



# EVALUATION OF QUERY GENERATORS FOR ENTITY SEARCH ENGINES

Stefan Endrullis, Andreas Thor, Erhard Rahm

University of Leipzig

<http://dbs.uni-leipzig.de>

Aug. 24, 2009

# ENTITY SEARCH ENGINE (ESE)

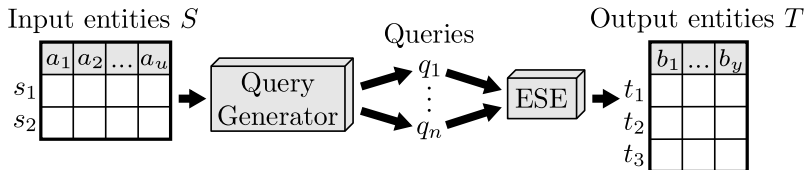
Examples: Google Scholar, Amazon Advanced Search for DVDs

- domain specific search engine
- web interface for querying certain type of entities

<b>Author</b>	Return articles written by	<input type="text"/>
		e.g., "PJ Hayes" or McCarthy
<b>Publication</b>	Return articles published in	<input type="text"/>
		e.g., J Biol Chem or Nature

- set of predicates  $p_1, \dots, p_m$
- corresponding search values  $v_1, \dots, v_m$
- query interpretation:  $q \approx \bigwedge_{1 \leq i \leq m, v_i \neq \epsilon} p_i(v_i)$
- result: top-k entities

Common data integration task: find specific set of entities at ESE



Relevant for

- offline data integration (e.g. data warehousing)
- online data integration (e.g. mashups)

Examples:

- list of products → find prices
- list of research papers of venue → enrich with citation counts

- Challenges
- Concept of Query Generators
- Evaluation of Query Generators
- Initial Evaluation for Bibliographic Domain
- Summary and Outlook

## Problem:

- find given set of entities at an entity search engine

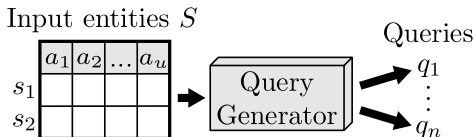
## Challenges:

- how to find effective and efficient search queries?
- myriad of different queries (predicate combinations)
- data quality problems (misspelling, missing values, synonyms, ...)
- limited number of requests per time interval
- limited number of entities per result page

# CONCEPT OF QUERY GENERATORS

Query Generator:

- function that maps a set of entities (multi-valued relation) to a set of search queries



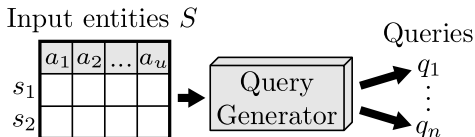
4 general properties based on generic data structures:

- attribute-predicate mapping
- search value generation
- partitioning: naive, frequent value
- aggregation: query combination

# CONCEPT OF QUERY GENERATORS

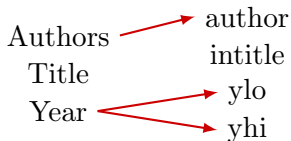
Query Generator:

- function that maps a set of entities (multi-valued relation) to a set of search queries



4 general properties based on generic data structures:

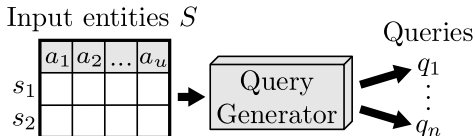
- **attribute-predicate mapping**
- search value generation
- partitioning: naive, frequent value
- aggregation: query combination



# CONCEPT OF QUERY GENERATORS

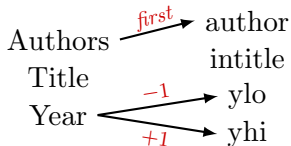
Query Generator:

- function that maps a set of entities (multi-valued relation) to a set of search queries



4 general properties based on generic data structures:

