
Top-k Query Processing with Probabilistic Guarantees

Martin Theobald

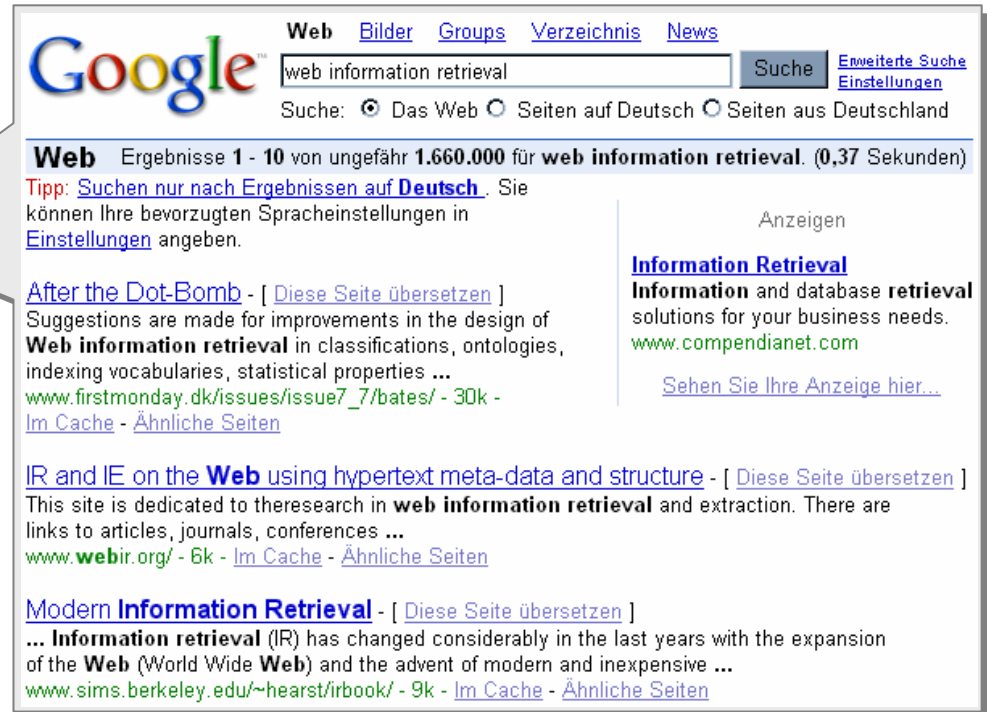
mtb@mpi-sb.mpg.de

International Max Planck Research School for Computer
Science, Saarbrücken



Motivation

- Focus of this work:
- Efficient top-k query evaluation in the style of current web search engines
- **Probabilistic score predictions for early candidate pruning**
- **Efficient queue management for candidate documents**



The screenshot shows a Google search interface. The search bar contains the text "web information retrieval". Below the search bar, there are navigation links for "Web", "Bilder", "Groups", "Verzeichnis", and "News". The search results are displayed in a list format. The first result is titled "After the Dot-Bomb" and includes a snippet about suggestions for improvements in the design of web information retrieval. The second result is titled "IR and IE on the Web using hypertext meta-data and structure" and includes a snippet about links to articles, journals, and conferences. The third result is titled "Modern Information Retrieval" and includes a snippet about the expansion of the World Wide Web and the advent of modern and inexpensive information retrieval technologies. On the right side of the search results, there is a section for "Anzeigen" (Ads) with a link to "Information Retrieval" solutions for business needs.

- *M. Theobald, G. Weikum, and R. Schenkel:*
“Top-k Query Evaluation with Probabilistic Guarantees.”
VLDB Conference 2004, Toronto, Canada



Moore's Law and the Web

“At our rate of technological development and advances in the semiconductor industry, the **complexity of integrated circuits doubles every 18 months.**”

[Moore '65]

Hardware	Annual Growth
CPU Performance	60%
Disk Storage Space	110%
Random Access Time of Dynamic RAM's	10%
Disk Random IO's (Latency)	8%
Data	
# Web Pages	~120%
Multimedia Data on the Web	~300%

