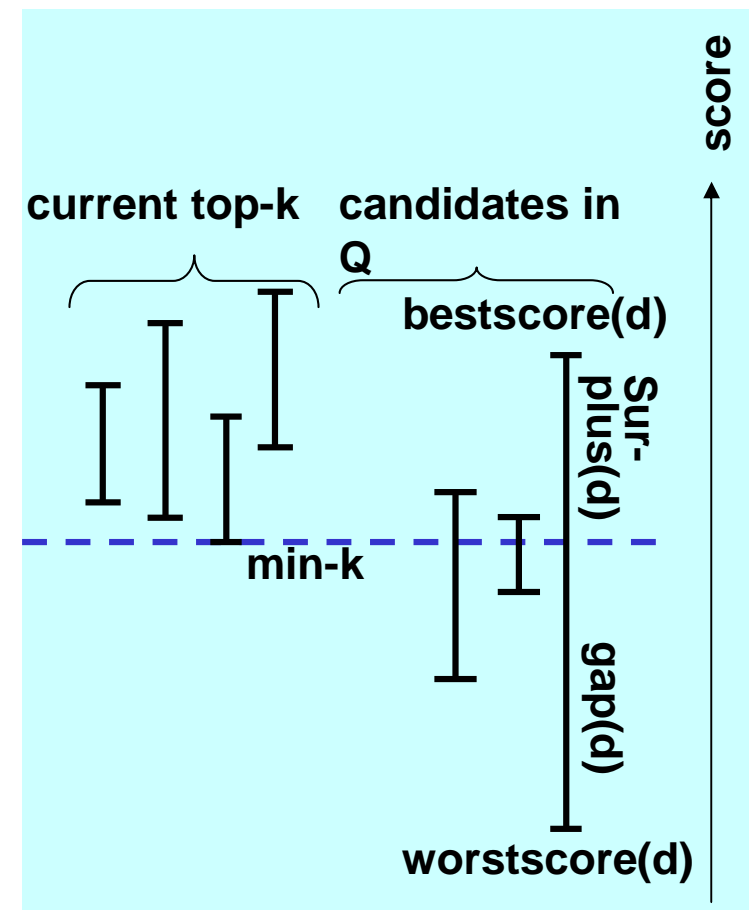
with $surplus(d) = bestscore(d) - \text{min-k}$

$gap(d) = \text{min-k} - worstscore(d)$

$\mu_i = E[\text{score}(j) \mid j \in [pos_i, pos_i + b_i]]$

$$benefit(b_1..b_m) = \sum_{d \in T \cup Q} benefit(d, b_1..b_m)$$

weighs documents and dimensions in benefit function



Random-Access Scheduling: Heuristics

Perform additional RAs when helpful

- 1) to increase min-k (increase worstscore of $d \in \text{top-k}$) or
- 2) to prune candidates (decrease bestscore of $d \in Q$)

For 1) Top Probing:

- perform RAs for current top-k (whenever min-k changes),
- and possibly for best d from Q
(in desc. order of bestscore, worstscore, or $P[\text{score}(d) > \text{min-k}]$)

For 2) 2-Phase Probing:

perform RAs for all candidates at point t

total cost of remaining RAs = total cost of SAs up to t

(motivated by linear increase of SA-cost(t) and sharply decreasing remaining-RA-cost(t))

Top-k Queries over Web Sources

Typical example:

Address = „2590 Broadway“ and Price = \$ 25 and Rating = 30
issued against mapquest.com, nytoday.com, zagat.com

Major complication:

some sources do not allow sorted access
highly varying SA and RA costs

Major opportunity:

sources can be accessed in parallel

→ extension/generalization of TA

distinguish S-sources, R-sources, SR-sources

Source-Type-Aware TA

For each R-source $S_i \in S_{m+1} \dots S_{m+r}$ set $high_i := 1$

Scan SR- or S-sources $S_1 \dots S_m$

Choose SR- or S-source S_i for next sorted access

for object d retrieved from SR- or S-source L_i do {

$E(d) := E(d) \cup \{i\}; high_i := si(q,d);$

$bestscore(d) := aggr\{x_1, \dots, x_m\}$ with $x_i := si(q,d)$ for $i \in E(d)$, $high_i$ for $i \notin E(d)$;

$worstscore(d) := aggr\{x_1, \dots, x_m\}$ with $x_i := si(q,d)$ for $i \in E(d)$, 0 for $i \notin E(d)$; }

Choose SR- or R-source S_i for next random access

for object d retrieved from SR- or R-source L_i do {

$E(d) := E(d) \cup \{i\};$

$bestscore(d) := aggr\{x_1, \dots, x_m\}$ with $x_i := si(q,d)$ for $i \in E(d)$, $high_i$ for $i \notin E(d)$;

$worstscore(d) := aggr\{x_1, \dots, x_m\}$ with $x_i := si(q,d)$ for $i \in E(d)$, 0 for $i \notin E(d)$; }

current top-k := k docs with largest worstscore;

min-k := minimum worstscore among current top-k;

Stop when $bestscore(d \mid d \text{ not in current top-k results}) \leq min-k$;

Return current top-k;

essentially NRA with choice of sources

Strategies for Choosing the Source for Next Access

for next sorted access:

Escore(Li) := expected si value for next sorted access to Li

(e.g.: $high_i$)

rank(Li) := $w_i * \text{Escore}(\text{Li}) / c_{SA}(\text{Li})$

// w_i is weight of Li in aggr

// $c_{SA}(\text{Li})$ is source-specific SA cost

choose SR- or S-source with highest rank(Li)

for next random access (probe):

Escore(Li) := expected si value for next random access to Li

(e.g.: $(high_i - low_i) / 2$)

rank(Li) := $w_i * \text{Escore}(\text{Li}) / c_{RA}(\text{Li})$

choose SR- or R-source with highest rank(Li)

or use more advanced statistical score estimators

The Upper Strategy for Choosing Next Object and Source (Marian et al.: TODS 2004)

idea: eagerly prove that candidate objects cannot qualify for top-k

for next random access:

among all objects with $E(d) \neq \emptyset$ and $R(d) \neq \emptyset$

choose d' with the highest $\text{bestscore}(d')$;

if $\text{bestscore}(d') < \text{bestscore}(v)$ for object v with $E(v) = \emptyset$ then

perform sorted access next (i.e., don't probe d')

else {

$\Delta := \text{bestscore}(d') - \text{min-k}$;

if $\Delta > 0$ then {

consider Li as „redundant“ for d' if for all $Y \subseteq R(d') - \{Li\}$

$$\sum_{j \in Y} w_j * \text{high}_j + w_i * \text{high}_i \geq \Delta \Rightarrow \sum_{j \in Y} w_j * \text{high}_j \geq \Delta ;$$

choose „non-redundant“ source with highest $\text{rank}(Li)$ }

else choose source with lowest $c_{RA}(Li)$;

};

The Parallel Strategy pUpper (Marian et al.: TODS 2004)

*idea: consider up to $MPL(L_i)$ parallel probes to the same R-source L_i
choose objects to be probed based on
bestscore reduction and expected response time*

for next random access:

probe-candidates := m objects d with $E(d) \neq \emptyset$ and $R(d) \neq \emptyset$
such that d is among the m highest values of $bestscore(d)$;

for each object d in probe-candidates do {

$\Delta := bestscore(d) - \min-k$;

if $\Delta > 0$ then {

choose subset $Y(d) \subseteq R(d)$ such that $\sum_{j \in Y} w_j * high_j \geq \Delta$
and expected response time

$\sum_{L_j \in Y(d)} (|\{d' \mid bestscore(d') > bestscore(d) \text{ and } Y(d) \cap Y(d') \neq \emptyset\}|$
 $* c_{RA}(L_j) / MPL(L_j))$

is minimum };

};

enqueue $probe(d)$ to $queue(L_i)$ for all $L_i \in Y(d)$
with expected response time as priority;

Experimental Evaluation

pTA:
parallelized TA
(with asynchronous probes,
but same probe order as TA)

synthetic data

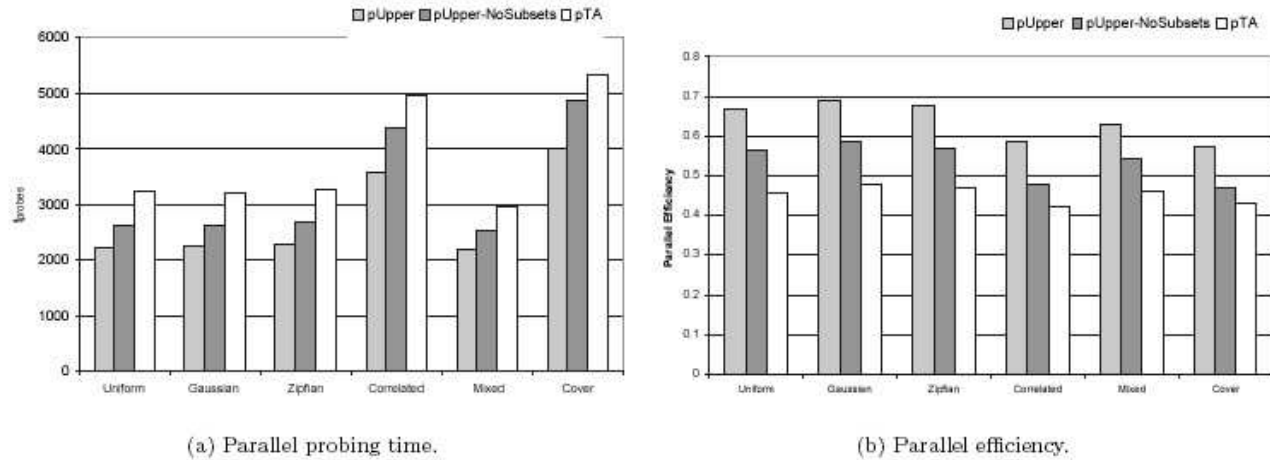
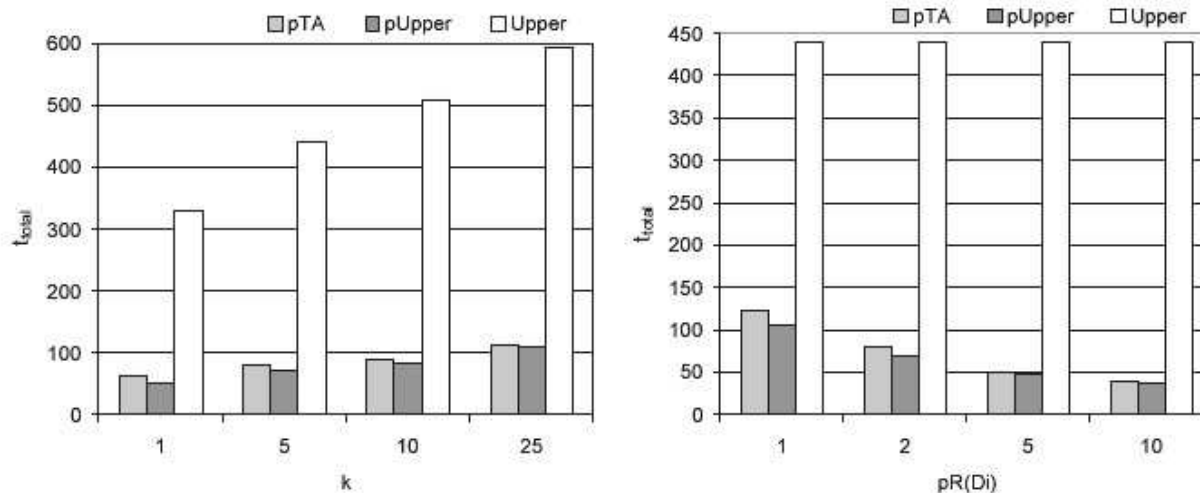