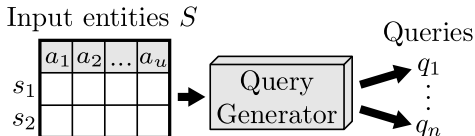
- attribute-predicate mapping
- **search value generation**
- partitioning: naive, frequent value
- aggregation: query combination



# CONCEPT OF QUERY GENERATORS

Query Generator:

- function that maps a set of entities (multi-valued relation) to a set of search queries



4 general properties based on generic data structures:

- attribute-predicate mapping
- search value generation
- **partitioning**: naive, frequent value
- aggregation: query combination

Input entities  $S$

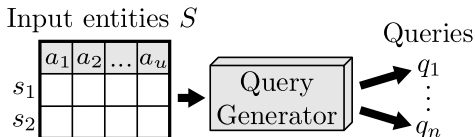
	$a_1$	$a_2$	...	$a_u$
$s_1$				
$s_2$				
$s_3$				

partition  $S$   
into  $S_1, \dots, S_n$

# CONCEPT OF QUERY GENERATORS

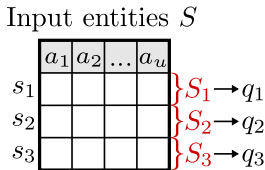
Query Generator:

- function that maps a set of entities (multi-valued relation) to a set of search queries



4 general properties based on generic data structures:

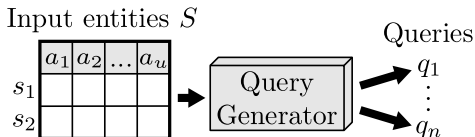
- attribute-predicate mapping
- search value generation
- partitioning: **naive**, frequent value
- aggregation: query combination



# CONCEPT OF QUERY GENERATORS

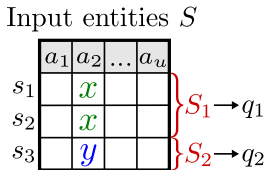
Query Generator:

- function that maps a set of entities (multi-valued relation) to a set of search queries



4 general properties based on generic data structures:

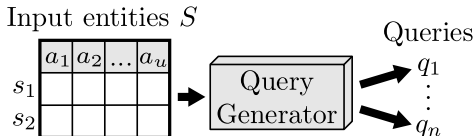
- attribute-predicate mapping
- search value generation
- partitioning: naive, **frequent value**
- aggregation: query combination



# CONCEPT OF QUERY GENERATORS

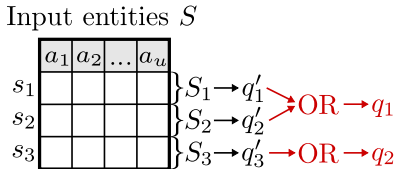
Query Generator:

- function that maps a set of entities (multi-valued relation) to a set of search queries



4 general properties based on generic data structures:

- attribute-predicate mapping
- search value generation
- partitioning: naive, frequent value
- **aggregation: query combination**



# NAIVE QUERY GENERATOR

Properties:

- one query per input entity
- OR aggregation to reduce number of queries
  - not applicable to any ESE

Example:

Authors	Title
{Smith, Jones}	The question to 42
{Williams, Smith}	Don't Panic

Authors  $\xrightarrow{\text{keywords}}$  author  
Title  $\xrightarrow{\text{keywords}}$  intitle

Queries:

- $q_1 = \text{intitle}(\text{question } 42)$
- $q_2 = \text{intitle}(\text{panic})$

# NAIVE QUERY GENERATOR

Properties:

- one query per input entity
- OR aggregation to reduce number of queries
  - not applicable to any ESE

Example:

Authors	Title
{Smith, Jones}	The <b>question</b> to 42
{Williams, Smith}	Don't <b>Panic</b>

Authors  $\xrightarrow{\text{keywords}}$  author  
Title  $\xrightarrow{\text{keywords}}$  intitle

Queries:

- $q_1 = \text{intitle}(\text{question } 42)$
- $q_2 = \text{intitle}(\text{panic})$

