Ranked Search on Data Graphs

Ramakrishna R. Varadarajan

Doctoral Dissertation Defense

FLORIDA INTERNATIONAL UNIVERSITY, School of Computing and Information Sciences, Miami.



Roadmap

- Problem Statement & Motivation
- Data Model
- State of Art Graph Search Methods
- Query-Specific Summarization
- Composed Pages Search
- Explaining & Reformulating Authority-Flow Queries
- Graph Information Discovery (GID)
- Comparing Top-k XML Lists
- Acknowledgements

Roadmap

- Problem Statement & Motivation
- Data Model
- State of Art Graph Search Methods
- Query-Specific Summarization
- Composed Pages Search
- Explaining & Reformulating Authority-Flow Queries
- Graph Information Discovery (GID)
- Comparing Top-*k* XML Lists
- Acknowledgements

Problem Statement & Motivation

- Graph-structured databases becoming a commonplace.
- Need for Efficient & High Quality search & retrieval.
- Common Graph Models:
 - -World Wide Web (unstructured)
 - □ Nodes Pages
 - **Edges** Hyperlinks

-Relational Databases (structured)

- □ Nodes Tuples
- **General Edges Primary/Foreign key relationships**

-XML (semi-structured)

- □ Nodes XML elements
- Edges Intra-document links (IDREFs), Inter-document links (Xlinks)

Problem Statement & Motivation

- *Keyword Search* most effective & dominant information discovery method.
- Success of search engines confirm this.
- <u>Key Advantages</u>:
 - Simplicity (ease of use).
 - **Query interface is flexible.**
 - □ No prior knowledge about structure of underlying data.
 - **Queries can be imprecise**
- Recently applied over Structured (databases) & Semistructured Data (XML).

Goals of the Dissertation

<u>Goal</u> - To facilitate *user-friendly & high-quality ranked search on data graphs* by providing solutions for:

- <u>**Result Discovery</u>** (Composed Pages Search [SIGIR'06,TKDE'08], GID [EDBT'09], Reformulating Authority-Flow queries [ICDE'08]).</u>
- <u>**Result Ranking</u>** (GID [EDBT'09], Composed Pages Search [SIGIR'06, TKDE'08], Reformulating Authority-Flow Queries [ICDE'08], Query-Specific Summarization [CIKM'05,CIKM'06]).</u>
- <u>**Result Presentation</u>** (Explaining Authority-Flow Queries [ICDE'08], Query-specific Summarization [CIKM'05,CIKM'06]).</u>
- **Evaluation of Ranked Results** (Comparing Top-k XML lists).

Roadmap

Problem Statement & Motivation

• Data Model

- State of Art Graph Search Methods
- Query-Specific Summarization
- Composed Pages Search
- Explaining & Reformulating Authority-Flow Queries
- Graph Information Discovery (GID)
- Comparing Top-*k* XML Lists
- Acknowledgements



Slide 7

Web Graph (directed, un-weighted) Web Pages (nodes) & Hyperlinks (edges)

 (v1) Admission Office - Florida International University. (v2) (v3) Florida International University, a member of the State University System of Florida, is a fully accredited comprehensive, multi-campus urban research institution located in Miami, Florida (more). (v4) Open House, Latest Scholarships, The Honors College. (v5) Page 1 (v6) http://www.fiu.edu/orgs/admiss/index.htm (v7) 	 (v1) Financial Aid Office. (v2) International Students. (v3) (v4) (v5) Graduate Assistance - Graduate students pursuing a master or doctoral degree may qualify for assistantships, fellowships and other awards offered through individual academic units
<pre>(v1) Florida International University: A Brief History. (v2) ACHIEVING THE UNIVERSITY'S VISION 1986 TO THE PRESENT. (v3) (v4) (v5) Page 4 (v6) http://www.fiu.edu/docs/brief history4.htm (v7) (v8) (v9) In early 1996 in order to secure resources for faculty</pre>	 (v9) http://www.fiu.edu/orgs/finaid/programs/programs.htm (v10) (v11) Institutional Programs Academic Merit Assistance The University's commitment to academic excellence Additional awards for outstanding high school seniors include the Presidential Scholars, Academic Excellence Scholars, Valedictorian and Salutatorian Scholarships
<pre>(v9) In early 1990, In order to secure resources for faculty endowed chairs, scholarships and facilities, the University launched its first major capital campaign, the \$65 million Campaign for FIU This goal that was achieved in early 2001, more than one year ahead ahead of schedule. (v10) (v11) (v12) Some 30 years after opening its doors, FIU is within reach of attaining its foremost goal - to become one of the nation's top, urban, public research universities</pre>	<pre>(v1) Florida International University: Miami's Public Research University. (v2) (v3) ADMISSIONS: GRADUATE - ARTS & CULTURE - ATHLETICS - BISCAYNE BAY CAMPUS PHONEBOOK -REGISTRATION - STUDENT GOVERNMENT - UNIVERSITY GRADUATE SCHOOL - VALUES STATEMENT - WEB MAIL- SEARCH FIU - ADVANCED SEARCH. (v4) (v5) Page 3 (v5)</pre>
LEGEND Node in Web Graph (Web Pa Hyperlinks	(v5) <u>http://www.fiu.edu/index.html</u>

