

WETSUIT: An Efficient Mashup Tool for Searching and Fusing Web Entities



Stefan Endrullis, Andreas Thor, Erhard Rahm

Database Group, University of Leipzig
http://dbs.uni-leipzig.de

Motivation

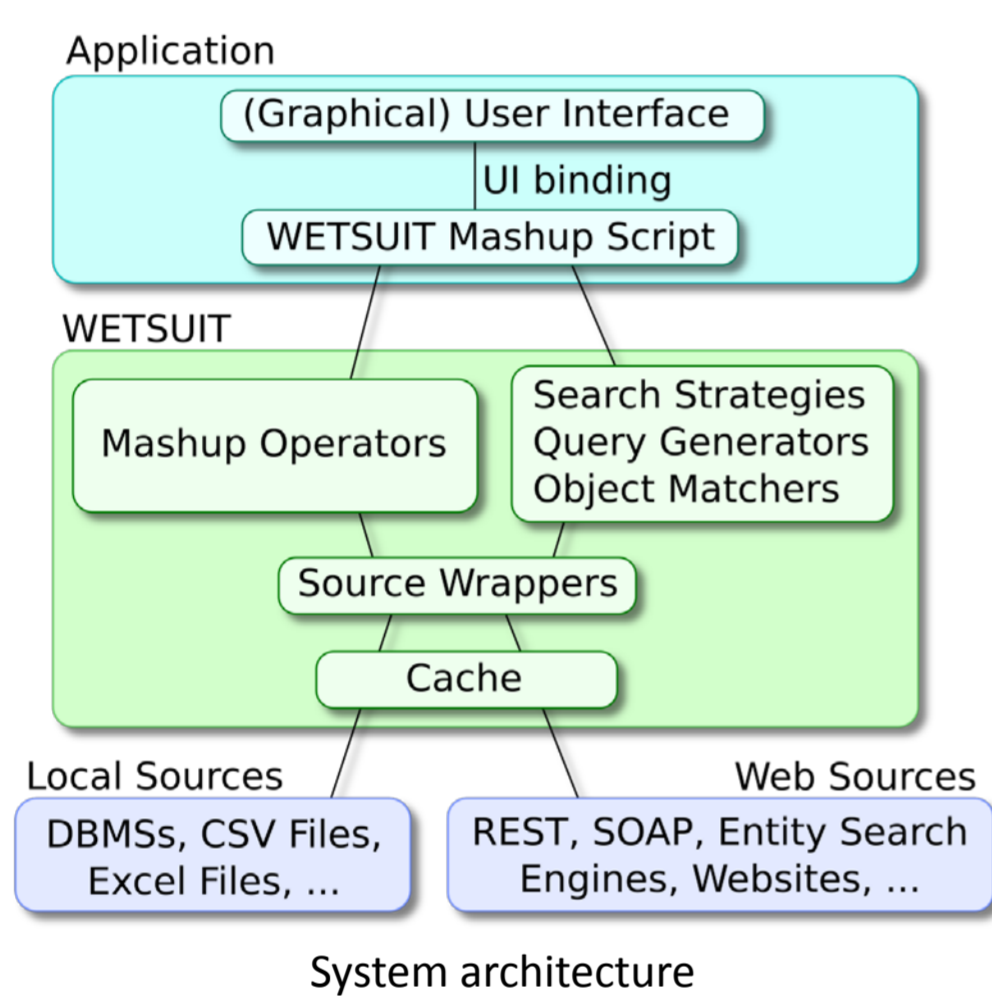
What are the top cited papers of the VLDB 2002?

