

2018-06-20

# Core techniques of QA Systems over KBs a Survey

Core techniques of  
QA Systems over KBs  
a Survey  
Guillermo Echegoyen Blanco

# Core techniques of QA Systems over KBs a Survey

## └ Overview

- 1 Intro
- 2 Tasks
- 3 Question Analysis
- 4 Phrase Mapping
- 5 Disambiguation
- 6 Query Construction
- 7 Conclusions
- 8 References

# Core techniques of QA Systems over KBs a Survey

└ Intro

└ Intro

Intro

- A Question Answering System should be able to:  
*Understand a Natural Language Question so as to be able to answer based on some pre-known data.*
- Typically involves accepting a question and generating a SparQL query capable of extracting the information which answers the user question.
- QALD benchmark
- WebQuestions benchmark
- SimpleQuestions benchmark

Perspective taken by most of the QA Systems evaluated on QALD  
Systems are tested like black boxes, so it is not clear enough with techniques work well on each part.

2018-06-20

# Core techniques of QA Systems over KBs a

## Survey

└ Tasks

└ Tasks

Tasks

- Question Analysis
- Phrase Mapping
- Disambiguation
- Query Construction

# Core techniques of QA Systems over KBs a Survey

## └ Question Analysis

### └ Question Analysis #1

- (is it a Which, What. . . question).
- (is it in English, French. . . ).

Analyze syntactic features to extract meaningful information:

- Type of question
- Multilinguality
- Correspondance to KB entities/classes.
- Tokens in the sentence and it's relations.
- Useless words in the sentence.

# Core techniques of QA Systems over KBs a Survey

## └ Question Analysis

### └ Question Analysis #2

Techniques based on:

- Recognizing Named Entities
- Segmenting with POS\* Tags
- Identifying dependencies using parsers

POS Tag: Part-Of-Speech Tag

- Recognizing Named Entities consists in finding the entities corresponding to parts of the phrase (eg: Europe dbr:European\_Union): Which token correspond to which resource in the KB
- Segmenting is like tokenization of different parts of the string, where the tag is usually universal
- Dependencies refer to parts of the phrase which depend upon others, direct complement, adjective, subjective noun. . .

# Core techniques of QA Systems over KBs a Survey

## └ Question Analysis

### └ Question Analysis #3 - Recognizing named

Identify Named Entities and map to resource in KB

- **NER Tools:** Tools from NLP. **Stanford NER Tool.** Domain specific, **low precision 51%** (hu2014a)
- **N-Gram:** Map n-grams to KB entities. Adv: Each NE can be recognized in the KB, disadvantage: Disambiguation explodes (**too much candidates**). (SINA: shekarpour2015a, CASIA: hu2014a)
- **Entity Linking Tools:** **DBpedia Spotlight** (daiher2013a), **DBpedia Lookup and AIDA** (yosef2011a). Recognize NE and find the underlying KB resource, disambiguating on the way. Adv: All-in-one. Disadv: Limited service, KB dependant.

Identify tokens in the sentence that refer to a resource in the KB, discarding useless words.

- When grouping n-grams, if an entity is found, the n-gram is considered, else more n-grams are tried.

Propose n-grams with attention mechanism?

# Core techniques of QA Systems over KBs a Survey

## └ Question Analysis

## └ Question Analysis #4 - Segmenting using POS

Identify which phrase correspond to instances, properties, classes... and which is irrelevant.

- Handmade rules: Regular expressions depending on question type, structure... (PowerAqua Lopez2012a, Treo Freitas2014a, DEANNA yalpa2013a). Disadv: regex built by hand.

WRB VBD DT NNP NNP VBN .  
When was the European Union founded ?

