



Web Science - Investigating the Future of Information and Communication Why Search the Past?

- Historical information needs (an analyst working on the social reactions on the net after 9/11)
- To find relevant resources that do not exist anymore
- To discover trends, opinions, etc. for a certain time period (what people think about UFOs at the beginning of millenium?)







Web Science - Investigating the Future of Information and Communication Time-travel Queries

- "Queries that combine the content and temporal predicates" [Berberich et al., SIGIR 2007]
- Interesting query types
 - Point-in-time
 - "euro 2012 articles" @ 1/July/2012
 - Interval
 - "euro 2012 articles" between 01.06-01.07 2012
 - Durability in a time interval [Hou U et al., SIGMOD 2010]:
 - "search engine research papers" that are in top-10 results for
 75% of the time between years 2000 and 2012



Web Science - Investigating the Future of Information and Communication Web Science - Investigating the Future of Information and Communication For an operation of the current of the current of the current of the current version of the current version, denoted by the current version, denoted by the current of the current version, denoted by the current version, denoted by

[Berberich et al., SIGIR 2007









e WIKI UKGOV # Postings Ratio # Postings Ratio - 8,647.996,223 100.00% 7,888,560,482 100.00% 0.00 7,769,776,831 89.84% 2,926,731,708 37.10% 0.01 1,616,014,825 18.69% 744,438,831 9,44% 0.05 556,204,068 6.43% 259,947,199 3.00% 0.10 379,962,802 4.39% 187,387,342 2.38% 0.25 252,581,230 2.92% 158,107,198 2.00% 0.50 203,269,464 2.35% 155,434,617 1.97%	Berger		Web Science – Inve It	wor	he Future of Info	rmation and (Communicatio
ϵ # Postings Ratio # Postings Ratio - 8.647.996.223 100.00% 7.888.560.482 100.00% 0.00 7.769.776.831 89.84% 2.926.731.708 37.10% 0.01 1.616.014.825 18.69% 744.438.831 9.44% 0.05 556.204.068 6.43% 259.947.199 3.30% 0.10 379.962.802 4.39% 187.387.342 2.38% 0.25 252.581.230 2.92% 158,107.198 2.00% 0.50 203.269.464 2.35% 155.434.617 1.97%			WIK	[]	UKGO	ov	
- 8,647,996,223 100.00% 7,888,560,482 100.00% 0.00 7,769,776,831 89,84% 2,926,731,708 37,10% 0.01 1,616,014,825 18,69% 744,438,831 9.44% 0.05 556,204,068 6.43% 259,947,199 3.30% 0.10 379,962,802 4.39% 187,373,422 2.38% 0.25 252,581,230 2.92% 158,107,198 2.00% 0.50 203,269,464 2.35% 155,434,617 1.97%		ϵ	# Postings	Ratio	# Postings	Ratio	
0.00 7,769,776,831 89.84% 2,926,731,708 37.10% 0.01 1,616,014,825 18.69% 274,438,831 9.44% 0.05 556,204,068 6.43% 259,947,199 3.30% 0.10 379,962,802 4.39% 157,387,342 2.38% 0.25 252,581,230 2.92% 158,107,198 2.00% 0.50 203,269,464 2.35% 155,434,617 1.97%		-	8,647,996,223	100.00%	7,888,560,482	100.00%	
0.01 1,616,014,825 18,69% 744,438,831 9,44% 0.05 556,224,068 6,439% 187,387,342 2.38% 0.25 252,581,230 2.92% 158,107,198 2.00% 0.50 203,269,464 2.35% 155,434,617 1.97%		0.00	7,769,776,831	89.84%	2,926,731,708	37.10%	
U.05 320,204,008 0.43% 229,947,199 3.30% 379,962,802 4.39% 187,387,342 2.38% 0.25 252,581,230 2.92% 158,107,198 2.00% 0.50 203,269,464 2.35% 155,434,617 1.97% 1 1 1 1 1 1 1 1 1		0.01	1,616,014,825	18.69%	744,438,831	9.44%	
0.10 379,902,802 4.39% 157,357,342 2.38% 0.25 (52,581,230 2.92% 158,107,198 2.00% 0.50 252,581,230 2.92% 158,434,617 1.97% 15,434,617 1.97% 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1		0.05	556,204,068	6.43%	259,947,199	3.30%	
0.25 2/2,351,250 2/2,357,250 195,107,155 2/0070 0.50 2/2,351,250 2/3,269,464 2.35% 155,434,617 1.97%		0.10	379,962,802	4.39%	187,387,342	2.38%	
1 1		0.25	203 269 464	2.92% 2.35%	155 434 617	2.00%	
		Rel	Relative Recall @ 10 (WIK) -+ Kendall's @ 10 (WIK) -+ Kendall's @ 10 (WKK) ->	1 + +	xxa IIx 0]	Relative Recal Kendall 166 Relative Recall	
					[Be	erberich et	al., SIGIR 2









- Given a space-bound of 3 (i.e., 3x of the optimal space), close-to-optimal performance is achievable!
 - (Reminder: Partitioning is not very cheap!)

