

Caching query-biased snippets for efficient retrieval

Salvatore Orlando⁺, <u>Raffaele Perego</u>^{*}, Fabrizio Silvestri^{*} *ISTI - CNR, Pisa, Italy ⁺Università Ca' Foscari Venezia, Italy

Diego Ceccarelli, Claudio Lucchese, Salvatore Orlando, Raffaele Perego, Fabrizio Silvestri: Caching query-biased snippets for efficient retrieval. EDBT 2011: 93-104

Search Engine Results Page

	Web Images Video Local Shopping News More	
	edbt2011	lon-k re
	Search: () the Web () only in UK () only in Ireland	for t
TITLE	EDBT/ICDT 2011	
SNIPPET	March 21-25, 2011 Uppsala, Sweden. EDBT: 14th International Conference Extending Database Technology, March 22-24 The EDBT series of conference	e on es is an
	established	

Dingming Wu

edbticdt2011.it.uu.se - Cached

URL

I was visiting The Hong Kong Polytechnic University collaborating with Man Lung Yiu from March 1st to August 31st in 2010. I was an external reviewer of **EDBT2011**, VLDB J ... www.cs.aau.dk/~dingming - Cached

BEWEB 2011 (BEWEBworkshop) on Twitter

BEWEB 2011 will be held on March 25 in conjuntion with EDBT/ICDT 2011 conf. in Uppsala (Sweden) http://bit.ly/9U188T #edbt2011 #beweb2011 6:22 AM Nov 15th via web twitter.com/BEWEBworkshop - Cached

EDBT 2011 : 14th International Conference on Extending ...

https://cmt.research.microsoft.com/EDBT2011 Publication and IP issues Papers will be made available in online conference proceedings, possibly hosted within the ACM ... www.wikicfp.com/cfp/servlet/event.showcfp?eventid=10620&... - Cached

sults relevant he query

Search Engine Results Page

	Web Images Video Local Shopping News More -	
	edbt2011	
	Search: 💿 the Web 🔘 only in UK 🔘 only in Ireland	
	EDBT/ICDT 2011	
SNIPPEI	March 21-25, 2011 Uppsala, Sweden. EDBT: 14th International Conference on Extending Database Technology, March 22-24 The EDBT series of conferences is an	
	established edbticdt2011.it.uu.se - Cached	

Very important to estimate relevance of the result for the user's query

- high-quality, query-biased
- expensive to generate



Different queries....

en.wikipedia.org/wiki/Uppsala University Hospital - 59k - Cache

	Web Immagini Video Pagine Gialle Shopping Notizie	Altro -	Web Immagini Video Pagine Gialle
	university of Uppsala		Uppsala university science
l			
	Cerca: 💿 nel Web 🔘 nei siti in italiano		Cerca: • nel Web • nei siti in italiano
	Uppsala University Sweden - Traduci Uppsala University is one of northern Europe's most highly ranked univers Uppsala University pursues research in a wide range of fields that helps u www.uu.se/en - Cache Research - Uppsala universitet - Traduci Uppsala University pursues research in a wide range of fields that he Hallberg, Vice Chancellor of Uppsala University. Breadth, diversity, a www.uu.se/en/node4 - Cache	sities JS Ips us and	Research - Uppsala universitet - Traduci Uppsala University pursues research in a wide Stockholm-Uppsala Life Science. Codex - rules www.uu.se/en/node4 - Cache Master Programme in Computer Sci University Sweden - Traduci I saw that Uppsala was the oldest universi Uppsala University), students will study Ac www.uu.se/en/node605?pKod=TDV2M&lass
	Uppsala University - Wikipedia, the free encyclopedia - Tradu History Administration and organisation Campus Student life The university rose to pronounced significance during the rise of Sweden as a great power at the end of the 16th century and was then given a relative financial stability with the large	UPPSALA UNIVERSITET	University of Uppsala: Information from Uppsala, University of, at Uppsala, Sweden; for Shopping. Sports. Technology. Travel & Places. www.answers.com/topic/uppsala-university - 2
	en.wikipedia.org/wiki/Uppsala_University - 159k - Cache University of Uppsala: Information from Answers.com - Tradu Uppsala, University of, at Uppsala, Sweden; founded 1477 by Sten Sture and Archbishop Jakob Ulvsson University of Uppsala. Home > Library www.answers.com/topic/uppsala-university - 222k - Cache	ci a, the Elder, >	International Science Programs (ISP), U <u>Traduci</u> International Science Programme, Uppsala Uni Medical Cell Biology, Uppsala University (Prof. www.sasnet.lu.se/ispuppsala.html - <u>Cache</u>
	Uppsala universitetsbibliotek - Traduci Uppsala University Library Catalogue " New library catalogue " Test: M European Environmental Policy " More news items Uppsala University www.ub.uu.se/eindex.cfm - 119k - <u>Cache</u>	lanual of Library	Uppsala University - Wikipedia, the free History Administration and organisation Car life The university rose to pronounced significance Sweden as a great power at the end of the 16th then given a relative financial stability with the la
	History Present facilities Organization Uppsala Care Uppsala University Hospital in Uppsala, Sweden is a teaching hospital for the Uppsala University Faculty of Medicine and the Nursing School, Uppsala University Hospital is owned and		UU/IT/Research in Computing Science Publications and reports from the Computing Sc Uppsala University. Department of Information