Sample Document

(vo) Brain chip offers hope for paralyzed

(v1) A team of neuroscientists have successfully implanted a **chip** into the **brain** of a quadriplegic man, allowing him to control a computer.

(v2) ...

(v3) ...

(v4)...

(v5) BrainGate offers the possibility of hitherto unimaginable levels of independence for the severely disabled.

(v6) ...

(v7) ...

(v8) ...

(v9) ...

(v10) Donoghue's initial **research**, published in the science journal Nature in 2002, consisted of attaching an implant to a monkey's **brain** that enabled it to play a simple pinball computer game remotely.

(v11) The four-millimeter square **chip**, which is placed on the surface of the motor cortex area of the **brain**, contains 100 electrodes each thinner than a hair which detect neural electrical activity. The sensor is then connected to a computer via a small wire attached to a pedestal mounted on the skull.

(v12) ...

(v13)...

(v14)...

(v15) "Here we have a **research** participant who is capable of controlling his environment by thought alone -something we have only found in science fiction so far," said Friehs.

(v16) ...

• Page Graph (undirected, weighted)

- Text Fragments (nodes)
- Semantic links (edges)
- ✓ Parsing delimiter NewLine.
- ✓ Text Fragments Paragraphs.
- ✓ 17 text fragments (vo…v16).
- ✓ 17 nodes in Document Graph.



Data Graph (directed, unweighted)

- Tuples (nodes) & primary/foreign key relationships (edges).
- Each node represents an object & has a role.
- Each edge is labeled with its role.
- ✤ Richer semantics metadata.



Paper Authors="H. Gupta, V.

• Schema Graph

Describes the structure of the data graph.

High-level view of Data graph.



Paper Authors="J. Gray, A.

Roadmap

- Problem Statement & Motivation
- Data Model
- State of Art Graph Search Methods
- Query-Specific Summarization
- Composed Pages Search
- Explaining & Reformulating Authority-Flow Queries
- Graph Information Discovery (GID)
- Comparing Top-k XML Lists
- Acknowledgements

State of Art Graph Search Methods

Keyword Proximity Search (as a black box)



State of Art Graph Search Methods

Authority Flow-Based Search (as a black box)



Roadmap

- Problem Statement & Motivation
- Data Model
- State of Art Graph Search Methods

Query-Specific Summarization

- Composed Pages Search
- Explaining & Reformulating Authority-Flow Queries
- Graph Information Discovery (GID)
- Comparing Top-*k* XML Lists
- Acknowledgements

Motivation

Locating relevant information on the Web is hard.

- **Summaries** are **helpful** because:
 - Provide a Quick preview of the document.
 - Allow users to quickly decide relevance.
 - Save user's browsing time.
- Two categories of summaries:
 - *Query-Independent –* Most of prior works.
 - *Query-Specific* Applicable to web search engines.
- Success of *Web search engines* Query specific **snippets** are important.

<u> </u>		Web	<u>Images</u>	<u>Video</u>	<u>News</u>	<u>Maps</u>	more	<u>»</u>
(-000)	e	brain c	hip resear:	ch				Search

Web Results 1 - 10 of about 4,740,000 for brain chip research. (0.30 seconds)

CNN.com - Brain chip research aims for future movement - Mar 1, 2006

Matthew Nagel awoke from a two-week coma in the summer of 2001 to learn he was paralyzed from the neck down.

www.cnn.com/2006/TECH/02/22/brain.gate/index.html - 44k - Cached - Similar pages

CNN.com - Brain chip offers hope for paralvzed - Oct 20, 2004

A team of neuroscientists have successfully implanted a **chip** into the **brain** of a quadriplegic man, allowing him to control a computer.