Figure 19: Effect of the attribute score distribution on performance.



(a) Parallel time t_{total} as a function of k ($pR(D_i) = 2$). (b) Parallel time t_{total} as a function of $pR(D_i)$ ($k = 5$).

real Web sources

- SR: superpages (Verizon yellow pages)
- R: subwaynavigator
- R: mapquest
- R: altavista
- R: zagat
- R: nytoday

from: A. Marian et al., TODS 2004

Figure 24: Effect of the number of objects requested k (a) and the number of accesses per source $pR(D_i)$ (b) on the performance of *pTA*, *pUpper*, and *Upper* over real web sources.

3.4 Index Organization and Advanced Query Types

Richer Functionality:

- **Boolean combinations of search conditions**
- **Search by word stems**
- **Phrase queries and proximity queries**
- **Wild-card queries**
- **Fuzzy search with edit distance**

Enhanced Performance:

- **Stopword elimination**
- **Static index pruning**
- **Duplicate elimination**

Boolean Combinations of Search Conditions

combination of AND and ANDish: $(t_1 \text{ AND } \dots \text{ AND } t_j) t_{j+1} t_{j+2} \dots t_m$

- TA family applicable with mandatory probing in AND lists
→ RA scheduling
- (worstscore, bestscore) bookkeeping and pruning
more effective with “boosting weights” for AND lists

combination of AND, ANDish and NOT:

NOT terms considered k.o. criteria for results

TA family applicable with mandatory probing for AND and NOT

→ RA scheduling

combination of AND, OR, NOT in Boolean sense:

- best processed by index lists in DocId order
- construct operator tree and push selective operators down;
needs good query optimizer (selectivity estimation)

Search with Morphological Reduction (Lemmatization)

Reduction onto grammatical ground form:

nouns onto nominative, verbs onto infinitive,
plural onto singular, passive onto active, etc.

Examples (in German):

- „Winden“ onto „Wind“, „Winde“ or „winden“
depending on phrase structure and context
- „finden“ and „gefundenes“ onto „finden“,
- „Gefundenes“ onto „Fund“

Reduction of morphological variations onto word stem:

flexions (e.g. declination), composition, verb-to-noun, etc.

Examples (in German):

- „Flüssen“, „einflößen“ onto „Fluss“,
- „finden“ and „Gefundenes“ onto „finden“
- „Du brachtest ... mit“ onto „mitbringen“,
- „Schweinkram“, „Schweinschaxe“ and „Schweinebraten“
onto „Schwein“ etc.
- „Feinschmecker“ and „geschmacklos“ onto „schmecken“

Stemming

Approaches:

- Lookup in comprehensive lexicon/dictionary (e.g. for German)
- Heuristic affix removal (e.g. Porter stemmer for English):
remove prefixes and/or suffixes
based on (heuristic) rules

Example:

stresses → stress, stressing → stress, symbols → symbol
based on rules: sses → ss, ing → ε, s → ε, etc.

The benefit of stemming for IR is debated.

Example:

Bill is operating a company.

On his computer he runs the Linux operating system.