# NAIVE QUERY GENERATOR

Properties:

- one query per input entity
- OR aggregation to reduce number of queries
  - not applicable to any ESE

Example:

Authors	Title
{Smith, Jones}	The <b>question</b> to 42
{Williams, Smith}	Don't <b>Panic</b>

Authors  $\xrightarrow{\text{keywords}}$  author  
Title  $\xrightarrow{\text{keywords}}$  intitle

Queries:

- $q_1 = \text{intitle}(\text{question } 42)$
- $q_2 = \text{intitle}(\text{panic})$

# FREQUENT VALUE QUERY GENERATOR

## Properties:

- input analysis  $\rightarrow$  identify frequent item sets (FIS) on attribute values
- variation of Apriori algorithm
- partitioning based on FISs with largest coverage  $\rightarrow$  reduce #queries
- applicable to any ESE

## Example:

Authors	Title
{Smith, Jones}	The question to 42
{Williams, Smith}	Don't Panic

Authors  $\xrightarrow{id}$  author  
Title intitle

## Queries:

- $q_1 = \text{author}(\text{Smith})$

# FREQUENT VALUE QUERY GENERATOR

## Properties:

- input analysis  $\rightarrow$  identify frequent item sets (FIS) on attribute values
- variation of Apriori algorithm
- partitioning based on FISs with largest coverage  $\rightarrow$  reduce #queries
- applicable to any ESE

## Example:

Authors	Title
{ <b>Smith</b> , Jones}	The question to 42
{Williams, <b>Smith</b> }	Don't Panic

Authors  $\xrightarrow{id}$  author  
Title intitle

## Queries:

- $q_1 = \text{author}(\text{Smith})$

# FREQUENT VALUE QUERY GENERATOR

## Properties:

- input analysis  $\rightarrow$  identify frequent item sets (FIS) on attribute values
- variation of Apriori algorithm
- partitioning based on FISs with largest coverage  $\rightarrow$  reduce #queries
- applicable to any ESE

## Example:

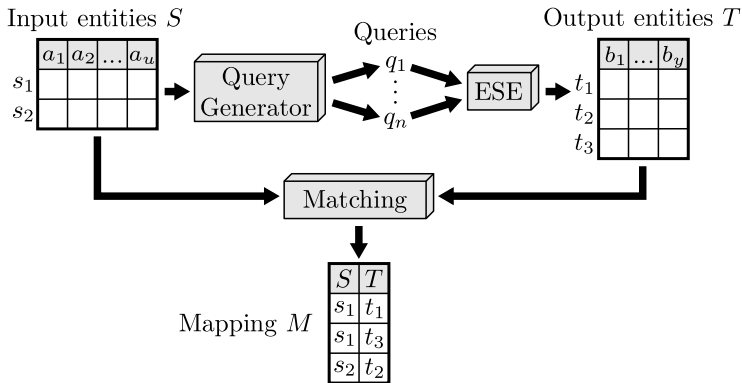
Authors	Title
{ <b>Smith</b> , Jones}	The question to 42
{Williams, <b>Smith</b> }	Don't Panic

Authors  $\xrightarrow{id}$  author  
Title intitle

## Queries:

- $q_1 = \text{author}(\text{Smith})$

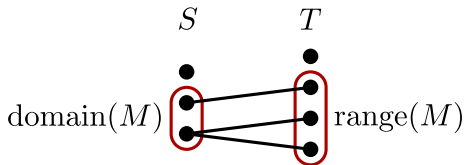
# EVALUATION OF QUERY GENERATORS



- automatic evaluation of several query generators for different datasets
  - matching between  $S$  and  $T \rightarrow$  identify same real world objects
  - query statistics ( $S, T, M, \#$ requests) stored in data warehouse to perform evaluation afterwards

# EVALUATION MEASURES

Mapping example:



Evaluation measures:

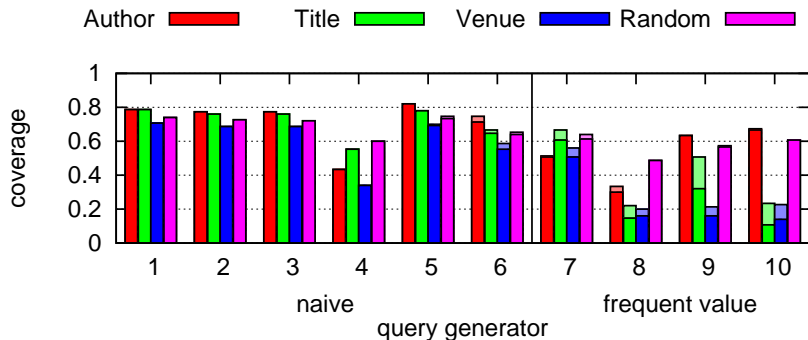
$$\textit{precision} := \frac{|\text{range}(M)|}{|T|} = \frac{3}{4}$$

$$\textit{coverage} := \frac{|\text{domain}(M)|}{|S|} = \frac{2}{3}$$

$$\textit{efficiency} := \frac{|\text{domain}(M)|}{\textit{requests}} = 2$$

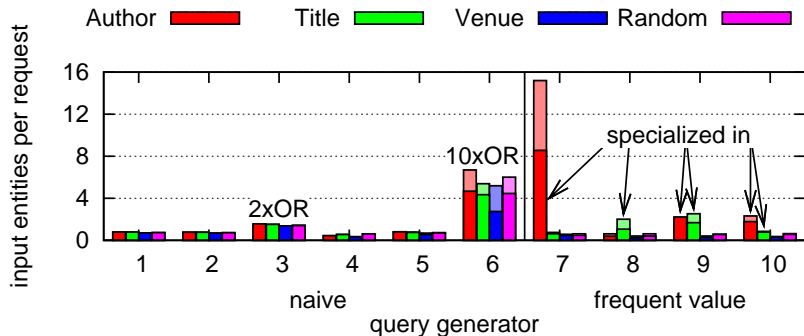
- ESE: Google Scholar, at most 300 results ( $3 \times 100$ )
- 10 query generators
- 60 ( $= 4 \cdot 3 \cdot 5$ ) datasets based on DBLP:
  - 4 categories (Author, Title, Venue, Random)
  - 3 dataset sizes (5, 30, 100)
  - 5 different datasets for each combination

# INITIAL EVALUATION - COVERAGE



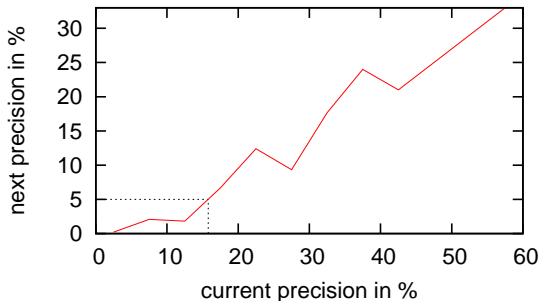
- naive QGs:
  - better coverage than frequent value QGs (except #4 which is too precise)
- frequent value QGs:
  - partially high coverage variations between different categories

# INITIAL EVALUATION - EFFICIENCY



- aggregation raises efficiency
- frequent value QGs:
  - efficiency strongly depends on dataset characteristics (category)

# ESTIMATED PRECISION FOR THE NEXT REQUEST



- raise efficiency by following the next link depending on its estimated precision (estimation based on current precision)

Example:

- aim: next link precision  $\geq 5\%$
- conclusion: follow next link only iff current precision  $\geq 15.8\%$

## Summary:

- approach for flexible query generators
- approach for automatic evaluation of query generators
- initial evaluation for bibliographic domain

## Future Work:

- adaptive search strategies to automatically select most efficient query generator (combination)
- query generators for additional domains and ESEs