Figure: POS tagging from the Stanford POS Tagger

The general strategy with POS tags is to identify some reliable POS tags to recognize entities relations and classes. Regex over those POS tags



# Core techniques of QA Systems over KBs a Survey

## └ Question Analysis

## └ Question Analysis #5 - Segmenting using POS

- Learning rules: **Machine Learning** approach, train over corpus (Xue [xu2014a](#), UTQA [pouran2016a](#), very good results). Disadv: **training corpus needed**.

are VB CB are am E.B S.I B.B .  
By which countries was the European Union founded ?

[Figure](#) Question annotated with **CoNLL IOB** format

Combine POS tags, NER tags, no handmade rules

2018-06-20

# Core techniques of QA Systems over KBs a Survey

└ Question Analysis

└ Question Analysis #6 - Parsers

Question Analysis #6 - Parsers

Dependency grammars: **Stanford dependency parser**, word dependencies. Adv: can extract relations along with it's arguments (gRinner zou2014a, PATTY nakashole2012a)

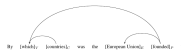


# Core techniques of QA Systems over KBs a Survey

## └ Question Analysis

### └ Question Analysis #7 - Parsers

Phrase Dependencies and DAGs: Dependencies between phrases. SHIFT-REDUCED parser. Disadv: **parser trained on dataset (Xie et al 2014a)**.



- DAG based parser operates on a phrase level, dependency grammars on a word level.
- DAG uses POS tags as features

# Core techniques of QA Systems over KBs a Survey

## Question Analysis

## Question Analysis #8 - Summary

[illegible]

2018-06-20

# Core techniques of QA Systems over KBs a Survey

└ Question Analysis

└ Question Analysis #9 - Summary

## Which techniques to choose?

- Xier (**trained DAG**) reports best results on QALD 4.1 & 5
- gAnswer (**Dependency grammars**) reports fastest results on QALD 3 & 4
- UTQA (**Learned POS tags**) reports best results on QALD 6

Machine Learning approach: Can be fast enough and there is plenty of data available.

# Core techniques of QA Systems over KBs a Survey

## └ Phrase Mapping

### └ Phrase Mapping #1

Find the resources in the KB with the highest probability that maps to the phrase.

Problems:

- String similarity
- Semantic similarity
- Language

- In this phase: consider the phrase mapping problem when  $s$  and  $\text{label}(r)$  have only a semantic relation
- String similarity: very similar words, different meaning (which, witch)
- Semantic similarity: words with related semantic meaning but different writing (king, queen)
- PATTY is used in Xser

# Core techniques of QA Systems over KBs a Survey

## └ Phrase Mapping

### └ Phrase Mapping #2

- Database with lexicalization: WordNet, Wiktionary, PATTY. Expand the phrase with synonyms and use that for search. Adv: High number of candidates, disadv: Big search space, not very useful for domain specific mappings.
- Mappings using large texts: word2vec semantics reflected in the associated vector. Adv: aids in lexical gap, string similarity and semantic similarity, disadv: needs training on large texts, noisy, performance.

- PATTY is a database with relational lexicalization, uses pattern synsets (is album, [[num]] album by)
- A possible advantage of using PATTY is that response text could be "easily" constructed.

# Core techniques of QA Systems over KBs a Survey

## └ Phrase Mapping

### └ Phrase Mapping #3

- KB Labels: Search in the labels provided by the KB's entity (a1)
- Redirects: Follow the owl:sameAs links (gAnswer zou2014a)
- Extracted knowledge: From the previous phase (gAnswer zou2014a). Relations and arguments.
- Wikipedia specific: **DBpedia Lookup**, **Wikimedia Miner Tool** (gAnswer zou2014a, Xue xu2014a, zhu-a)

Labels can be a powerful resource, smart indexing techniques can be used here (gAnswer)



- Phrase Mapping

## Phrase Mapping #4 - Summary

- Distributional semantics word2vec
- Indexed labels: Lucene index

[illegible]

# Core techniques of QA Systems over KBs a Survey

└ Phrase Mapping

└ Phrase Mapping #5 - Summary

## Which techniques to choose?

The best results here depend on the previous step and the computing resources available.