- Tasks involved:
 - Find DBLP publications for VLDB 2002
 - Find corresponding Scholar publications
 - Determine Matches
 - Aggregate Scholar citations
 - Print the top 10
- Requires efficient entity search: e.g. start with one venue query, refine with a keyword query for each DBLP publication
- Requires entity resolution: e.g. compare pub. titles and authors
- Present intermediate results to reduce waiting time for user

WETSUIT

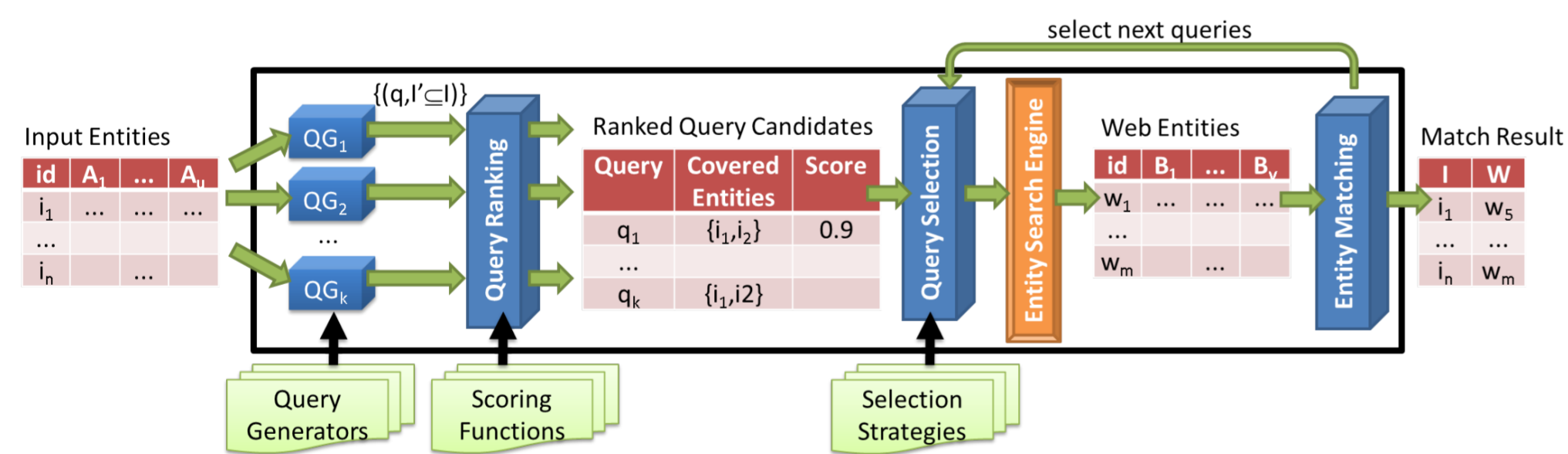
What is WETSUIT?

- WETSUIT: Web Entity Search and fUson Tool
- Efficient and powerful mashup framework



Features

- Efficient execution of mashup workflows
 - Efficient and effective entity retrieval from entity search engines using **advanced search strategies** which exploit multiple query generators
 - may increase recall while reducing the number of search queries



- Pipeline parallelism:** all workflow operators work in parallel as long as they have input entities to process
- Data parallelism:** each operator may process multiple entities simultaneously
- Entity resolution** to deal with dirty data
- Fast presentation of first results**
 - By streaming partial results immediately to the next operators
 - **Result refinement in background**
 - Classical blocking operators (e.g. filterTop, groupBy, aggregate) stream preliminary results, i.e., they may update or revoke them later
- Mashups definition language **embedded in Scala**
 - Concise and well-readable syntax
 - High extensibility** due to Scala and WETSUIT's object oriented design
 - Simple integration of existing Java/Scala libraries in WETSUIT mashups (e.g. for entity extraction from HTML documents) due to Java Bytecode compatibility
 - Well supported by development tools (e.g. Eclipse, IntelliJ IDEA, maven, ant)
- Automatic GUI generation** based on workflow definition
- Open Source** (GNU Affero General Public License (AGPL) Version 3)