Web Science - Investigating the Future of Information and Communication
 Time-interval Queries?
 When more than one partitions should be accessed, wasted I/O due to repeated postings! (e.g., 3x more postings in SB!)
 Example

Web Science - Investigating the Future of Information and Communication Time-point Queries: Result

[Berberich et al., SIGIR 2007]

- Time-point queries can be handled like this:
- Example



Web Science - Investigating the Future of Information and Communication Solution Approaches

Partition selection [Anand et al., CIKM 2010] Can we avoid accessing all partitions related to a given query?

Document partitioning (sharding) [Anand et al., SIGIR 2011] Can we partition postings in a list in a different way i.e., other than using time information?









- At each step:
 - Select the partition with highest B(P) / C(P)
 - Update benefits of the unselected partitions

[Anand et al., CIKM 2010]







- No of new documents that appear in the intersection of the partitions in x
- Modify algorithm to update benefits and costs after each picking a tuple **x**.

[Anand et al., CIKM 2010]











- For each partition of each term:
 - create a flat file storing KMV synopses (5% and 10% samples)

Index	UKGOV	NYT	WIKI
Fixed-7	11G	13G	13G
Syno sis Index - 5% sam le	146MB	134MB	146MB
Syno sis Index - 10% sam le	291MB	258MB	290MB
Fixed-30	4.4G	3.5G	6.3G
Syno sis Index - 5% sam le	61MB	39MB	75MB
Syno sis Index - 10% sam le	122MB	74MB	149MB









[Anand et al., SIGIR 2011]





Web Science - Investigating the Future of Information and Communication Performance

- Sharding improves query processing times w.r.t. temporal partitioning or no partitioning et al.
- Merging shards results shorter QP times
- Time measurements are taken on warm caches!!!

[Anand et al., SIGIR 2011]

Web Science - Investigating the Future of Information and Communication

- Work from Polytechnique Ins. of NYU
- Focus on the index size

Approaches up to now consider each version of a document separately: no special attention on the overlap between versions





- Key ideas
 - Small integers can be represented with smaller codes
 - Doc ids are not so small: instead, compress the gaps between the ids
 - Term frequencies are already small
 - Example









 Reduces the number of postings but increases the document space! (theoretically, up to N² virtual documents!)

```
[He et al., CIKM 2009]
```























	,
Partition-focused	Size-focused
– Time information kept	 – Time information kept
in postings	separately
- Redundancy solved	 Redundancy solved by
by partitioning	2-level compression
 Process only relevant 	 Process compact blocks
partitions	Hierarchical partitions
- Lossy compression of	– Lossless compression
versions (coalescing)	of versions









Web Science - Investigating the Future of Information and Communication Searching the past

- Problem statements
 - Time must be explicitly modeled in order to increase the effectiveness
 - Time uncertainty should be taken into account
 Two temporal expressions can refer to *the same* time period even though they are *not equally written*
- Example
 - Given the query "**Independence Day 2011**", a retrieval model relying on term-matching will fail to retrieve documents mentioning "July 4, 2011"





Web Science - Investigating the Future of Information and Communication Current approach

- Previous time-aware ranking methods follow two main approaches
 - 1. Mixture model: linearly combining *textual* and *temporal* similarity
 - 2. Probabilistic model: generating a query from the *textual part* and *temporal part* of a document independently

[Kanhabua et al., SIGIR 2011]



























• 4	Web Science - Releva 12 future-relate	ance judgr ance judgr	f Information and Commun Ments	nication
	POLITICS	ENVIRONMENT	SPACE	
	president election	global warming	Mars	
	Iraq war	energy efficiency	Moon	
	SCIENCE	PHYSICS	HEALTH	
	earthquake	particle Physics	bird flue	
	tsunami	Big Bang	influenza	
	BUSINESS	SPORT	TECHNOLOGY	
	subprime	Olympics	Internet	
	financial crisis	World cup	search engine	
		ſŀ	Kanhabua et al., SIGI	3 2011a



- Results
 - Topic features play an important role in ranking
 - Features in top-5 features with *lowest weights* are entity-based features
- Open issues
 - Extract predictions from other sources, e.g., Wikipedia, blogs, comments, etc.
 - Sentiment analysis for future-related information

[Kanhabua et al., SIGIR 2011a]



- Yue et al., SIGIR 2007] Visong Yue, Thomas Finley, Filip Radinski, Thorsten Joachims: A support vector method for optimizing average precision. SIGIR 2007: 271-278
- [Zhang, ICML 2004] Tong Zhang: Solving large scale linear prediction problems using stochastic gradient descent algorithms. ICML 2004

