

Jamie Callan Language Technologies Institute Carnegie Mellon University callan+@cs.cmu.edu



Case 1: Large Commercial Information Provider

Task: Provide rapid access to several terabytes of information

3

- Difficulty: Long queries, many Web-search optimizations not applicable
- Each information source is stored separately
 - E.g., potentially hundreds or thousands of databases
- Simplifies maintenance, increases speed
- · Searching all databases is computationally expensive
- Missing relevant databases is costly
- Customers don't find what they need
- One organization, many distinct databases
- How does a customer know which to search?



© 2002, Jamie Callar



Case Summaries

- 1. Large Information Provider (West)
- · Different information stored in many different databases
- Each database managed by the same organization
- 2. Customer Support Bureau
- Information stored in a few different databases
- · Each database managed by the same organization
- 3. Competitive Analysis
- · Different information stored in many different places

7

- Each database managed independently
- 4. Distributed Wireless Networks
- · User has no access to centralized resources

Multi-Database Solutions: The Browsing Model (#1)



© 2002, Jamie Callar











Single Site / LAN / Few Sites: Relevant Document Distribution (RDD)

Assumption: New queries often look a lot like old queries Approach:

- Find a few old queries that are similar to the new query

 similarity is determined by vector similarity of query term
- Look up the effectiveness of each database for these old queries
- Compute an estimated effectiveness of each database for the new query as the average of its effectiveness for the similar old queries

This method is designed for environments where there is little variety in queries and new collections are not added often

(Voorhees, et al., 1995) 15



© 2002, Jamie Calla





Directory Model: Summary of Techniques for a Few Sites

19

- Single Site / LAN / Few Sites
- Relevant document distribution (RDD)
- Query Clustering

Common Assumptions:

- · Few collections
- Relatively stable collections,
- · Relatively stable information needs (queries)
- Considerable overlap in query contents
- · Good training data

Directory Model : Techniques for Many Sites

- Many collections (hundreds, thousands, tens of thousands, ...)
 Perhaps dispersed widely, managed independently, expensive, ...
- Collections are homogeneous and heterogeneous
- Queries are heterogeneous and can't be predicted in advance
- Training data is not practical
 - Many collections
 - Collections constantly being added, deleted, changed
- Solution:
- Content-based database selection (gGlOSS, CORI, CVV, Napster, ...)

20

© 2002, Jamie Callan

Page 5





1987 WSJ (132 MB)	1991 Patent (254 MB)	1989 AP (267 M	(B)
stobb (1)	sto (1)	sto (7)	<u> </u>
stochast (1)	stochast (21)	sto1 (4)	Word + frequenc
stock (46704)	stochiometr (1)	sto3 (1)	models are
stockad (5)	stociometr (1)	stoaker (1)	most common
stockard (3)	stock (1910)	stoand (1)	(gGlOSS, CORI,
stockbridg (2)	stockbarg (30)	stober (6)	language models
stockbrok (351)	stocker (211)	stocholm (1)	C v v,)
stockbrokag (1)	stockholm (1)	stock (28505)	
stockbrokerag (101)	stockigt (4)	stock' (6)	
stockdal (8)	stockmast (3)	stockad (35)	
stockhold (970)	stockpil (7)	stockard (12)	







T	A	Tana	A	T	A	Tana	A
Term	AVG U	Term	AVG tr	Term	Avg tr	lerm	AVG
project	0.750	visual	5.273	articles	4.121	alaiog	3.515
excel	0.750	Deta	4.986	setup	4.094	command	3.504
office	8.565	service	4.983	maii	4.067	following	3.387
works	7.389	basic	4.903	users	4.042	windows	3.369
server	7.271	file	4.867	set	3.948	new	3.369
word	7.221	nt	4.845	application	3.919	settings	3.317
table	6.639	field	4.729	product	3.890	example	3.152
printer	6.507	access	4.554	menu	3.840	version	3.147
foxpro	6.486	print	4.550	text	3.717	message	3.119
database	6.117	data	4.322	software	3.621	information	3.076
microsoft	5.736	internet	4.268	code	3.617	select	3.072
object	5.637	error	4.217	name	3.611		
user	5.297	box	4.213	system	3 544		

Directory Model: Evaluation						
Query 63: Iden	ntify a machine translation	system being	developed o	r		
marketed in any	country. Identify the deve	loper or vend	lor, name the	Aucznar		
system, and ide	itily one or more leatures of	or the system.	F	Key		
Rank	Collection	Score	RelDocs	9		
1	1989-90 Ziff Davis	0.571	153			
2	1989-90 Ziff Davis	0.569	34			
3	1991-92 Ziff Davis	0.563	77			
4	1991 Patent	0.407	0			
5	1989 AP	0.404	1			
6	1988 AP	0.401	4			
7	1988 WSJ	0.399	1			
8	1987 WSJ	0.399	0			
9	1988 Federal Register	0.341	0			
10	1989 Federal Register	0.337	0			
	28			© 2002, Jamie Callan		



avano, et al., 199 29

© 2002, Jamie Calla

Techniques for Many Sites: CORI (Inference Networks)