www.cnn.com/2004/TECH/10/20/explorers.braingate/ - 42k - Cached - Similar pages

BBC NEWS | Health | Brain chips could help paralysed

Brain. The **chip** contains tiny spikes which will extend into the **brain** ... It is hoped the **research** - which until now has been carried out on animals - could ... news.bbc.co.uk/2/hi/health/3632855.stm - 34k - Cached - Similar pages

Motivation

Query-Specific Summaries

Drawbacks of current approaches:

- Ignores semantic relations between keywords in the document.
- Association between query keywords is unclear.
- Follows a naïve approach for query-specific summarization.

Summarization research till date:

- Mostly Query-Independent.
- Not applicable for web search.

- Document $\rightarrow graph$
- We call it *Document Graph*.

Three Steps

Step 1: **Preprocess**

• Build a document graph, *G*. (extract semantic relations between text fragments)

Step 2: Summary Generation (keyword proximity search)

• Given a query *Q* and a document graph *G*,

Summaries \rightarrow *Spanning Trees* that cover all keywords.

Step 3: Rank spanning trees.

Sample Document

Example

Document Graph



something we have only found in science fiction so far," said Friehs.

(v16) ...



Brain chip offers hope for paralyzed.

L Donoghue's initial **research** published in the science journal Nature in 2002 consisted of attaching an implant to a monkey's **brain** that enabled it to play a simple pinball computer game remotely.

User Surveys[CIKM'05,CIKM'06]

	Google Desktop		MSN E	esktop	Our Approach		
Queries	D1	D2	D1	D2	D1	D2	
1	2.33	3.67	2.33	3.67	4.87	3.67	
2	2.00	3.33	2.00	3.00	4.33	3.33	
3	3.00	2.67	0.67	3.00	4.93	4.00	
4	1.67	2.67	1.67	3.00	4.67	4.00	
5	2.00	1.67	3.00	1.00	4.00	3.67	

Queries	Document D1	Document D2
1	Microsoft worm protection	IT Research awards
2	Anti-virus protection	Algorithms development research
3	Recovering worm deleted files	Software projects
4	Worm affected agencies	Large research grants
5	Deleted computer software	Computer network security project

Slide 19

Performance Experiments

News articles from science section of cnn.com



Average times to calculate node weights

Number of keywords	2	3	4	5
Time (msec)	5.31	9.37	11.50	17.33

Average ranks of Top-1 Algorithms

Number of keywords	2	3	4	5
Top-1 Enumeration Algorithm.	1.4	1.8	2.1	2.78
Top-1 Expanding Search Algorithm.	1.1	1.3	1.4	1.8

Roadmap

- Problem Statement & Motivation
- Data Model
- State of Art Graph Search Methods
- Query-Specific Summarization
- Composed Pages Search
- Explaining & Reformulating Authority-Flow Queries
- Graph Information Discovery (GID)
- Comparing Top-*k* XML Lists
- Acknowledgements

http://www.cis.fiu.edu/advisement.php



Motivation

Consider a Keyword Query – "Ph.D. admission requirements & fellowships"

Each Web page only "*partially*" answers the user query.

Basic unit for search, retrieval & ranking individual web page.

Current web search engines don't answer such queries completely.

Motivation

• In WWW - **Information** is typically **distributed** across web pages & are **hyperlinked**.

Current Web Search Engines:

- Basic unit for search & retrieval **individual web page**.
- Return a list of *individual* web pages ranked by relevance.
- This degrades the quality of search results:
 - Especially for Long & Uncorrelated (multi-topic) Queries.
 - Results are not descriptive enough.
 - Does not completely satisfy users information need.
 - Users spend more time searching for relevant information.

We want to extract & stitch together "pieces" of relevant information.

Greatly reduces user browsing time !

CS and IT (Single, Double, Minor, Combined Undergraduate/Graduate)

- Undergraduate Catalog
- CS (Masters, Ph.D.)
- Graduate Catalog

Program Information

Assistanship, Scholarship

http://www.cis.fiu.edu/advisement.php

http://www.cis.fiu.edu/assistantship.php

FIU's University Graduate School (UGS) offers highly competitive Presidential Fellowships (PF) and Presidential Enhancement Fellowships (PEF) to outstanding incoming Ph.D. students. UGS pays each PF \$20,000 plus an in-state tuition waiver for 2 years. SCIS supplements each PF with additional \$5000 annually as a reward. UGS pays each PEF \$10,000 plus in-state tuition waiver. SCIS supplements each PEF with additional \$15,000 annually as a reward.

