# Question Answering
# using Structured and Unstructured Data

Denis Savenkov
denis.savenkov@emory.edu
Emory University

advisor: Eugene Agichtein
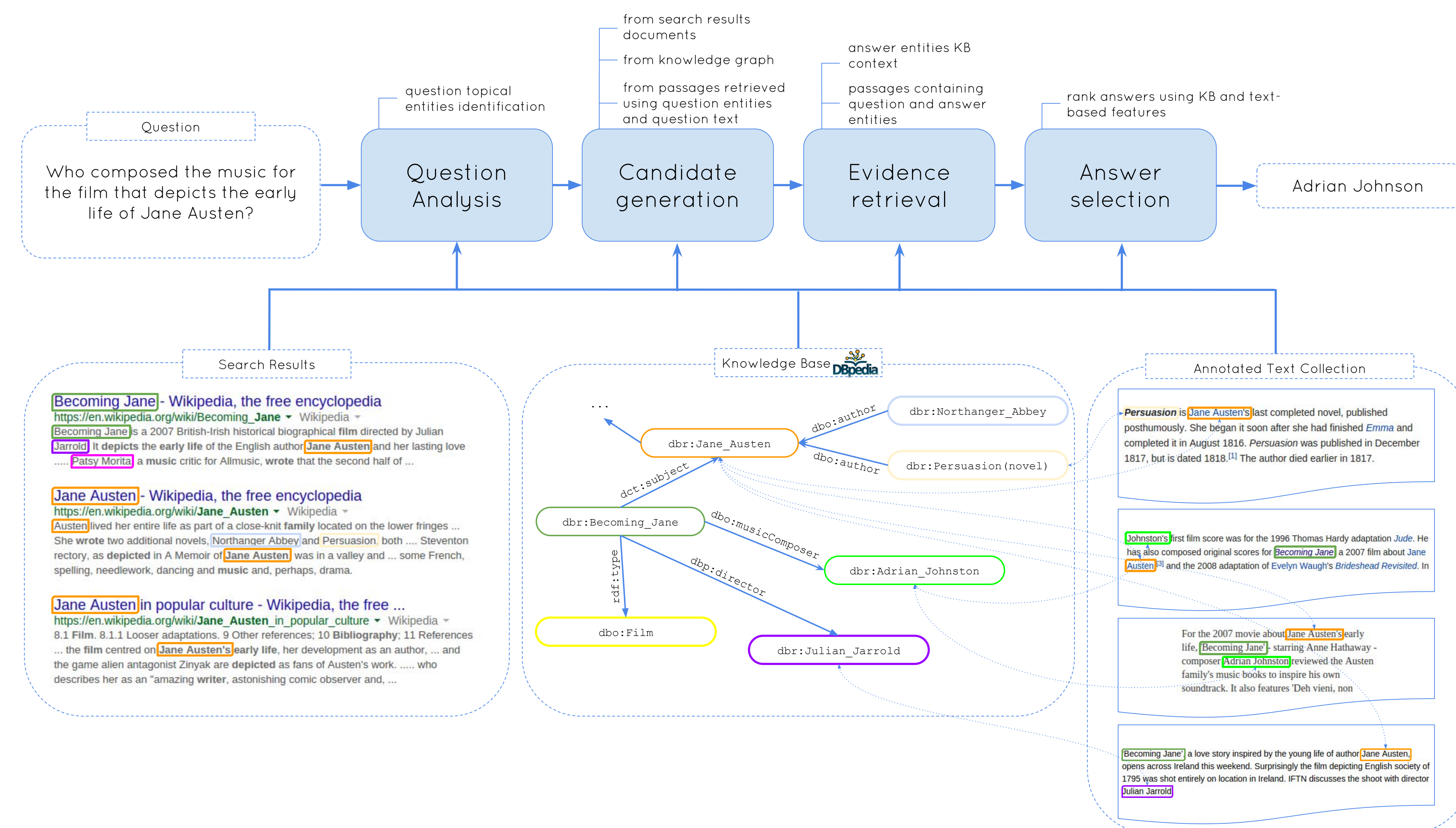eugene@emory.edu
Emory University

## Problem

➔ This thesis investigates methods for combining structured and unstructured data for answering a variety of user questions