Assumption: The inference net retrieval model can apply to collections (represent each collection as a "big document")

Approach:

- Represent each collection by an inference net (standard inference net model)
- · Standard tf.idf similarity measure of query to each collection
- · Rank collections by their similarity to each query

Automatic creation, consistent, dynamic grouping, good for most information needs, works well for most collections

30

(Callan, et al., 1995; Callan, 2000) © 2002, Junic Callan









	Directory Model: Merging Results	
Returned document rankings	 Global idf + simple score merge e. Effective, but limited to a single organization Domload & rerank documents at client Costly in communication, computation Eurstic reranking of scores at client e. Grective (for Inquery) but ad-hoc Durge-specific, semi-supervised learning E. Gesource selection service assists in merging E. Guernet state-of-the-art 	

















41

© 2002, Jamie Callar

Message-Passing Model: Content-Based Routing (KaZaA)

• Some nodes with spare bandwidth are designated "supernodes" - Supernodes learn contents of other nodes in a region of the network

• Messages are routed first to nearest supernodes

- They route messages to nodes containing desired content

Advantages:

- Very robust, no central site vulnerable to attack
- Less network traffic, better scaling
- More consistent answers to identical requests

• Disadvantages:

- Each node must have a basic understanding of the content in each message it routes

42

© 2002, Jamie Callan

Message Passing Models: **Summary of Techniques**

43

Approaches:

- Content-free routing
- · Content-based routing

Common Assumptions:

- No centralized directories
- Possibly dynamic "regional" directories
- · Direct communication among nodes
- · No assumption about database contents
- · No assumptions about range of queries

Handling Multiple Databases: **Summary of Approaches** Web-search model • Directory models - Minimal user effort

- Handle a wide range of text

- Some handle large-scale problems

Some commercial options (e.g., Verity), but still developing

- Handles a wide range of information

© 2002, Jamie Callar

- Free options (Gnutella, FreeNet) - Some commercial options (FastTrack)

- Some are high precision

information

Message-passing model

- Minimal user effort

Scale is an issue

types

- Minimal user effort Doesn't handle some common situations
- Handles large-scale problems
- Many commercial options (e.g., Google, AltaVista)
- Browsing model
- Considerable manual effort for provider and user
- Handles a wide range of information types Scale is an issue
- Many commercial options (but most of the work is manual)

44

Distributed Retrieval: A Return to Cases

1. Very Large Database (West)

Many different databases, each DB managed by same organization

- 2. Customer Service BureauFew different databases, each DB managed independently
- 3. Competitive Analysis
- · Many different databases, each DB managed independently

4. Distributed wireless networks Nodes in range of other nodes, but not central site, DBs managed independently

45

Which techniques are appropriate for each case?

© 2002, Jamie Callan

For More Information

- Available at http://www.cs.cmu.edu/~callan/
 J. Callan. "Distributed information retrieval." In W. B. Croft, editor, Advances in Information Retrieval (pp. 127-150). Kluwer Academic Publishers. 2000.
- •

- A. Sham, Lyskinowski miormation retrieval, "in W. B. Crolt, editor, Advances in Information Retrieval (pp. 127-150), Winewr Academic Publishers. 2000.
 J. Calam and M. Cannell, "Query-based sampling of text databases," ACM Transactions on Information Systems, 19(2), pp. 97-130. ACM. 2001.
 J. Calam, M. Connell, and A. Du, "Automatic discovery of language models for text databases," In Proceedings of the AGMSGMOD International Conference on Management of Data. (pp. 479-490), Philadelphia: ACM. 1999.
 J. Calam, K. Gonell, and M. B. Croft, W. B. "Searching distributed collections with inference networks," In Proceedings of the EdMSGMOD International ACM SIGIR Conference and Development in Information Retrieval, (pp. 21-28), Seattle, WA, 1995.
 J. Lyck, O. Connell, and Y. L. Callam, "Clicicion selections and results merging with topical pion and the Noveledge Management (ICML) (pp. 22-289), CAUM. 2000.
 J. Arowell, J.C. French, J. Callam, "Clocedings of the 2^{ed} Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 21-23), Seattle, WA, 1992.
 J. Arowell, J.C. French, J. Callam, "Clocedings of the 2^{ed} Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (pp. 232-239), 2000.
 J. Guana, "Lings mayneled data and regression to merge search engine results," in forcedenges of the Yoneyri Fifth Annual International ACM SIGIR Conference on Research and Development, International Conference on Research and Development in Information Conference on Research and Development in Information Retrieval (pp. 232-239), 2000.
 J. Arowell, S. Chalm, "Cling sampled data and regression to merge search engine results," in forcedenges of the Yoneyri Fifth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (pp. 232-239). •

46

© 2002, Jamie Callar