Phrase Queries and Proximity Queries

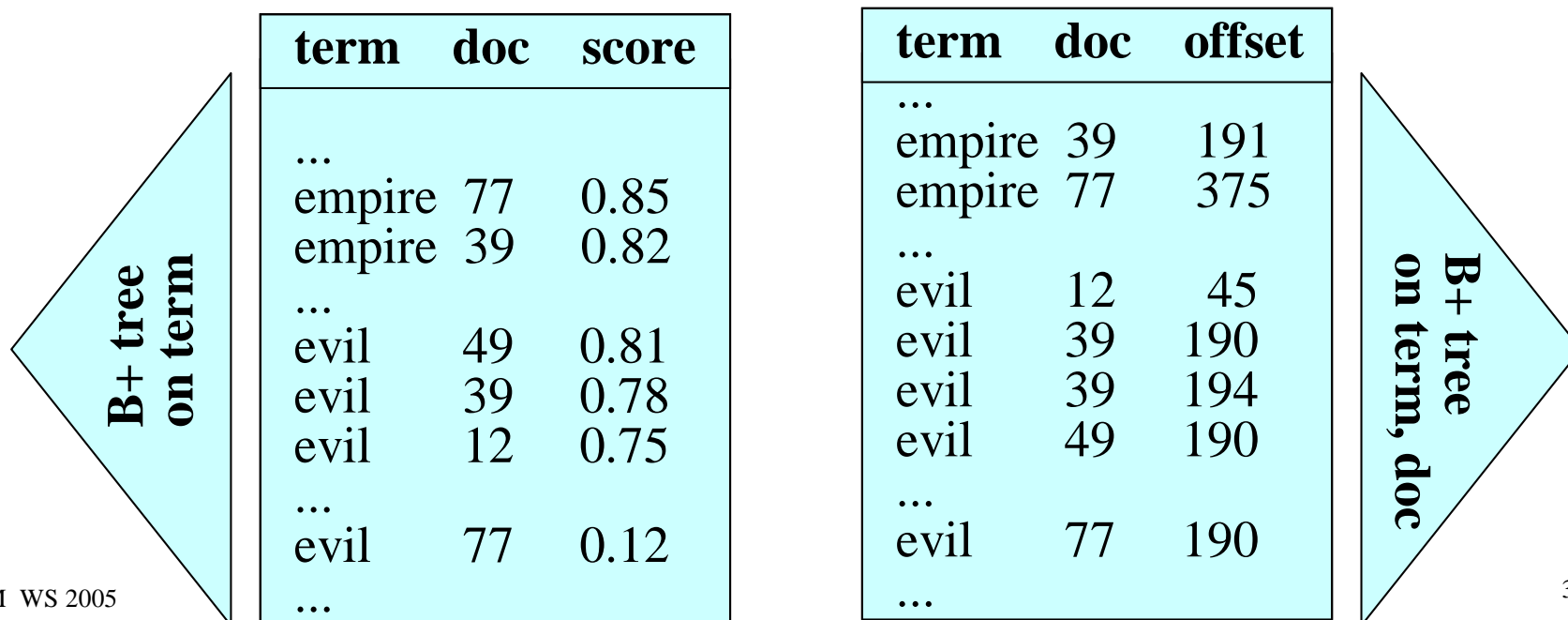
phrase queries such as:

„George W. Bush“, „President Bush“, „The Who“, „Evil Empire“, „PhD admission“, „FC Schalke 04“, „native American music“, „to be or not to be“, „The Lord of the Rings“, etc. etc.

difficult to anticipate and index all (meaningful) phrases
sources could be thesauri (e.g. WordNet) or query logs

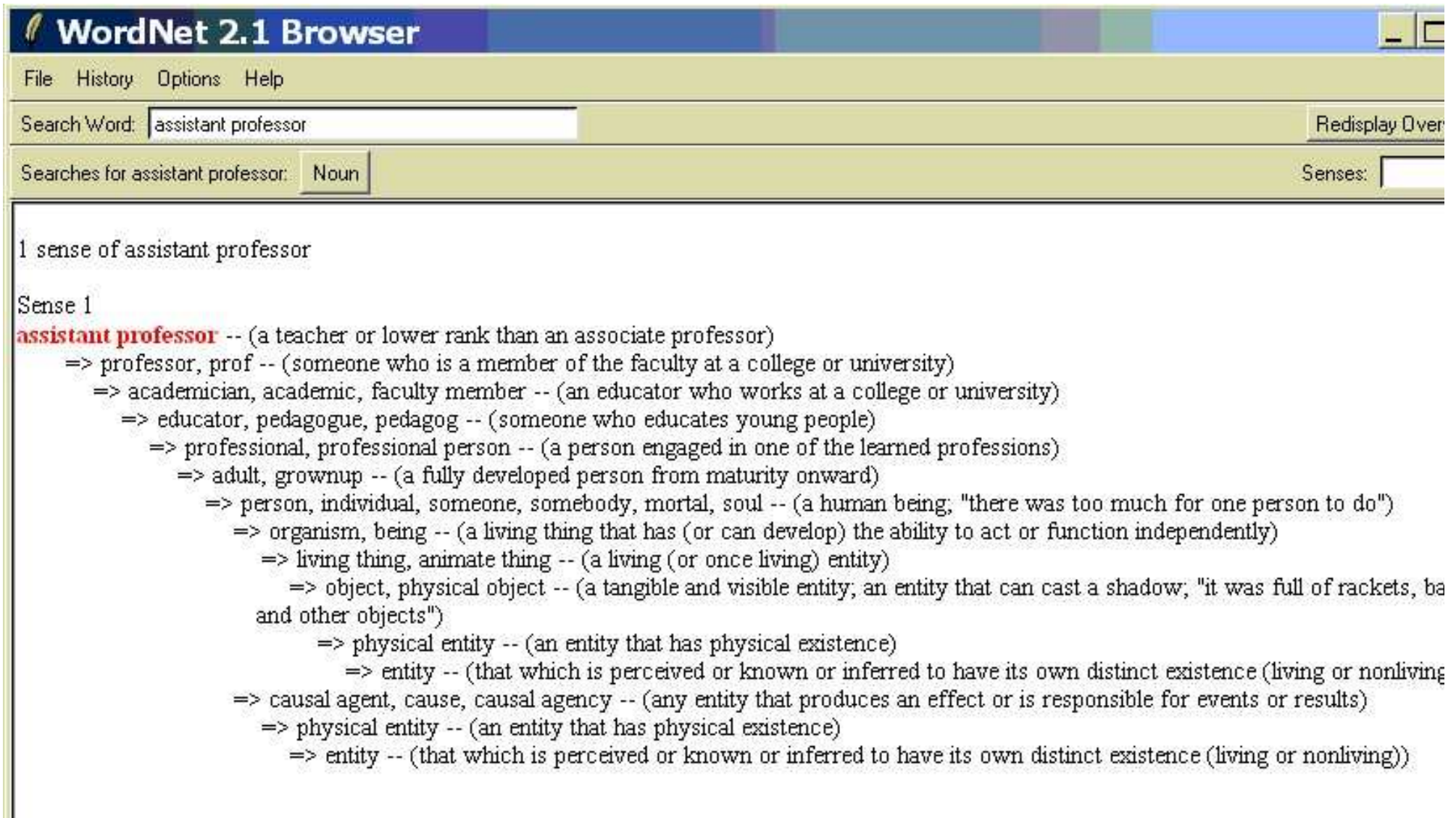
→ standard approach:

combine single-term index with separate position index



Thesaurus as Phrase Dictionary

Example: WordNet (Miller/Fellbaum), <http://wordnet.princeton.edu>



The screenshot shows the WordNet 2.1 Browser interface. The title bar reads "WordNet 2.1 Browser". The menu bar includes "File", "History", "Options", and "Help". The search bar contains the text "assistant professor" and a "Redisplay Over" button. Below the search bar, it shows "Searches for assistant professor: Noun" and "Senses: []". The main content area displays "1 sense of assistant professor" and "Sense 1". The definition for "assistant professor" is provided, followed by a list of related terms and their definitions.

1 sense of assistant professor

Sense 1

assistant professor -- (a teacher or lower rank than an associate professor)

- => professor, prof -- (someone who is a member of the faculty at a college or university)
- => academician, academic, faculty member -- (an educator who works at a college or university)
- => educator, pedagogue, pedagog -- (someone who educates young people)
- => professional, professional person -- (a person engaged in one of the learned professions)
- => adult, grownup -- (a fully developed person from maturity onward)
- => person, individual, someone, somebody, mortal, soul -- (a human being; "there was too much for one person to do")
- => organism, being -- (a living thing that has (or can develop) the ability to act or function independently)
- => living thing, animate thing -- (a living (or once living) entity)
 - => object, physical object -- (a tangible and visible entity; an entity that can cast a shadow; "it was full of rackets, ba and other objects")
 - => physical entity -- (an entity that has physical existence)
 - => entity -- (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))
- => causal agent, cause, causal agency -- (any entity that produces an effect or is responsible for events or results)
- => physical entity -- (an entity that has physical existence)
- => entity -- (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

Biword and Phrase Indexing

build index over all word pairs:

**index lists (term1, term2, doc, score) or
for each term1 nested list (term2, doc, score)**

variations:

- **treat nearest nouns as pairs,
or discount articles, prepositions, conjunctions**
- **index phrases from query logs, compute correlation statistics**

query processing:

- **decompose even-numbered phrases into biwords**
- **decompose odd-numbered phrases into biwords
with low selectivity (as estimated by $df(\text{term1})$)**
- **may additionally use standard single-term index if necessary**

Examples:

to be or not to be → (to be) (or not) (to be)

The Lord of the Rings → (The Lord) (Lord of) (the Rings)

N-Gram Indexing and Wildcard Queries

Queries with wildcards (simple regular expressions),
to capture mis-spellings, name variations, etc.

Examples:

Brit*ney, Sm*th*, Go*zilla, Marko*, reali*ation, *raklion

Approach:

- decompose words into N-grams of N successive letters
and index all N-grams as terms
- query processing computes AND of N-gram matches

Example (N=3):

Brit*ney → Bri AND rit AND ney

**Generalization: decompose words into frequent fragments
(e.g., syllables, or fragments derived from mis-spelling statistics)**

Refstring Indexing (Schek 1978)

In addition to indexing all N-grams for some small N (e.g. 2 or 3), determine frequent fragments – refstrings $r \in R$ – with properties:

- $df(r)$ is above some threshold θ
- if $r \in R$ then for all substrings s of r : $s \notin R$ unless $df(s|\neg r) = |\{\text{docs } d \mid d \text{ contains } s \text{ but not } r\}| \geq \theta$

Refstring index build:

1) Candidate generation \rightarrow preliminary set R :

generate strings r with $|r| > N$ in increasing length, compute $df(r)$;
remove r from candidates if $r=xy$ with $df(x) < \theta$ or $df(y) < \theta$

2) Candidate selection: consider candidate $r \in R$ with $|r|=k$ and sets

$left(r) = \{xr \mid xr \in R \wedge |xr|=k+1\}$, $right(r) = \{ry \mid ry \in R \wedge |ry|=k+1\}$,

$left^-(r) = \{xr \mid xr \notin R \wedge |xr|=k+1\}$, $right^-(r) = \{ry \mid ry \notin R \wedge |ry|=k+1\}$

select r if $weight(r) = df(r) - \max\{leftf(r), rightf(r)\} \geq \theta$ with

$leftf(r) = \sum_{q \in left(r)} df(q) + \sum_{q \in left^-(r)} \max\{leftf(q), rightf(q)\}$ and

$rightf(r) = \sum_{q \in right(r)} df(q) + \sum_{q \in right^-(r)} \max\{leftf(q), rightf(q)\}$