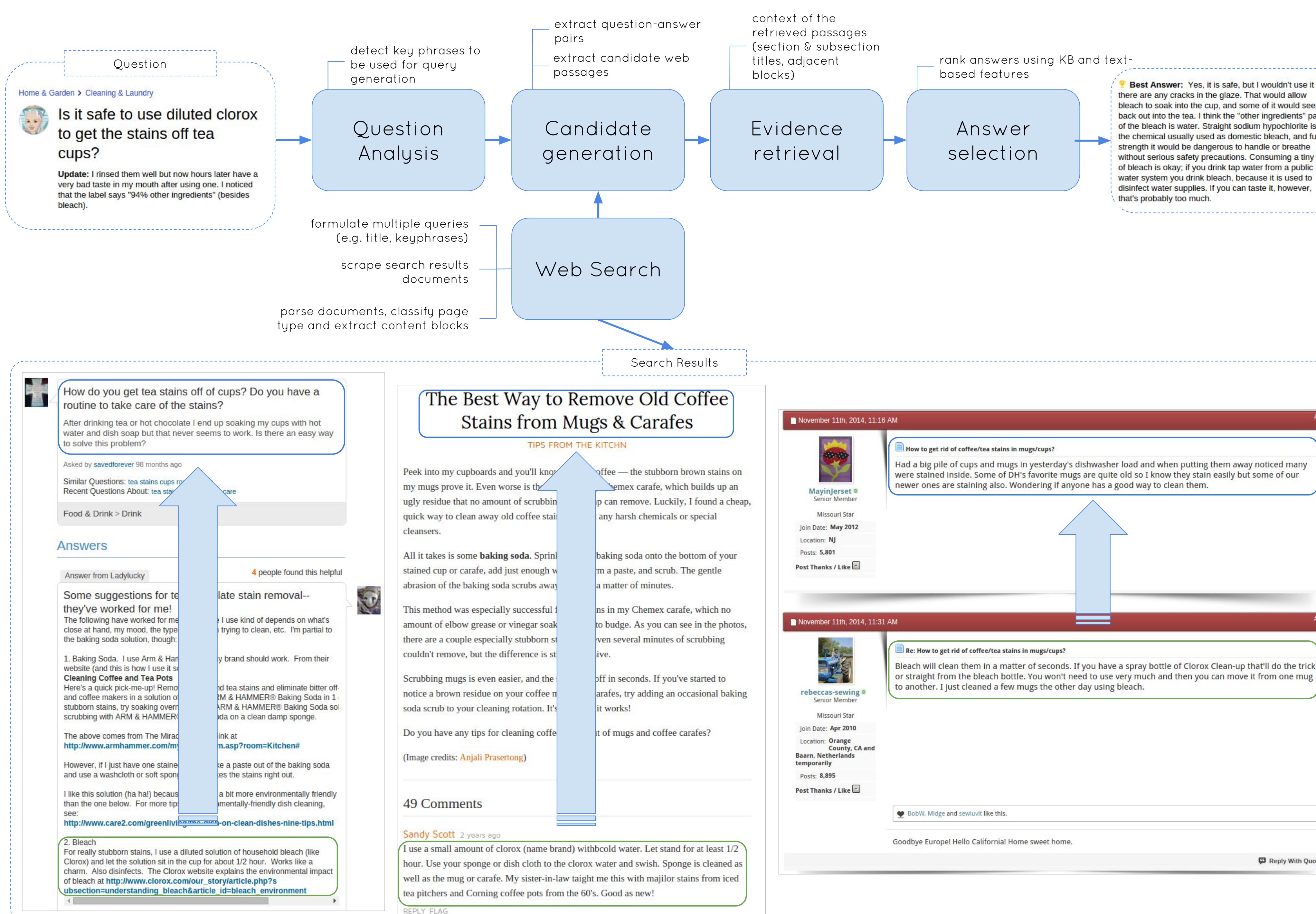| | Unstructured data | Structured data |
|---|---|---|
| | *Text collections* | *Knowledge Bases (KB)* |
| factoid questions | + easy to match against question text <br> + cover a variety of different information types <br> − each text phrase encodes a limited amount of information about mentioned entities | + aggregate all the information about entities <br> + allow complex queries over this data using special languages (e.g. SPARQL) <br> − hard to translate natural language questions into special query languages <br> − KBs are incomplete (missing entities, facts and properties) |
| | *Text collections* | *Question-Answer pairs* |
| non-factoid questions | + contain relevant information to a big chunk of user needs <br> − hard to match an answer to a question (lexical gap) | + easy to find a relevant answer by matching the corresponding questions <br> − cover a smaller subset of user information needs |

## Research Questions

1. What types of questions can be answered using text, KB or a combination of both?

2. How does semantic annotation of unstructured data compare to information extraction for question answering?
   - information extraction for KB construction vs open information extraction vs unstructured data annotation

3. How does a combination of structured and unstructured data sources improve each of the main QA system components: question analysis, candidate generation, evidence extraction and answer selection?

## Combination of Text & KB for factoid QA



## Web page structure for non-factoid QA



## Proposed Experiments

*Factoid questions*
➔ Datasets:
  ○ TREC QA [RQ1]
  ○ WebQuestions [RQ1]
  ○ Develop a new dataset [RQ1]
    ✓ questions derived from Yahoo! Answers WebScope collection
    ✓ answers are KB entities
    ✓ heuristic-based filtering: single sentence, no personal pronouns, no comparative and superlative adjectives, not future questions, and no keywords (e.g. recommend, suggest, will, etc)
    ✓ crowdsourcing for filtering to keep factoid, non-subjective questions and labeling answer entities, based on the answer text provided by a CQA user
➔ Compare against existing text-based, KB, Open IE and hybrid methods [RQ2, RQ3]
➔ Error analysis [RQ1, RQ3]

*Non-factoid questions*
➔ Datasets:
  ○ TREC LiveQA 2015 and the new 2016 shared task
➔ Comparison with last year system that doesn't use the information on the structure of web pages and other LiveQA participants [RQ3]

## Expected Results

1. new factoid question answering dataset
2. new approach to combining unstructured text and structured KB data, that:
   ✓ improves QA precision
   ✓ improves recall, including questions that could only be answered with a combination of text and KB data
3. new system for non-factoid question answering with improved performance due to better utilization of the information provided in web documents

## Open Questions

➔ How to incorporate other available sources of factual information, including semi-structured, e.g. tables, diagrams, etc?
➔ Non-factoid questions often include unique context information, that makes reusing information non effective. How a system can generate the answer given all potentially useful extracted information?
➔ How to construct KBs for non-factoid information needs?