

# Efficient Top-K Query Processing for Text, Semistructured-, and Structured Data

Martin Theobald

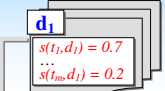
mtb@mpi-inf.mpg.de

1

## Threshold Algorithms for Top-K Retrieval

R.Fagin et al., PODS '01  
Balke et al., VLDB '00

Corpus:  $d_1, \dots, d_n$



Query:  $q = (t_1, t_2, t_3)$

INDEXING

$L_i$	d78	d23	d10	d1	d88	...
L1	d78	d23	d10	d1	d88	...
L2	d64	d23	d10	d10	d78	...
L3	d10	d78	d64	d99	d34	...

1.  $NRA(q, L)$ : scan all lists  $L_i$  ( $i = 1..m$ ) in parallel & consider doc  $d$  at pos;  $E(d) := E(d) \cup \{i\}$ ;
2.  $E(d) := E(d) \cup \{i\}$ ;
3.  $high = s_i(d)$ ;
4.  $worstscore(d) := \{s_i(d) \mid i \in E(d)\}$ ;
5.  $bestscore(d) := \text{worstscore}(d) + \sum \{high_i \mid i \in E(d)\}$ ;
6. if  $worstscore(d) > \text{min-k}$  then add  $d$  to top-k
7.  $\text{min-k} := \min\{\text{worstscore}(d') \mid d' \in \text{top-k}\}$ ;
8. else if  $bestscore(d) > \text{min-k}$  then candidates := candidates  $\cup \{d\}$ ;
9. if  $\max\{\text{bestscore}(d') \mid d' \in \text{candidates}\} \leq \text{min-k}$  then exit;

Rank	Doc	Worst-score	Best-score
1	d10	(2.1)	2.1
2	d78	1.4	2.0
3	d23	1.4	1.8
4	d64	1.2	2.0

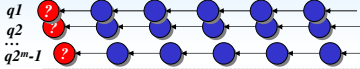
3

## Prob-k – Efficient Queue Management

Theobald, Schenkel, Weikum, VLDB'04

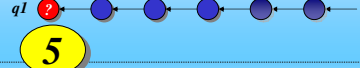
### Prob-conservative

- $2^m - 1$  queues per query & use bestscore( $d$ ) as priority
- Group candidates by remainder sets  $\{1..m\} - E(d)$
- Top candidate dominates all candidates in each queue
- Test top candidate only
- For all queues  $q$ : Drop queue  $q$ , if  $P[\text{top}(q) \text{ can qualify for top-k}] \leq \epsilon$



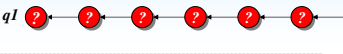
### Prob-smart

- 1 bounded queue per query & use bestscore( $d$ ) as priority
- Merge all candidates by their best-scores
- No dominating candidate in terms of score prediction
- Update & rebuild entire queue periodically
- Stop heuristically, if  $P[\text{top}(q) \text{ can qualify for top-k}] \leq \epsilon$



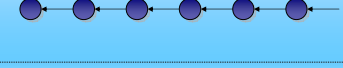
### Prob-progressive

- 1 queue per query & use bestscore( $d$ ) as priority
- Merge all candidates by their best-scores
- No dominating candidate wrt. score prediction
- Test all candidates periodically
- For all candidates  $d$  in  $q$ : Drop individual candidate  $d$ , if  $P[d \text{ can qualify for top-k}] \leq \epsilon$



### Prob-aggressive

- No queue
- Use efficient hash-joins only
- Consider virtual candidate  $d_v$  with  $E(d_v) = \emptyset$
- $d_v$  dominates all yet unseen candidates
- Stop heuristically, if  $P[d_v \text{ can qualify for top-k}] \leq \epsilon$



5

## TopX – Efficient Support for XML IR

Theobald, Schenkel, Weikum, VLDB'05

### Supports the XPath 2.0 Full-Text and NEXI query languages

```
//article[about(./[abskwd], "genetic algorithm")]  
//body[sec[about(simulated annealing, "fuzzy")]]
```

### Probabilistic candidate pruning

- Combined aggregated score- and XML-specific selectivity predictor

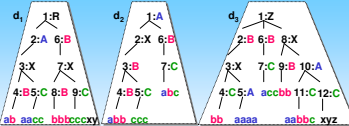
### XML-specific scoring model

- Okapi BM-25 extension for semistructured data using individual element statistics
- Optional redundant full-content indexing
- Accelerates descendants axis between tag/term pairs and/or tag paths

### Supports multiple index structures

- Pre/Post-Order or Data Guides
- Dynamic query rewriting wrt. to individual predicate selectivities