Shopping Notizie Altro -

wide range of fields that helps us ... rules and regulations for research ...

r Science 2010/2011 - Uppsala

iversity in the Nordic countries and had ... (at dy Advanced Computer Science Studies in ... 1&lasar=10/11 - <u>Cache</u>

from Answers.com - Traduci den; founded 1477 by Sten Sture, ... Science. laces. Q & A. University of Uppsala ... sity - 222k - Cache

SP), Uppsala University, Sweden. -

a University ... Swedish collaborator: Dept. of (Prof. Erik Gylfe) ...

ne free encyclopedia - Traduci

ance during the rise of 16th century and was the large... y - 159k - <u>Cache</u>



ence - Traduci

www.csd.uu.se - Cache

ng Science Division. ... Copyright © 2011 nation Technology. ...

Same title, same URL, same snippet (A)



en.wikipedia.org/wiki/Uppsala University Hospital - 59k - Cache

www.csd.uu.se - Cache

Shopping

Notizie Altro -



Same title, same URL, different snippets (B)



Shopping Notizie Altro -



Web Searching process



WSE Front End caching

Occurrences of the most popular queries



WSE Front End caching



SERPs' caching does not help in the cases A and B above, when the same or a similar snippet is generated for the same doc, but for slightly different queries



Doc Repository Caching strategies



Doc Repository Caching

Given q and docID, the DR has to: 1) retrieve document content and URL, 2) generate an effective query-biased snippet

- A DR cache can substantially reduce disk accesses
- Two cache organizations discussed in literature:
 - DR_{Cachedoc} storing integral copies of most accessed documents*
 - DR_{Cachesurr} storing surrogates of most accessed documents**

We propose DR_{CacheSsnip} whose entries store supersnippets built on the basis of the past queries submitted to the WSE

*Andrew Turpin, Yohannes Tsegay, David Hawking, and Hugh E. Williams. Fast generation of result snippets in web search. SIGIR '07 **Y. Tsegay, S. J. Puglisi, A. Turpin, and J. Zobel, Document Compaction for Efficient Query Biased Snippet Generation. ECIR '09

The Query Log used

MSN RFP 2006 query log

- queries from a US Microsoft search site sampled over one month (May 2006)
- submitted to Yahoo! BOSS to retrieve top-10 results

	Query Log	
 9,000,000 queries 4,447,444 distinct queries 25,897,247 distinct URLS -> high URLs sharing! 	D1 training test	1
	D2 training test	94