FIU's UGS offers a limited number of Dissertation Year Fellowships (DYF) to outstanding final year Ph.D. students. Each DYF carries a stipend of \$18,000 plus an in-state tuition waiver; SCIS supplements the DYF to maintain each Ph.D. student.s current stipend level.

"Fellowships"

http://www.cis.fiu.edu/programs/grad/info/7.php

7 Ph.D. Program

7.1 Admission to Ph.D. Program

The requirements for admission to the doctoral program in Computer Science are:

- A baccalaureate or Master's degree in Computer Science, or equivalent degree in a related field as judged by the School's Graduate Committee.
- A minimum of a "B" average on all upper division work and acceptable courses in Calculus and Statistics.
- GRE general test score of 1120 (verbal and guantitative combined), with a minimum guantitative score of 650.
- Three letters of recommendation from persons in a position to judge the applicant's potential for advanced graduate study in computer science.
- Approval of the Graduate Committee.
- Applicants whose native language is not English must score a total of 80 on the iBT TOEFL or 6.3 overall on the IELTS, TOEFL = Test of English as a Foreign Language (www.toefl.org). IELTS = International English Language Testing System (www.ielts.org).

"Ph.D. Admission requirements"

Web Graph (crawled)





v5 Node in Page Graph (text fragment) LEGEND

0.1

0.001

0.078

- Semantic links

0.066

0.001

We extract & stitch together pieces of information.

 v^7

0.001

In contrast to previous works, we go beyond page granularity.

Presentation & Ranking of Composed Pages

First ranking principle - search results involving fewer pages are ranked higher.

Score (R) =
$$\sum_{p \in R} \frac{Score(p)}{PR(p)}$$

Second ranking principle - Search results of same page size, rank according to the involved page spanning trees.

Scores of PSTs are combined using a monotone combining function.



User Surveys [SIGIR'06,TKDE'08]

Keyword Queries	Google Search	Heuristic Expanding Search
Undergraduate Housing safety	2.06	3.41
Graduate financial aid regulations	2.41	3.59
Computer Science Internship opportunities	2.88	3.65
Campus Safety requirement regulations	2.24	3.35
Biomedical Research fellowship eligibility	1.24	3.35
Undergraduate Summer athletics accomplishments	2.25	4.5
Physics alumni achievements	3.25	3.00
Electrical transfer student eligibility	2.66	4.66
Freshman internship opportunities	1.66	4.66
Mechanical Graduate admission policies	1.66	4.66
Average Rating	2.44	3.88

Slide 27

Performance/Quality Experiments

Dataset: Crawled FIU web-pages

Nodes (web pages) – 25,108 & Edges (hyperlinks) - 137,929





(a) Performance with changing k (with m = 2) Execution time for Top-k Search Results. (b) Performance with changing m (with k = 25)



(a) Spearman's rho vs. Top-k (with m = 2)



Quality of Algorithms.

Roadmap

- Problem Statement & Motivation
- Data Model
- State of Art Graph Search Methods
- Query-Specific Summarization
- Composed Pages Search
- Explaining & Reformulating Authority-Flow Queries
- Graph Information Discovery (GID)
- Comparing Top-*k* XML Lists
- Acknowledgements

A Quick Introduction to Authority Flow Ranking:

Consider a Bibliographic Data Graph of papers & citations –

<u>Simple Ranking Strategy</u>: Papers ranked by citation count (vote).

Drawback: Each citation is given equal importance.

<u>A better ranking Strategy</u>: Papers ranked by

- Number of citations with each citation counted according to its importance.
- Importance of each citation, determined by the paper importance. (recursive in nature)
- Evenly divide the "propagated" importance to the cited papers.

 $PR(u) = \sum_{v \in B_{u}} \frac{PR(v)}{L(v)}$ System tunable for Global/Query-specific Importance.



• **ObjectRank** Ranks Objects According to Probability of Reaching Result Starting from Base Set

Motivation – ObjectRank [VLDB'04]

Authority Transfer Data Graph (Keyword Query: [OLAP])



Slide 32

Slide 33

VH2 Database have edges of different types.

Different authority flows through various edges...

The authority transfer rates, which are shown at the bottom, show the maximum ratio of a node's authority transfered over edges of this type.