Options (can be combined together):

- KB's Labels
- Redirects
- word2vec
- PATTY

Mixed approach with all the possible methods?? Maximize:

- performance
- KB independance

# Core techniques of QA Systems over KBs a Survey

## └ Dissambiguation

### └ Disambiguation #1

QA systems generate lots of possible interpretations due to language ambiguities and search process.

- Find univocally the resource that maps to the requested question.

Basic approach (local disambiguation):

- String or semantic similarity to resource label.
- Consistency check between the properties and their arguments.

- String or semantic similarity to resource: (include)
- Consistency check between the properties and their arguments: (exclude).
- Local disambiguation excluded, all systems do it. Example "Who is the director of The Lord Of the Rings?", with no information associated with the director resource, it is not possible.

# Core techniques of QA Systems over KBs a Survey

## └ Dissambiguation

### └ Disambiguation #2 - Graph Search

Disambiguation carried out in the KB search step

- Subgraph matching against the KB ([gAnswer2014a](#) does it on phrase mapping). Represent the question as a dependency graph and find an isomorphic subgraph in KB. Adv: very fast. Disadv: disambiguation carries over. (**high precision, low recall**)
- [SemSek](#) [aggarwal2012a](#) and [Treo](#) [freitas2014a](#) do it only with recognized instances (during question analysis phase). (**low precision, high recall**)

Assume that all relational phrases can be deduced from the question.

- gAnswer scores each possible match proportionally to the distance between labels and resources, searches both in edges and nodes.
- PowerAqua explores only the most probable mappings, doing a balance between recall and precision based on the question analysis.
- SemSek and Treo do not explore the relations, but the attached properties of the different interpretations on the KB found matches.

2018-06-20

# Core techniques of QA Systems over KBs a Survey

└ Dissambiguation

└ Dissambiguation #3 - Graph Search

Dissambiguation #3 - Graph Search

**gAnswer** searches in the edges and vertices, **SemSel**, **Tree** search on instances and properties attached.

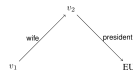


Figure: Subgraph generated for the question "Who is the wife of the president of the EU?"

# Core techniques of QA Systems over KBs a Survey

## └ Dissambiguation

## └ Dissambiguation #4 - Hidden Markov Model

- Assumption: This means that the appearance of a resource at time  $t$  depends only on the appearance of a resource at  $t-1$
- Random variables

"By which countries was the EU founded?"  
Two stochastic processes:

- **Hidden (dissambiguation)**  $X_{t \in \mathbb{N}}$ :  $\{db\text{:Country}, db\text{:Euro}, db\text{:EuropeanUnion}, db\text{:founded}, db\text{:establishedEvent}\}$ .
- **Observed (question tokens)**,  $Y_{t \in \mathbb{N}}$ :  $\{\text{"countries"}, \text{"EU"}, \text{"founded"}\}$ .



# Core techniques of QA Systems over KBs a Survey

└ Dissambiguation

└ Dissambiguation #5 - HMM

The problem is reduced to find the most probable set of states. Extra parameters:

- Initial probability  $P(X_0 = x)$  for  $x \in X$
- Transition probability  $P(X_t = x_1 | X_{t-1} = x_2)$  for  $x_1, x_2 \in X$
- Emission probability  $P(Y_t = y | X_t = x)$  for  $x \in X, y \in Y$

It is not necessary to know the dependency between different resources, just the available resources.

# Core techniques of QA Systems over KBs a Survey

└ Dissambiguation

└ Disambiguation #6 - HMM

**SINA (shekarpour2015a): slow**

- Emission: string similarity between label and segment.
- Initial & Transition: estimated based on the distance of the resource in the KB and popularity.