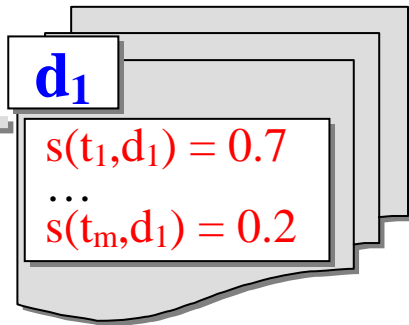
- **Google's** index today contains 4.285.199.774 documents in $\gg 1$ TB
 - *Less than 1s average response time*
 - *Server farm with 10,000+ units*
→ *massive redundancy*

→ *Efficient top-k retrieval using mainly sequential IO's to save hardware resources*



Fagin's NRA_[PODS '01] at a Glance

Corpus: d_1, \dots, d_n

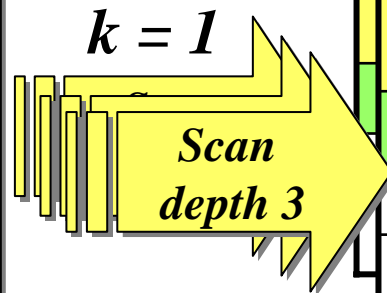


Query: $q = (t_1, t_2, t_3)$

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

Inverted Index

t_1	d78 0.9	d23 0.8	d10 0.8	d1 0.7	d88 0.2	...
t_2	d64 0.8	d23 0.6	d10 0.6	d10 0.2	d78 0.1	...
t_3	d10 0.7	d78 0.5	d64 0.4	d99 0.2	d34 0.1	...



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

STOP!

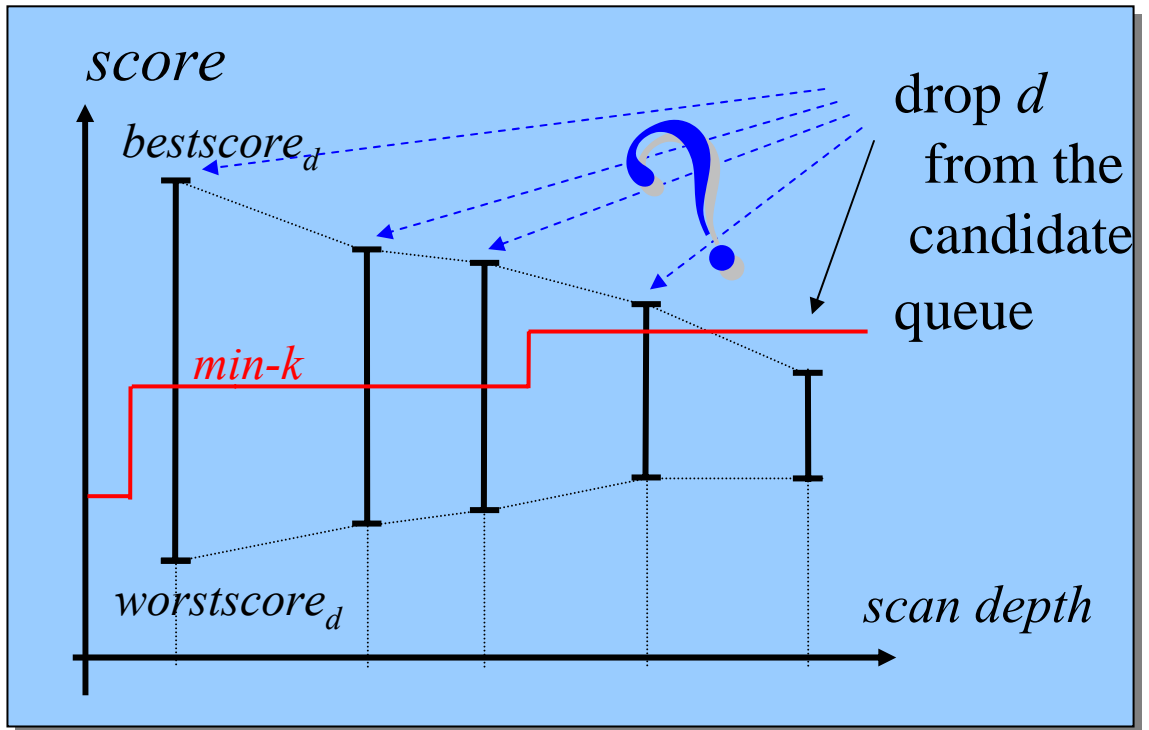


How can we be even faster?