P->P edge has higher rate than the others because...

Another difference from the way that Web-search engines use PageRank is that we have keyword-specific ObjectRanks

Now assume we have the keyword query OLAP...

In contrast to PageRank on the Web, we can do keyword specific ObjectRanks because (a) smaller size dbs and (b) exploit schema properties to optimize algorithm.

Vagelis, 3/2/2004

Motivation:

Top 20 results for keywords: semistructured data

1.	Object Exchange Across Heterogeneous Information Sources. ICDE 1995. Hector Garcia-Molina, Jennifer Widom et al.
2.	The Lorel Query Language for Semistructured Data. Int. J. on Digital Libraries 1997. Janet L. Wiener, Jennifer Widom et al.
3.	Querying Semi-Structured Data. ICDT 1997. Serge Abiteboul
4.	A Query Language and Optimization Techniques for Unstructured Data. SIGMOD Conference 1996. Gerd G. Hillebrand, Susan B. Davidson et a
5.	DataGuides: Enabling Query Formulation and Optimization in Semistructured Databases. VLDB 1997. Jennifer Widom, Roy Goldman
6.	Mining Association Rules between Sets of Items in Large Databases. SIGMOD Conference 1993. Arun N. Swami, Tomasz Imielinski et al.
7.	Lore: A Database Management System for <u>Semistructured Data</u> . SIGMOD Record 1997. Dallan Quass, Roy Goldman et al.
8.	Querying Semistructured Heterogeneous Information. DOOD 1995. Dallan Quass, Anand Rajaraman et al.
9.	The TSIMMIS Project: Integration of Heterogeneous Information Sources. IPSJ 1994. Hector Garcia-Molina, Kelly Ireland et al.
10.	Semistructured Data. PODS 1997. Peter Buneman
11.	Storing Semistructured Data with STORED. SIGMOD Conference 1999. Dan Suciu, Alin Deutsch et al.
12.	Adding Structure to Unstructured Data. ICDT 1997. Dan Suciu, Susan B. Davidson et al.
13.	Querying Object-Oriented Databases. SIGMOD Conference 1992. Michael Kifer, Yehoshua Sagiv et al.
14.	Querying the World Wide Web. PDIS 1996. Alberto O. Mendelzon, George A. Mihaila et al.
15.	Mediators in the Architecture of Future Information Systems. IEEE Computer 1992. Gio Wiederhold
16.	From Structured Documents to Novel Query Facilities. SIGMOD Conference 1994. Vassilis Christophides, Sophie Cluet et al.
17.	Representing and Querying Changes in <u>Semistructured Data.</u> ICDE 1998. Serge Abiteboul, Jennifer Widom et al.
18.	Relational Queries Computable in Polynomial Time Information and Control 1986. Neil Immerman
19.	Extracting Schema from Semistructured Data. SIGMOD Conference 1998. Serge Abiteboul, Svetlozar Nestorov et al.
20.	Programming Constructs for Unstructured Data, DBPL 1995, Dan Suciu, Susan B. Davidson et al.

Drawbacks:

• Many Top results don't contain query terms in them.

• It is *not obvious* why the results are <u>relevant</u> or <u>important</u> to the query.

• <u>Reason</u> – ranking primarily based on structure & not content.

Motivation:

Limitations of Authority Flow Systems(*ObjectRank*[VLDB04]):

- No way to *explain* to the user why a particular result is relevant/important to the query.
- Authority transfer rates have to be set manually by a domain expert.
- No *query reformulation* methodology to refine results based on user-preferences.

Our Focus

• Typed domain-specific data graphs (Web search - out of scope)

- <u>Problem</u> Given a **target object** *T*, explain user why it received a high rank (or score).
- <u>*Our Solution*</u> Display an **explaining sub-graph** of Authority transfer data graph, for *T*.
- Explaining sub-graph contains:
 - All Edges & corresponding Nodes that transfer authority to T.
 - Edges are annotated with amount of authority flow.
- <u>Steps</u>:
 - Construction Stage (using Bi-directional Breadth-First Search)
 - Flow Adjustment Stage (Adjust original authority flows most challenging)

Explaining & Reformulating Authority-Flow Queries [ICDE'08] Traditional Query Reformulation Methods:

- Well studied in Traditional IR (Salton, Buckley 1990)
- Query Expansion was the dominant strategy (ignores link-structure)
- Term selection, re-weighting, query expansion [SIGIR94, TREC95].