**RTV (giannone2013a): inaccurate**

- Emission: word embeddings
- Initial & Transition: uniform across all resources

The slow part is the distance to the resource in KB (first approach), biiiig search space



# Core techniques of QA Systems over KBs a Survey

└ Dissambiguation

└ Dissambiguation #7 - ILP & MLN

## ILP Optimization problem

- **DEANNA (yahya2013a)** Dependencies between the segments have to be computed in the question analysis phase. **slow, low precision & recall**

## Markov Logic Network

- **CASIA (hu2014a)** Hard constraints like ILP, soft constraints flexibility **training needed low precision & recall**

Ambiguity during phrase mapping and segmentation

ILP: Dependencies between the segments have to be computed in the question analysis phase.

Boolean variables to indicate:

- if a segment of a question is chosen or not
- if a resource corresponding to a segment is chosen
- if a segment corresponds to a property or an instance

Optimization function terms:

- Increase if the label of a resource is similar to the corresponding segment
- Increase if the two selected resources often occur in the same context
- Maximize the number of selected segments.

# Core techniques of QA Systems over KBs a Survey

└ Dissambiguation

└ Dissambiguation #8 - Structured Perceptron

Considering:

- Similarity of the phrase and the corresponding resource
- Popularity of a label for a resource
- Compatibility of the range and domain of a property with the arguments.

**Xser (xu2014a)** Solves ambiguity *fast*, *training needed*



2018-06-20

# Core techniques of QA Systems over KBs a Survey

└ Dissambiguation

└ Dissambiguation #9 - Summary

## Which techniques to choose

- Best results QALD 6, reported from UTQA ([pouran2016a](#)), whose method is unknown.
- Fastest method is gAnswer ([zou2014a](#)) which disambiguates in the phrase mapping step. ([subgraph matchin](#))
- Best results in QALD 4.1 & 5 by Xser ([zou2014a](#)) which uses a (Perceptron)

# Core techniques of QA Systems over KBs a Survey

## └ Query Construction

## └ Query Construction #1 - Issues

Construct a **SPARQL** query that reflects user question and gets the answer.

**Semantic Gap:** Issues with how the information is encoded in the KB. One cannot deduce how the information is stored from the question.

"Which countries are in the European Union?"

Could be encoded as:

*dbr:Greece dtp:member dbr:European\_Union*

*dbr:France dtp:member dbr:European\_Union*

or as:

*dbr:Greece dct:subj dbc:Member\_states\_of\_the\_European\_Union*

*dbr:France dct:subj dbc:Member\_states\_of\_the\_European\_Union*

**How to search correctly**

2018-06-20

# Core techniques of QA Systems over KBs a Survey

└ Query Construction

└ Query Construction #2

## Approaches

- Using templates
- Using information from the question analysis
- Using Semantic Parsers
- Using Machine Learning
- Using semantic information

# Core techniques of QA Systems over KBs a Survey

## └ Query Construction

## └ Query Construction #3 - Templates

Templates with parts of the query to be filled, in general by triples.

- **QAKIS** (calerio2012a) select queries with only one triple.
- **ISOFT** (park2014a) ASK over one triple, simple SELECT, COUNT and ORDER BY or FILTER.
- **PowerAqua** (lopez2012a) reduces the question to one or two triples ( $\leq 2$  predicates).

Very restricted questions, language is too rich, disambiguation is key

# Core techniques of QA Systems over KBs a Survey

## Query Construction

### Query Construction #4 - Question Analysis

Most systems get the form of the query in the question analysis and phrase mapping step.

- Freya, Intu3 (dima2014a) resources extracted in the phrase mapping step are combined into triples.
- DEANNA (yahya2013a) regex over POS tags in analysis step mapped to resources in phrase mapping step. ILP in disambiguation step to get the triples.
- gAnswer, QAnswer, RTV, SemGraphQA (zou2014a, ruseti2015a, giamone2013a, beaumont2015a) extract all the possible information from the dependency graph.

Question analysis because is done in that step

- gAnswer: The graph takes the form of the final query, resources associated with nodes and edges are fetched from the KB and used in the query.
- QAnswer: Scan dependency tree to find subgraph tokens corresponding to resources, many graphs, local disambiguation to get the best ones. Top ranked graph chosen for query.
- RTV dependency graph -  $\hat{\gamma}$  ordered list of alternated properties and non properties. Resources searched and disambiguated with HMM.