QP decomposes term into small number of refstrings contained in t

Fuzzy Search with Edit Distance

Idea:

tolerate mis-spellings and other variations of search terms
and score matches based on editing distance

Examples:

- 1) query: Microsoft
fuzzy match: Migrosoft
score ~ edit distance 3
- 2) query: Microsoft
fuzzy match: Microsiphon
score ~ edit distance 5
- 3) query: Microsoft Corporation, Redmond, WA
fuzzy match at token level: MS Corp., Redmond, USA

Similarity Measures on Strings (1)

Hamming distance of strings $s_1, s_2 \in \Sigma^*$ with $|s_1|=|s_2|$:
number of different characters (cardinality of $\{i: s_{1_i} \neq s_{2_i}\}$)

Levenshtein distance (edit distance) of strings $s_1, s_2 \in \Sigma^*$:
minimal number of editing operations on s_1
(replacement, deletion, insertion of a character)
to change s_1 into s_2

For $\text{edit}(i, j)$: Levenshtein distance of $s_1[1..i]$ and $s_2[1..j]$ it holds:

$$\text{edit}(0, 0) = 0, \text{edit}(i, 0) = i, \text{edit}(0, j) = j$$

$$\text{edit}(i, j) = \min \{ \text{edit}(i-1, j) + 1, \\ \text{edit}(i, j-1) + 1, \\ \text{edit}(i-1, j-1) + \text{diff}(i, j) \}$$

with $\text{diff}(i, j) = 1$ if $s_{1_i} \neq s_{2_j}$, 0 otherwise

→ efficient computation by dynamic programming

Similarity Measures on Strings (2)

Damerau-Levenshtein distance of strings $s_1, s_2 \in \Sigma^*$:
minimal number of replacement, insertion, deletion, or
transposition operations (exchanging two adjacent characters)
for changing s_1 into s_2

For edit (i, j) : Damerau-Levenshtein distance of $s_1[1..i]$ and $s_2[1..j]$:

$$\text{edit}(0, 0) = 0, \text{edit}(i, 0) = i, \text{edit}(0, j) = j$$

$$\text{edit}(i, j) = \min \{ \text{edit}(i-1, j) + 1, \\ \text{edit}(i, j-1) + 1, \\ \text{edit}(i-1, j-1) + \text{diff}(i, j), \\ \text{edit}(i-2, j-2) + \text{diff}(i-1, j) + \text{diff}(i, j-1) + 1 \}$$

with $\text{diff}(i, j) = 1$ if $s_{1_i} \neq s_{2_j}$, 0 otherwise

Similarity based on N-Grams

Determine for string s the set of its N-Grams:

$$G(s) = \{\text{substrings of } s \text{ with length } N\}$$

(often trigrams are used, i.e. $N=3$)

Distance of strings s_1 and s_2 :

$$|G(s_1)| + |G(s_2)| - 2|G(s_1) \cap G(s_2)|$$

Example:

$$G(\text{rodney}) = \{\text{rod, odn, dne, ney}\}$$

$$G(\text{rhodnee}) = \{\text{rho, hod, odn, dne, nee}\}$$

$$\text{distance}(\text{rodney}, \text{rhodnee}) = 4 + 5 - 2 \cdot 2 = 5$$

Alternative similarity measures:

Jaccard coefficient: $|G(s_1) \cap G(s_2)| / |G(s_1) \cup G(s_2)|$

Dice coefficient: $2 |G(s_1) \cap G(s_2)| / (|G(s_1)| + |G(s_2)|)$

N-Gram Indexing for Fuzzy Search

Theorem (Jokinen and Ukkonen 1991):

**for query string s and a target string t ,
the Levenshtein edit distance is bounded by the N-Gram overlap:**

$$edit(s, t) \leq d \Rightarrow |Ngrams(s) \cap Ngrams(t)| \geq |s| - (N - 1) - dN$$

**→ for fuzzy-match queries with edit-distance tolerance d ,
perform top-k query over Ngrams,
using count for score aggregation**

Phonetic Similarity (1)

Soundex code:

Mapping of words (especially last names) onto 4-letter codes such that words that are similarly pronounced have the same code

- first position of code = first letter of word
- code positions 2, 3, 4 (a, e, i, o, u, y, h, w are generally ignored):

b, p, f, v	→ 1	c, s, g, j, k, q, x, z	→ 2
d, t	→ 3	l	→ 4
m, n	→ 5	r	→ 6
- Successive identical code letters are combined into one letter (unless separated by the letter h)

Examples:

Powers → P620 , Perez → P620

Penny → P500, Penee → P500

Tymczak → T522, Tanshik → T522

Phonetic Similarity (2)

Editex similarity:

edit distance with consideration of phonetic codes

For editex (i, j): Editex distance of s1[1..i] and s2[1..j] it holds:

$$\text{editex}(0, 0) = 0,$$

$$\text{editex}(i, 0) = \text{editex}(i-1, 0) + d(s1[i-1], s1[i]),$$

$$\text{editex}(0, j) = \text{editex}(0, j-1) + d(s2[j-1], s2[j]),$$

$$\text{editex}(i, j) = \min \{ \text{editex}(i-1, j) + d(s1[i-1], s1[i]), \\ \text{editex}(i, j-1) + d(s2[j-1], s2[j]), \\ \text{edit}(i-1, j-1) + \text{diffcode}(i, j) \}$$

with $\text{diffcode}(i, j) = 0$ if $s1_i = s2_j$,

1 if $\text{group}(s1_i) = \text{group}(s2_j)$, 2 otherwise

und $d(X, Y) = 1$ if $X \neq Y$ and X is h or w,

$\text{diffcode}(X, Y)$ otherwise

with group:

{a e i o u y}, {b p}, {c k q}, {d t}, {l r},

{m n}, {g j}, {f p v}, {s x z}, {c s z}

3.4 Index Organization and Advanced Query Types

Richer Functionality:

- **Boolean combinations of search conditions**
- **Search by word stems**
- **Phrase queries and proximity queries**
- **Wild-card queries**
- **Fuzzy search with edit distance**

Enhanced Performance:

- **Stopword elimination**
- **Static index pruning**
- **Duplicate elimination**

Stopword Elimination

Lookup in stopword list

(possibly considering domain-specific vocabulary,
e.g. „definition“ or „theorem“ in math corpus)

Typical English stopwords

(articles, prepositions, conjunctions, pronouns,
„overloaded“ verbs, etc.):

a, also, an, and, as, at, be, but, by,
can, could, do, for, from, go,
have, he, her, here, his, how,
I, if, in, into, it, its,
my, of, on, or, our, say, she,
that, the, their, there, therefore, they,
this, these, those, through, to, until,
we, what, when, where, which, while, who, with, would,
you, your

Static Index Pruning (Carmel et al. 2001)

Scoring function S' is an ε -variation of scoring function S if
 $(1-\varepsilon)S(d) \leq S'(d) \leq (1+\varepsilon)S(d)$ for all d

Scoring function S_q' for query q is (k, ε) -good for S_q if
there is an ε -variation S' of S_q such that
the top- k results for S_q' are the same as those for S' .

S_q' for query q is (δ, ε) -good for S_q if
there is an ε -variation S' of S_q such that
the top- δ results for S_q' are the same as those for S' ,
where top- δ results are all docs with score above $\delta \cdot \text{score}(\text{top-1})$

Given k and ε , prune index lists so as to guarantee (k, ε) -good results for all queries q with r terms where $r < 1/\varepsilon$.

→ for each index list L_i , let $s_{i(k)}$ be the rank- k score;
prune all L_i entries with score $< \varepsilon \cdot s_{i(k)}$

Efficiency and Effectiveness of Static Index Pruning

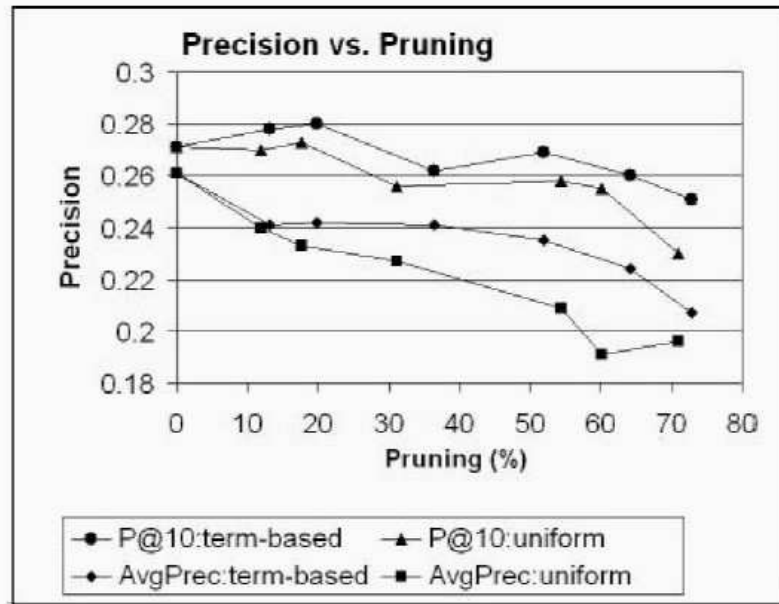


Figure 2: Precision of search results at varying levels of pruning.

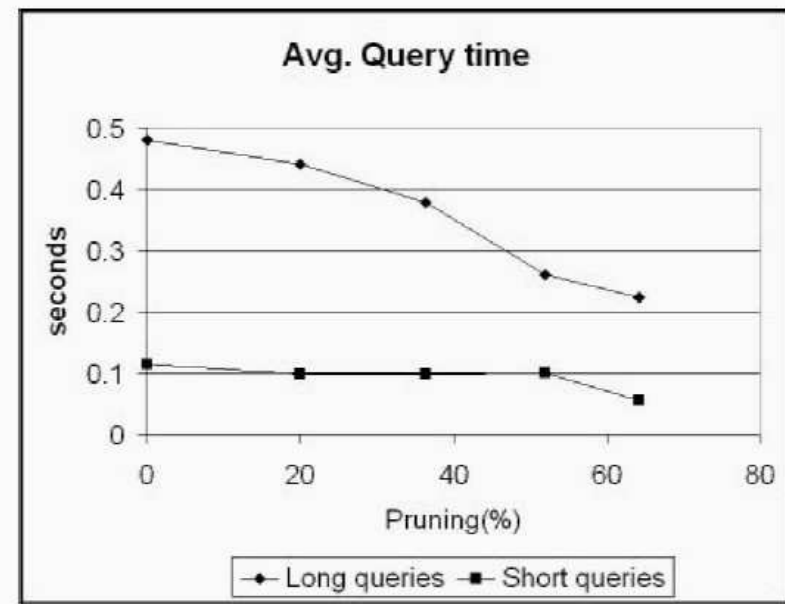


Figure 4: Average query processing time at varying levels of pruning.