OVERVIEW OF OUR REFORMULATION ALGORITHM:

- 1) System computes Top-k objects with high ObjectRank2 scores.
- 2) User marks relevant "feedback" objects.
- 3) Explaining sub-graphs of feedback objects are computed.
- 4) Reformulate based on (a) **Content** (b) **Structure** of the graph.
- **5)** Practically diameter is limited to a constant (L=3).

User Surveys [ICDE'08]

- <u>Dataset</u>: **DBLP** (Nodes 876,110 & Edges 4,166,626)
- <u>Query Reformulation types tested</u>:
 - Content-based Reformulations (tuning parameters $C_f=0.0 \& C_e=0.2$).
 - Structure-based Reformulations (tuning parameters $C_f=0.5 \& C_e=0.0$).
 - Content & Structure-based Reformulations ($C_f=0.5 \& C_e=0.2$).
- <u>2 stages of experiments</u>:
 - Evaluate Reformulation types (User Surveys using residual collection method).
 - Evaluate how close the trained authority transfer bounds are to the ones set by domain experts in ObjectRank [VLDB04].









Roadmap

- Problem Statement & Motivation
- Data Model
- State of Art Graph Search Methods
- Query-Specific Summarization
- Composed Pages Search
- Explaining & Reformulating Authority-Flow Queries
- Graph Information Discovery (GID)
- Comparing Top-*k* XML Lists
- Acknowledgements

Graph Information Discovery (GID) [EDBT'09]

MOTIVATION

- Consider a biologist's exploration as follows:
 - Starting from *genes* in Entrez Gene she follows *links* to Entrez Protein and then to PubMed.
 - Her objects of interest are papers in **PubMed**.
 - Wants to find **PubMed papers** of highest importance/relevance to keyword "**human**".



- Traverse paths Entrez Gene → Entrez Protein → PubMed.
- Compute **sub-graph**.
- Rank objects in sub-graph for query "*human*" using authority-flow.
- Filter and output Top-k PubMed Publications.



Graph Information Discovery (GID) [EDBT'09]

• <u>Limitations</u> of current graph querying Approaches:

- Support *extremes* of Query complexity :
 - <u>Plain keyword queries</u> limited query capability.
 - <u>Complex queries</u> too hard for users to learn & formulate queries.
 - <u>Fewer solutions in between</u>.
- DOES NOT support:
 - Customized or personalized ranking.
 - Sophisticated *ranking techniques* like authority flow.

• **<u>Objective</u>**: Create a graph querying framework -

- Easy to use & formulate Sophisticated graph queries.
- Rank results by customized or personalized criteria.
- Provides simple & flexible query interface.

• <u>Data Model</u>:

- A rich web of annotated & hyperlinked data entries.
- Includes *schema graph* and a *data graph*.

Graph Information Discovery (GID) [EDBT'09]

GID Query Syntax & Semantics

- A query *q* is a sequence $[r_1 > ... > r_m]$ of FILTERS r_i .
- A <u>Score assignment</u> function for a filter:
 - Assigns a score in [0,1] for each node.
 - Nodes with score o are eliminated (including its edges).



- A **filter** *r*={*R*,*N*,*S*} is the following 3-tuple:
 - R (Filter Selection Condition)
 - A Keyword Boolean expression (or)
 - An Attribute-value pair (or)
 - A Type (or)
 - A Path Expression
 - N (Boolean to specify if R needs to be negated)
 - S (Boolean Soft or a Hard Filter)

GID FILTER TYPES

• HARD FILTER

- *Score assignment function* is **Boolean** (assigns score 0 or 1 to nodes).
- Used to eliminate nodes (and their incident edges).
- Examples:
 - Keyword expression *E*: Score 1 for nodes satisfying *E*.
 - Type *T*: Score 1 for nodes of type *T*; (0 otherwise).

• <u>SOFT FILTER</u>

- Ranking is inherently fuzzy.
- *Score assignment function* could be *complex* (assigns score in [0,1])
 - Authority Flow function,
 - Keyword proximity function (or)
 - IR scoring function.



Experiments [EDBT'09]

Datasets:

•DBLP (Nodes - 876,110 & Edges - 4,166,626)

•**DS7** (Nodes - 699,199 & Edges - 3,533,756)



Quality Experiments of Path-Length-Bound Technique.