### Cost-based scheduling of random I/Os to test navigational query conditions



Example query:  $//A[/*a*/][B[/*b*/]C[/*c*/]]$

Tag	Term	DocId	Score	ElemId	Pre	Post
1	A	1	1.0	1	1	2
1	A	1	1.0	2	2	3
1	A	1	1.0	3	3	4
2	B	1	0.8	1	2	3
2	B	1	0.8	2	3	4
2	B	1	0.8	3	4	5
2	B	2	0.7	1	2	3
2	B	2	0.7	2	3	4
2	B	2	0.7	3	4	5
2	B	3	0.6	1	2	3
2	B	3	0.6	2	3	4
2	B	3	0.6	3	4	5
2	B	4	0.5	1	2	3
2	B	4	0.5	2	3	4
2	B	4	0.5	3	4	5
2	B	5	0.4	1	2	3
2	B	5	0.4	2	3	4
2	B	5	0.4	3	4	5
2	B	6	0.3	1	2	3
2	B	6	0.3	2	3	4
2	B	6	0.3	3	4	5
2	B	7	0.2	1	2	3
2	B	7	0.2	2	3	4
2	B	7	0.2	3	4	5
2	B	8	0.1	1	2	3
2	B	8	0.1	2	3	4
2	B	8	0.1	3	4	5
2	B	9	0.0	1	2	3
2	B	9	0.0	2	3	4
2	B	9	0.0	3	4	5

Sorted Block-Scans & Random Range-Scans

2

## Prob-k – Probabilistic Candidate Pruning

Theobald, Schenkel, Weikum, VLDB'04

### TA family of algorithms based on invariant:

$$\sum_{i \in E(d_j)} s_i(d_j) \leq \text{score}(q, d_j) \leq \sum_{i \in E(d_j)} s_i(d_j) + \sum_{i \in E(d_j)} \text{high}_i$$

- If  $\text{worstscore}(d_j) > \text{min-k} \rightarrow$  add  $d_j$  to top-k
- If  $\text{bestscore}(d_j) \leq \text{min-k} \rightarrow$  drop  $d_j$  from queue

### Probabilistic threshold test:

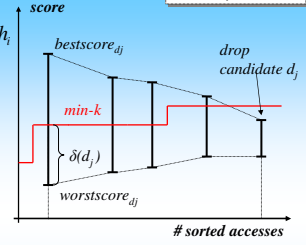
$$P(d_j) = P\left[\sum \{s_i \mid i \in E(d_j)\} > \delta(d_j)\right] < \epsilon$$

with  $\delta(d_j) := \text{min-k} - \sum \{s_i \mid i \in E(d_j)\}$

$\rightarrow$  Drop  $d_j$  from queue

### Score predictors for different distributions:

- Uniform in  $[0, \text{high}_j]$  or in  $[\text{low}_j, \text{high}_j]$
- Truncated Poisson over  $n_j$  equidistant values with  $P[d = v_j] = e^{-\alpha} \frac{\alpha^{j-1}}{(j-1)!}$
- Arbitrarily distributed and approximated by an equi-width histogram with cells  $j=1..c$  ranging over  $(l_j, h_j)$



## Self-tuning Incremental Query Expansions

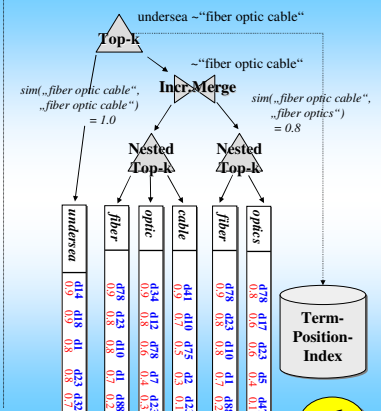
Theobald, Schenkel, Weikum, SigIR'05

### Incremental Merge operator

- Incrementally merge  $L_{j-1}, L_j$  in descending order of  $s_{ij}$  into a single virtual index list
- Seamless integration into top-k operator and single threshold termination
- Sharp semantics: "The best match out of  $c_1, \dots, c_k$ " instead of "Any or all of  $c_1, \dots, c_k$ "
- Modified score aggregation  $\text{score}(d_j) := \sum_{i=1}^m \max_{c_{ip} \in \text{exp}(c_j)} \{s(c_{ip}, d_j)\}$
- Incorporate additional similarity measures  $\rightarrow \text{sim}(c_i, c_j), s_{ij}$  (e.g., ontological similarities)
- Increased efficiency and retrieval robustness!

$\sim c = \{c_1, c_2, c_3\}$
$\text{sim}(c_1, c_1) = 1.0$ $c_1$ d78 d23 d10 d1 d88 0.9 0.8 0.8 0.7 0.2
$\text{sim}(c_1, c_2) = 0.9$ $c_2$ d64 d23 d10 d10 d78 0.8 0.6 0.6 0.2 0.1
$\text{sim}(c_1, c_3) = 0.5$ $c_3$ d10 d78 d64 d99 d34 0.7 0.5 0.4 0.2 0.1

### Nested Top-k operator



6

## Experiments

– Aquaint News Corpus: 528,155 docs; ~86,000,000 tuples; 1.9 GB  
– 50 "hard" queries taken from the TREC Robust track 2004