Evolution of a Candidate's Score

- Worst- and best-scores slowly converge to final score
 - Add d to top- k result, if $worstscore_d > min-k$
 - Drop d only if $bestscore_d < min-k$, otherwise keep it in candidate queue
- **Overly conservative threshold & long sequential index scans**

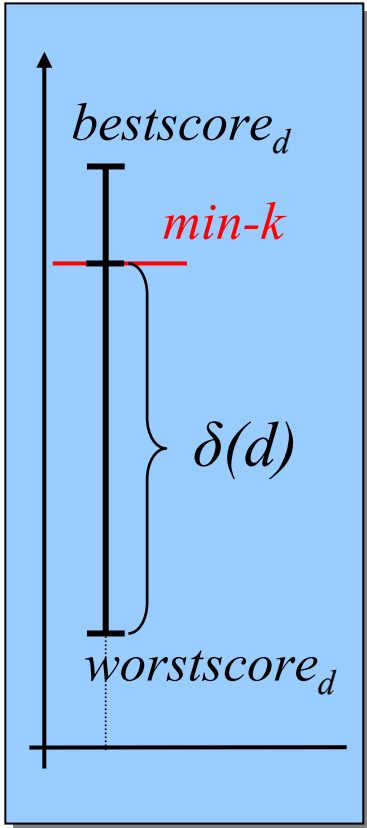


- **Approximate top-k**

“What is the **probability** that d qualifies for the top- k ?”



Safe Thresholding vs. Probabilistic Guarantees



- NRA based on invariant

$$\underbrace{\sum_{i \in E(d)} s_i(d)}_{worstscore_d} \leq s(d) \leq \underbrace{\sum_{i \in E(d)} s_i(d) + \sum_{i \notin E(d)} high_i}_{bestscore_d}$$

- Relaxed into **probabilistic threshold test**

$$p(d) := P \left[\sum_{i \in E(d)} s_i(d) + \sum_{i \notin E(d)} s_i(d) > \min_k \right] \leq \varepsilon$$

- Or equivalently, with $\delta(d) := \min_k - \sum \{s_i \mid i \in E(d)\}$

$$p(d) := P \left[\sum_{i \notin E(d)} s_i(d) > \delta(d) \right] \leq \varepsilon$$



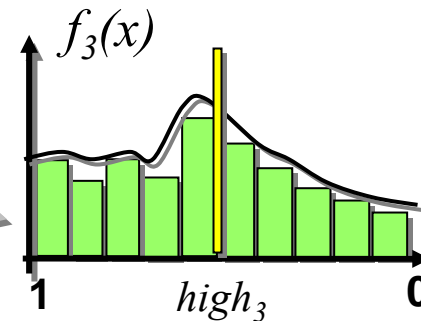
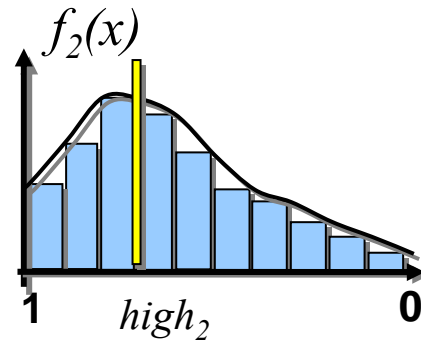
Probabilistic Score Predictions

How to estimate

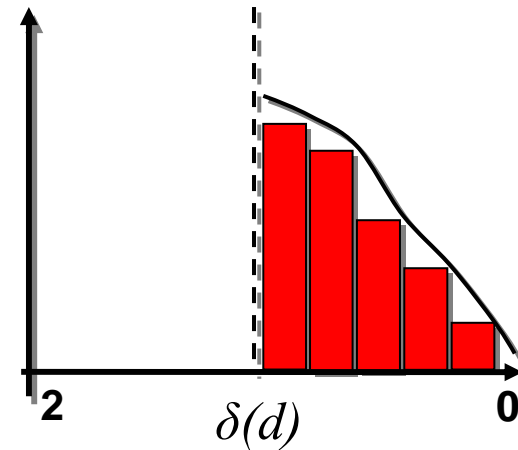
$$P \left[\sum_{i \notin E(d)} s_i(d) > \delta(d) \right]$$

Inverted Index

t_1	d78 0.9	d23 0.8	d10 0.8	d1 0.7	d88 0.2	...
t_2	d64 0.8	d23 0.6	d10 0.6	d11 0.2	d78 0.1	...
t_3	d1 0.7	d78 0.5	d64 0.4	d99 0.2	d34 0.1	...