Slide 44

Roadmap

- Problem Statement & Motivation
- Data Model
- State of Art Graph Search Methods
- Query-Specific Summarization
- Composed Pages Search
- Explaining & Reformulating Authority-Flow Queries
- Graph Information Discovery (GID)
- Comparing Top-k XML Lists
- Acknowledgements

- The notion of a "*top k list*" is ubiquitous in IR.
- **Objective**: Compare how similar/dissimilar two top-k Lists are based on.
 - Objects present in each list.
 - Ranking of objects.
- **Problem:** Define reasonable and meaningful distance measures between top k lists.
 - Compute a numeric distance value in [0,1]

• Applications:

- Compare different search engines or variations of it.
- Synthesize a good composite ranking function from several simpler ones (rank aggregation).
- Design a Meta Search engine.

- **Current State-of-Art**: Distance measures for permutations and Top-k Lists [Fagin et. al SIAM'03].
 - Spearman's footrule (L1 distance)
 - Spearman's rho (L2 distance)
 - Kentall Tau
- **Objective**: Distance Measures for top-k XML Lists.
- Can we adapt existing methods ?
- Drawbacks of **existing approaches**:
 - Each object in the top-k list is viewed as a WHOLE Object.
 - In XML Consider 2 sub-trees differing by a single node.
 - Matches are BOOLEAN (either match or no-match).
 - In XML Partial matching needed for accurate distance measures.

BACKGROUND (Comparing individual XML trees)

• Tree Similarity Measures:

- Tree-Edit, Tree-Alignment (general tree measures)
- XML-Specific measures
 - Nierman et al. WebDB 2002 (insert-tree, delete-tree)
 - Flesca et al. WebDB 2002 (Fourier transform based Similarity)
 - Buttler 2004 (Path Shingle based Similarity)
 - Helmer VLDB 2007 (Entropy based Similarity)
 - Tag based Similarity

• XML Lists Distance based on Total Mapping:

 $-XLDTM(La,Lb) = a \times MinMSD^{T}(La,Lb) + b \times PD^{T}(La,Lb,fmin^{T})$

(XML Similarity Component) (Position Component)

Datasets:

•DBLP (Elements - 7,137,933 & Average Depth - 1.90 & Max. Depth - 5)

•NASA (Elements - 791,923 & Average Depth - 5.58 & Max. Depth - 8)

XML Similarity Distance Position Distance



XLDTM Experiments on NASA Dataset.

XKeyword (XK)

XML Similarity Distance Position Distance

XKeyword (XK) Slide 49

Research (published/accepted)

- <u>Structure-Based Query-Specific Document Summarization</u>
 - Ramakrishna Varadarajan, Vagelis Hristidis
 - Published in ACM CIKM, 2005 (2-page poster)
- <u>Searching the Web Using Composed Pages</u>
 - Ramakrishna Varadarajan, Vagelis Hristidis, Tao Li
 - Published in ACM SIGIR, 2006 (2-page poster)
- <u>A System for Query-Specific Document Summarization</u>
 - Ramakrishna Varadarajan, Vagelis Hristidis
 - Published in ACM CIKM, 2006 (full paper)
- <u>Beyond Single-Page Web Search Results</u>
 - Ramakrishna Varadarajan, Vagelis Hristidis, Tao Li
 - Published in IEEE TKDE, 2008 (Journal paper)
- <u>Explaining and Reformulating Authority Flow Queries</u>
 - Ramakrishna Varadarajan, Vagelis Hristidis, Louiqa Raschid
 - Published in IEEE ICDE, 2008 (full paper)
- Flexible & Efficient Querying & Ranking on Hyperlinked Data Sources
 - **Ramakrishna Varadarajan**, Vagelis Hristidis, Louiqa Raschid, Maria-Esther Vidal, Luis Ibáñez, Héctor Rodríguez- Drumond
 - Accepted for publication in EDBT, 2009 (full paper)

Research (Current/Ongoing)