Queries	Distinct Queries	Distinct Urls (top 10)
1,000,000	601,369	4,317,603
500,000	$310,\!495$	$2,\!324,\!295$
8		
9,000,000	4,447,444	25,897,247
4,500,000	$2,\!373,\!227$	$12,\!134,\!453$

Fact 1: URLs occurrences in SERPs

Double Pareto distribution





Fact 2: top-1000 most frequently retrieved documents

- snippets generated for the same document are a few (<9)</p>
- a document may answer several different queries



Evidences

- 1. Relevant documents are characterized by a few snippets
- 2. Relevant documents are shared by different queries
- 3. Snippets are made of a few sentences
- 4. Different snippets from the same document share common sentences

Q: an Elephant with a Rhino? A: I haven't a goddamn clue.



Definition of QL-based supersnippet

Given the set Q_d of all the past queries in QL for which document d was returned, we define supersnippet ss_n of d the set of the n most frequent sentences occurring in the snippets $S_{d,a}$ generated for answering the queries in Q_d

Claim

Caching LRU supersnippets ss_n is an effective and efficient strategy for DR_{Cache}

Effectiveness: DR_{CacheSsnip} vs. Yahoo! BOSS



Efficiency

	FECache	DRCache	Size (in MB)	Hit Ratio DRCache
	#Entries			
	(Hit Ratio)			
		1	256M	0.083
	SERPs cache	dee	512M	0.1
	filtering out most	aoc	1024M	0.134
			2048M	0.17
	recently submitted		256M	0.104
	queries	eaurm	512M	0.135
		surr	1024M	0.177
0.001			2048M	0.23
			256M	0.139
20 1e-04 -	256K (0.38)	<u> </u>	512M	0.194
cy article	2501 (0.56)	335	1024M	0.266
6 1e-05			2048M	0.34
	—		256M	0.13
10-06		8810	512M	0.189
1e-08 1e-07 1e-06 1e-05 1e-04 0.001 0. Query rank before cache	01 0.1 1	3310	1024M	0.266
			2048M	0.334

SS Cache peculiarities

CacheLookup(q, DocID) if InCache(DocID) then SS <- RetrieveSS(DocID) snippet <- GenerateS(q,SS)</pre> if Quality(q,SS) < Threshold then snippet <- UpdateSS (SS, DocID, q)</pre> UpdateLRU(DocID) else SS < - nullsnippet <- UpdateSS (SS, DocID, q)

```
ReplaceLRU(SS, DocID)
```

SS Cache peculiarities

CacheLookup(q, DocID)

if InCache(DocID) then

SS <- RetrieveSS(DocID)

snippet <- GenerateS(q,SS)</pre>

if Quality(q,SS) < Threshold then

snippet <- UpdateSS (SS, DocID, q)</pre>

```
UpdateLRU(DocID)
```

else

```
SS < - null
```

```
snippet <- UpdateSS (SS, DocID, q)
```

```
ReplaceLRU(SS, DocID)
```



A "Quality" Miss may occur when the quality of the snippet generated from the available SS is below a given threshold

Quality Miss count

- We count a Quality Miss when even a single query term is not present in the snippet
 - Also Yahoo! BOSS snippets sometimes do not contain ALL query terms
 - misspellings, disjunctive queries, etc.
- We underestimate the Hit Ratio

What about considering a Quality Miss only when our snippets' quality is lower than Yahoo!'s one?

Refined Efficiency Measurement



Total Hit Ratio: ~62%

tio	DRCache
	0.21
	0.29
	0.36
	0.42
	0.2
	0.28
	0.356
	0.419
	0.19



Conclusion & Future Work

- Novel technique for scaling up search engine performance by snippet caching
 - motivated by the analysis of a large real-world query log
 - enabling the generation of effective snippets even for queries and URLs not previously seen
 - achieving high Hit Ratios (cumulatively up to 0.62!)
 - evaluated by means of large-scale experiments
- Open issues
 - study of dynamic vs static policies for managing supersnippets
 - quality of SSs suffers for aging?
 - study of diversification in supersnippets?
 - heuristics, machine learning approaches for managing Quality Misses