Convolution
($f_2(x), f_3(x)$)



- Approximate each basic distribution
 - Poisson estimator** or **pre-computed histograms**
 - Compute convolution over all $f_i(x)$ where $i \notin E(d)$
- Moment-Generating Functions & **Chernoff-Hoeffding bounds** [Siegel '95]



Is It Worth all the Effort?

● Non-negligible **prediction overhead**

- $2^m - 1$ possible convolutions
- Expensive computation
- Frequent predictor updates

● **Queue management**

as a key role for query evaluation

- Which candidates are tested?
- How often is a candidate tested?
- What actions are taken when a candidate fails the test?



Conservative Queuing

● *Prob-conservative*

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

● *Prob-progressive*

- 1 queue per query
- Merge all candidates by their best-scores
- No dominating candidate in terms of score prediction
- Test all candidates periodically
- For all candidates d in q :
Drop candidate d , if
 $P[d \text{ can qualify for top-}k] \leq \varepsilon$

- Stop by safe min- k threshold test
or when all queues are empty



Aggressive Queuing

● *Prob-smart*

- *1 bounded queue* per query
- Merge all candidates by their best-scores
- No dominating candidate in terms of score prediction
- Update & rebuild entire queue periodically

● *Stop heuristically, if*

$$P[\text{top}(q) \text{ can qualify for top-}k] \leq \varepsilon$$

● *Prob-aggressive*

- *No queue*
- Consider virtual candidate d_v with $E(d_v) = \emptyset$
- d_v dominates all yet unseen candidates

● *Stop heuristically, if*

$$P[d_v \text{ can qualify for top-}k] \leq \varepsilon$$



Experimental Setup

● Gov

- TREC-12 Web Track's .Gov collection
- 1.250.000 web documents (html, doc, pdf)
- 50 keyword queries from the Topic Distillation task, $m \leq 5$
- e.g. *“legalization marijuana”*

● XGov

- Gov with manual query expansion, $m \leq 20$
- e.g. *“legalization law marijuana cannabis drug abuse pot ...”*

● IMDB

- Structured (XML) version of the Internet Movie Data Base
- 375.000 movie files, 1.200.000 person files
- Mixed text and categorical attributes: Genre, Actor, Description
- e.g. *“Genre \supseteq {Western} \wedge Actor \supseteq {John Wayne, Katherine Hepburn} \wedge Description \supseteq {Sheriff, Marshall}”*



Baseline

$k=20$ & $\varepsilon = 0.1$

Gov

	# sorted accesses	execution time (sec.)	max. queue size	macro-avg. precision
<i>original NRA</i>	2,263,652	148.7	10,849	1.0
<i>Prob-con</i>	993,414	25.6	29,207	0.87
<i>Prob-pro</i>	1,659,706	44.2	6,551	0.87
<i>Prob-smart</i>	527,980	15.9	400	0.69
<i>Prob-agg</i>	20,435	0.6	0	0.42

XGov

<i>original NRA</i>	22,403,490	7,908	70,896	1.0
<i>Prob-con</i>	10,165,677	6,448	51,893	0.90
<i>Prob-pro</i>	20,006,283	1,791	12,435	0.95
<i>Prob-smart</i>	18,287,636	1,066	400	0.88
<i>Prob-agg</i>	133,745	2	0	0.35

IMDB

<i>original NRA</i>	1,003,650	201.9	12,628	1.0
<i>Prob-con</i>	463,562	17.8	14,990	0.71
<i>Prob-pro</i>	490,041	69.0	9,173	0.75
<i>Prob-smart</i>	403,981	12.7	400	0.54
<i>Prob-agg</i>	41,821	0.7	0	0.18



Conclusions & Ongoing Work

● Performance vs. Quality

- Speedup factor of up to **10** at ~ **80%** prec. (*Prob-smart*)
- Speedup factor of more than **100** at ~ **40%** prec. (*Prob-agg*)

● Semantic extensions

- Query-specific weights
- Efficient query expansions

● TREC (Text REetrieval Conference) - benchmark competition

- Robust Track – hard queries
- Web Track – *tf·idf* and global *PageRank* scores
- Terabyte Track – 480 GB web documents

● XML retrieval (XPath)

- Structural query conditions, path similarities & no predefined retrieval unit