<u>Comparing Top-k XML Lists</u>

- Ramakrishna Varadarajan, Fernando Farfan, Vagelis Hristidis
- Under Review in IEEE TKDE, 2009 (Journal paper)
- <u>Using Proximity Search to Estimate Authority Flow</u>
 - Vagelis Hristidis, Yannis Papakonstantinou, Ramakrishna Varadarajan
 - Under Review in IEEE TKDE, 2009 (Concise paper)
- <u>Information Discovery on Electronic Medical Records Using Authority-Flow</u> <u>Techniques</u>
 - Vagelis Hristidis, **Ramakrishna Varadarajan**, Paul Biondich, Redmond Burke, Michael Weiner
 - In preparation for **JAMIA**, 2009 (journal paper)
- <u>Electronic Health Records</u>
 - Fernando Farfan, Ramakrishna Varadarajan, Vagelis Hristidis
 - A book chapter under review in "Information Discovery on EHRs" book
- <u>Searching Electronic Health Records</u>
 - Ramakrishna Varadarajan, Vagelis Hristidis, Fernando Farfan
 - A book chapter under review in "Information Discovery on EHRs" book
- Web Information Extraction Using Visual Patterns

– **Ramakrishna Varadarajan**, Vijil Chenthamarakshan, Prasad Deshpande, Raghuram Krishnapuram

Roadmap

- Problem Statement & Motivation
- Data Model
- State of Art Graph Search Methods
- Query-Specific Summarization
- Composed Pages Search
- Explaining & Reformulating Authority-Flow Queries
- Graph Information Discovery (GID)
- Comparing Top-k XML Lists
- Acknowledgements

- SCHOOL OF COMPUTING AND INFORMATION SCIENCES (SCIS)
 - Consistent Graduate assistantships (TA & RA)
 - Awards recognizing student research.
 - Conference Travel Support.





Professor Masoud Milani



Professor Yi Deng

- FLORIDA INTERNATIONAL UNIVERSITY
 - Dissertation Year Fellowship (DYF)
 - FIU GSA for travel awards



Dissertation Committee (for great advise & supervision):



Professor Vagelis Hristidis (<u>Academic Advisor</u>)



Professor Shu-Ching Chen



Professor Tao Li



Professor Raju Rangaswami



Professor Kaushik Dutta

Research Collaborators (for their support & time):



Professor Louiqa Raschid

University of Maryland at College Park



Professor Gautam Das University of Texas at Arlington



Professor Maria-Esther Vidal

Universidad Simón Bolívar



Dr. Raghuram Krishnapuram

IBM India Research Lab

Members of Database & Systems Research Lab (DSRL):



Thanks !

Questions/Comments/

Suggestions ???

Related Work

Document Summarization

- ✓ Mostly Query-Independent
- ✓ Summarizing Web Pages
 - □ OCELOT Berger et.al [SIGIR2000] synthesizes summaries (*non-extractive*).
 - □ INCOMMENSENSE Paris et.al [CIKM2000] uses anchor text (ignores content).
- ✓ Splitting Web pages in to blocks
 - □ Song et.al [WWW2004] Block importance models (learning algorithms)
 - Cai et.al [SIGIR2004] Block level link analysis

✓ Document modeled as Graphs

- Lexrank [JAIR2004] : Sentence Centrality using link analysis.
- □ TextRank [EMNLP2004]: "representative" sentences using link analysis.

Keyword Search on Data Graphs

- BANKS [ICDE2002]: group-steiner tree problem
- DISCOVER [VLDB2002], DBXplorer [ICDE2002], IRStyle [VLDB2003].
- **XRANK** [SIGMOD2003], Xkeyword [ICDE2003] : search in XML documents.



Information Unit [WWW Conference 2001]

Tree of hyperlinked pages containing *ALL keywords* - "logical **Information Unit**" (page-level)

Traditional IR Ranking

- Term weighting State of art IR is based on tf *idf principle.
 - Okapi Formula (Modern IR overview Singhal [IEEE data bulletin 2001]).
 - Pivoted Normalized Weighting.

Link-Based Semantics

- 1. Web PageRank [WWW98], HITS [ACM Journal 99], Topic-Sensitive PageRank [WWW02]
- **2.** Database ObjectRank for the database [VLDB02].
- **3.** XML XRANK [SIGMODo3].

Authority Transfer Schema Graph

- Edges reflect the authority transfer rates.
- ✤Bi-directional authority transfer.
- Potentially different rates for each direction.

Authority Transfer Data Graph (directed, weighted)

Conference

Data graph edges labeled with authority transfer rates.



0.7 cites

Paper

0 cited

0.1

Year

0.2

0.2

Author

Target Object – "Modeling Multidimensional databases" paper for query "OLAP".

Explaining Sub-graph Creation

- Perform a BFS search in reverse direction from the target object. 1.
- Perform a BFS search in forward direction from base set objects 2. (authority sources).
- Sub-graph will contain all nodes/edges traversed in the forward 3. direction.