Mashup Operators

- Set-based; support streaming of partial and preliminary results

Selected mashup operators of WETSUIT

Operator	Definition
inputOne (<i>label</i>)	asks the user for one/multiple input value(s); <i>label</i> denotes the caption of the GUI component
inputMany (<i>label</i>)	
selectOne (<i>label</i>)	displays intermediate results and asks the user to select one/multiple of them
selectMany (<i>label</i>)	
outputOne (<i>label</i>)	presents results
outputMany (<i>label</i>)	
map (<i>f</i>)	maps each input entity to a new entity using the mapping function <i>f</i>
flatMap (<i>f</i>)	maps each input entity to a set of new entities using the mapping function <i>f</i>
filter (<i>cond</i>)	filters the input entities based on the filter condition <i>cond</i>
groupBy (<i>groupAttr</i>)	groups the input entities by <i>groupAttr</i> ; groupBy has to be followed by filterTop or aggregateValue
filterTop (<i>n</i>) by (<i>orderAttr</i>)	for each group/whole input set: filters the top <i>n</i> elements on <i>orderAttr</i>
aggregateValue (<i>value</i>) via (<i>agg</i> , [<i>rev</i>])	for each group/whole input set: aggregates <i>value</i> using aggregation function <i>agg</i> and its reverse operation <i>rev</i> (<i>rev</i> is needed if input entities can be revoked)
findAt (<i>ese</i>)	searches for the input entities at entity search engine <i>ese</i> using an implicitly defined search strategy
matchWith <i>set₂</i> using (<i>matcher</i>)	matches the input entities against entities of <i>set₂</i> by using the match strategy <i>matcher</i>
union, intersect, diff	classical set operators
join <i>set₂</i> on (<i>theta</i>)	θ -join of two input sets using the binary function <i>theta</i> as join condition
buffer (<i>minDelay</i> , <i>maxDelay</i>)	buffers new entities and entity revokes for at least <i>minDelay</i> and at most <i>maxDelay</i> milliseconds

Mashup Example: Online Citation Service

Screenshot

Journal/conference:

Select a journal/conference

FULL TITLE	TITLE	TYPE	YEAR
VLDB 2002, Proceedings of 28th International Conference on Very Large Data Bases	VLDB	Conference	2002
Efficiency and Effectiveness of XML Tools and Techniques and Data Integration	EEXTT	Conference	2003

DBLP pubs with citations

CITATIONS	TITLE	AUTHORS	YEAR
966	Approximate Frequency Counts over Data Streams	Gurmeet Singh Manku, Rajeev Motwani	2002
949	Streaming Queries over Streaming Data	Sirish Chandrasekaran, Michael J. Franklin	2002
921	COMA - A System for Flexible Combination of Queries	Hong Hai Do, Erhard Rahm	2002
853	Monitoring Streams - A New Class of Data Management Systems	Donald Carney, Ugur Çetintemel, Mihai Anitei	2002
694	XMark: A Benchmark for XML Data Management Systems	Albrecht Schmidt, Florian Waas, Martin Günzler	2002
641	Hippocratic Databases	Rakesh Agrawal, Jerry Kiernan, Ramakrishnan Srikant	2002

Workflow definition

```
inputOne ("Journal/conference: ", "VLDB 2002")
// search for corresponding DBLP journals or conferences
flatMap (name => Dblp.Jonf.where("full_title like ?", "%"+name+"%"))
selectOne ("Select one journal/conference")
flatMap (_.publications)
// search for the selected DBLP publications at Google Scholar
findAt (Scholar) // returns a set of correspondences (DBLP pub., GS pub., similarity)
// for each Scholar publication find the best matching DBLP publication
groupBy (_.range) filterTop 1 by (_.sim)
// summarize the Scholar citation counts for each DBLP publication
groupBy (_.domain) aggregateValue (_.range.citations) via (_.+, _-_)
outputMany ("DBLP pubs with citations")
```

Related Work

- S. Endrullis, A. Thor, and E. Rahm. Entity Search Strategies for Mashup Applications. In ICDE, 2012.
- A. Thor and E. Rahm. CloudFuice: A Flexible Cloud-Based Data Integration System. In ICWE, 2011.
- S. Endrullis, A. Thor, and E. Rahm. Evaluation of Query Generators for Entity Search Engines. In Int. Workshop on USETIM, 2009.