Special mention to Xser (best results), next slide



# Core techniques of QA Systems over KBs a Survey

## └ Query Construction

## └ Query Construction #5 - Question Analysis

**Xser (xu2014a)** 3 ML algorithms, two KB independent (on the question analysis phase), one KB dependant (on disambiguation step)

- First algorithm: determines segments of the question corresponding to variables, properties, instances and classes.
- Second algorithm: find dependencies between phrases. (**Stanford dependencies**, **PATTY**)
- Third algorithm: Disambiguation with a **Structured Perceptron**

# Core techniques of QA Systems over KBs a Survey

└ Query Construction

└ Query Construction #6 - Question Analysis

The problem with these methods is that they all **assume** that is **possible to deduce the structure of the SPARQL query from the structure of the question** without knowing how the knowledge is encoded in the KB.

2018-06-20

# Core techniques of QA Systems over KBs a Survey

└ Query Construction

└ Query Construction #7 - Semantic Parsing

Compose a grammar and use it to extract structure from the query.

- Grammatical Framework (GFMed marginian2017a)
- Feature-based Context Free Grammar (TR Discover, song2015a)
- Combinatorial Categorical Grammar (hakimov2015a)
- Lexical Tree Adjoint Grammar (TBSL unger2012a, BELA walter2012a)

# Core techniques of QA Systems over KBs a Survey

## Query Construction

## Query Construction #8 - Semantic Parsing

Question has to be well formulated. For each lexical item a corresponding semantic representation is needed. (ie married has to map with *dbo:spouse*). Learning corpus ([balkimov2015a](#)) or from POS tags ([unger2012a](#)). In general, **low recall**

Lexical item	Syntactic category	Semantic representation
Barack Obama	NP	dir :Barack.Obama
is	(S' NP) (S' NP)	At Is Fin
married to	(S NP <sub>1</sub> NP)	by Is dba :spouse[1, y]
Michelle Obama	NP	dir :Michelle.Obama

Figure: Semantic parsed question "Barack Obama is married to Michelle Obama"

2018-06-20

# Core techniques of QA Systems over KBs a Survey

└ Query Construction

└ Query Construction #9 - Machine Learning

## CASIA (h=2014a) (low recall & precision)

- Question Analysis step: extract features like position of a phrase and POS tags or the type of dependency in the dependency tree
- Phrase Mapping step: associate resources with phrase segments and extract more features
- Disambiguation step: MLN with extracted features to find most probable relation between segments and most probable mapping. **retrained for each KB**

2018-06-20

# Core techniques of QA Systems over KBs a Survey

## └ Query Construction

## └ Query Construction #10 - Semantic

SINA (shekarpour2015a), POMELO (hamon2014a), zhang2016a do not rely on the syntactic features of the question, instead the whole process is done based on the KB, just with semantic information.

Advantages:

- high recall, & precision

Disadvantages:

- computationally expensive
- does not respect user question syntax. No difference between "Who is the mother of Angela Merkel?" and "Angela Merkel is the mother of who?"

- Query Construction

## Query Construction #11 - Summary

2018-06-20

# Core techniques of QA Systems over KBs a Survey

└ Query Construction

└ Query Construction #12 - Summary

Which techniques to choose?

There is no clear way to do this. A good approach is to construct the query on the previous analysis steps, assuming structure can be extracted from the question.



# Core techniques of QA Systems over KBs a Survey

└ Conclusions

└ Conclusions

There exists many techniques for each part:

- Smart balance between different techniques lead to best results
- Tendency goes to Machine Learning

When possible, maximize:

- KB independence (pluggability)
- Performance (real time, scalability)
- Extracted Knowledge (implicit, enriches the context/domain)

- Data availability grows every day
- Domain and context are key, big future for attention based mechanisms