from: D. Carmel et al., Static Index Pruning for Information Retrieval Systems, SIGIR 2001

Duplicate Elimination (Broder 1997)

**duplicates on the Web may be slightly perturbed
crawler & indexing interested in identifying near-duplicates**

Approach:

- **represent each document d as set (or sequence) of shingles (N-grams over tokens)**
- **encode shingles by hash fingerprints (e.g., using SHA-1), yielding set of numbers $S(d) \subseteq [1..n]$ with, e.g., $n=2^{64}$**
- **compare two docs d, d' that are suspected to be duplicates by**

- **resemblance:**
$$\frac{|S(d) \cap S(d')|}{|S(d) \cup S(d')|}$$

- **containment:**
$$\frac{|S(d) \cap S(d')|}{|S(d)|}$$

- **drop d' if resemblance or containment is above threshold**

Min-Wise Independent Permutations (MIPs)

set of ids

17	21	3	12	24	8
----	----	---	----	----	---

$h_1(x) = 7x + 3 \pmod{51}$

20	48	24	36	18	8
----	----	----	----	----	---

 $h_2(x) = 5x + 6 \pmod{51}$

40	9	21	15	24	46
----	---	----	----	----	----

 \vdots
 $h_N(x) = 3x + 9 \pmod{51}$

9	21	18	45	30	33
---	----	----	----	----	----

compute N random permutations with:

MIPs vector: minima of perm.

8
9
⋮
9

N

MIPs (set1)		MIPs (set2)
8	↔	8
9	↔	24
33	↔	45
24	↔	24
36	↔	48
9	↔	13

estimated resemblance = 2/6

$$P[\min\{\pi(x) | x \in S\} = \pi(x)] = 1/|S|$$

MIPs are unbiased estimator of resemblance:

$$P[\min\{h(x) | x \in A\} = \min\{h(y) | y \in B\}] = |A \cap B| / |A \cup B|$$

MIPs can be viewed as repeated sampling of x, y from A, B

Efficient Duplicate Detection in Large Corpora

avoid comparing all pairs of docs

Solution:

- 1) for each doc compute shingle-set and MIPs
- 2) produce (shingleID, docID) sorted list
- 3) produce (docID1, docID2, shingleCount) table with counters for common shingles
- 4) Identify (docID1, docID2) pairs with shingleCount above threshold and add (docID1, docID2) edge to graph
- 5) Compute connected components of graph (union-find)
→ these are the near-duplicate clusters

Trick for additional speedup of steps 2 and 3:

- compute super-shingles (meta sketches) for shingles of each doc
- docs with many common shingles have common super-shingle w.h.p.

Additional Literature for Chapter 3

Top-k Query Processing:

- Grossman/Frieder Chapter 5
- Witten/Moffat/Bell, Chapters 3-4
- A. Moffat, J. Zobel: Self-Indexing Inverted Files for Fast Text Retrieval, TOIS 14(4), 1996
- R. Fagin, A. Lotem, M. Naor: Optimal Aggregation Algorithms for Middleware, J. of Computer and System Sciences 66, 2003
- S. Nepal, M.V. Ramakrishna: Query Processing Issues in Image (Multimedia) Databases, ICDE 1999
- U. Guentzer, W.-T. Balke, W. Kiessling: Optimizing Multi-Feature Queries in Image Databases, VLDB 2000
- C. Buckley, A.F. Lewit: Optimization of Inverted Vector Searches, SIGIR 1985
- M. Theobald, G. Weikum, R. Schenkel: Top-k Query Processing with Probabilistic Guarantees, VLDB 2004
- M. Theobald, R. Schenkel, G. Weikum: Efficient and Self-Tuning Incremental Query Expansion for Top-k Query Processing, SIGIR 2005
- X. Long, T. Suel: Optimized Query Execution in Large Search Engines with Global Page Ordering, VLDB 2003
- A. Marian, N. Bruno, L. Gravano: Evaluating Top-k Queries over Web-Accessible Databases, TODS 29(2), 2004

Additional Literature for Chapter 3

Index Organization and Advanced Query Types:

- Manning/Raghavan/Schütze, Chapters 2-6, <http://informationretrieval.org/>
- H.E. Williams, J. Zobel, D. Bahle: Fast Phrase Querying with Combined Indexes, ACM TOIS 22(4), 2004
- WordNet: Lexical Database for the English Language, <http://wordnet.princeton.edu/>
- H.-J. Schek: The Reference String Indexing Method, ECI 1978
- D. Carmel, D. Cohen, R. Fagin, E. Farchi, M. Herscovici, Y.S. Maarek, A. Soffer: Static Index Pruning for Information Retrieval Systems, SIGIR 2001
- G. Navarro: A guided tour to approximate string matching, ACM Computing Surveys 33(1), 2001
- G. Navarro, R. Baeza-Yates, E. Sutinen, J. Tarhio: Indexing Methods for Approximate String Matching. IEEE Data Engineering Bulletin 24(4), 2001
- A.Z. Broder: On the Resemblance and Containment of Documents, Compression and Complexity of Sequences Conference 1997
- A.Z. Broder, M. Charikar, A.M. Frieze, M. Mitzenmacher: Min-Wise Independent Permutations, Journal of Computer and System Sciences 60, 2000