

End-to-End Non-Factoid Question Answering with an Interactive Visualization of Neural Attention Weights



Andreas Rücklé and Iryna Gurevych
Ubiquitous Knowledge Processing Lab (UKP)
Department of Computer Science, Technische Universität Darmstadt
<https://www.ukp.tu-darmstadt.de>



1 Summary

A modular, extensible framework to visualize and compare attention mechanisms in (non-factoid) question answering

Non-Factoid Question Answering

- Questions that do *not necessarily* ask about facts; can often be answered with opinions, experiences, descriptions...
- We search for *existing* content on the web that can answer the question (e.g. by searching in CQA platforms)

Task: Answer Selection

- Given a question and a set of candidate answers, we rank the candidates for relevancy according to the question
- Attention-based models achieve state-of-the-art results

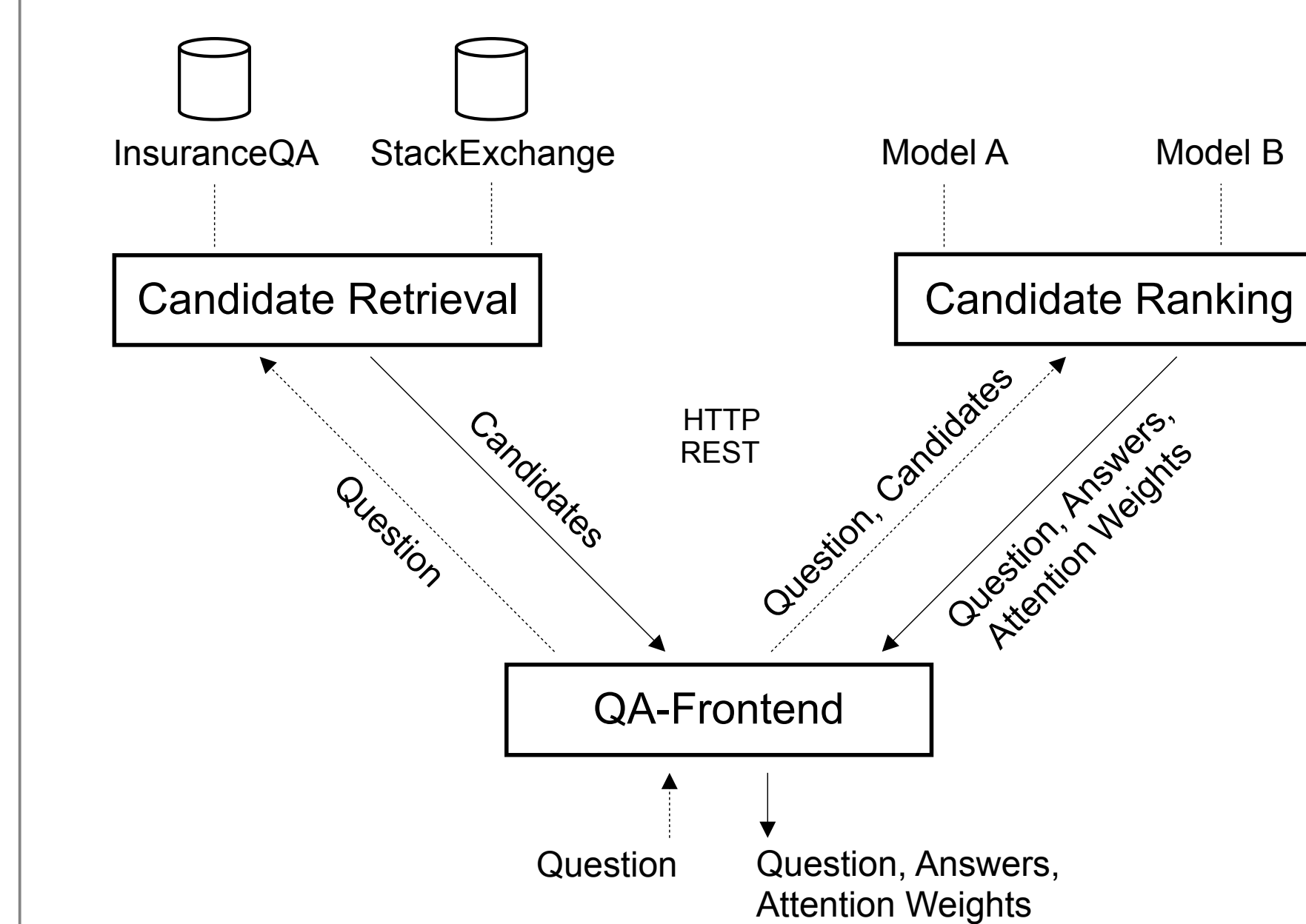
Challenge

- Understanding strengths and weaknesses of particular attention mechanisms is important
- Researchers usually plot attention weights for a number of predefined Q/A pairs within a dataset
- This is not interactive and makes the direct comparison of different approaches time-consuming

Our Contributions

1. We present a framework for the **interactive exploration** of different attention-based models
 - By transforming models to end-to-end QA systems
 - Supports one-way-, two-way-, and self-attention
2. Our UI allows **comparing** different approaches side by side

2 System Architecture



3 Candidate Retrieval

Approach

- We use ElasticSearch to index all answers of a dataset
- For an input question we query this index to obtain a set of candidate answers

Extensibility

- Our service implementation can easily be extended with new datasets through public interfaces
- We include readers for InsuranceQA and Stack-Exchange

4 Candidate Ranking

Approach

- We re-rank all candidates with (attention-based) NNs
- The results include attention weights for *every* Q/A pair → Question attention can be different for each answer

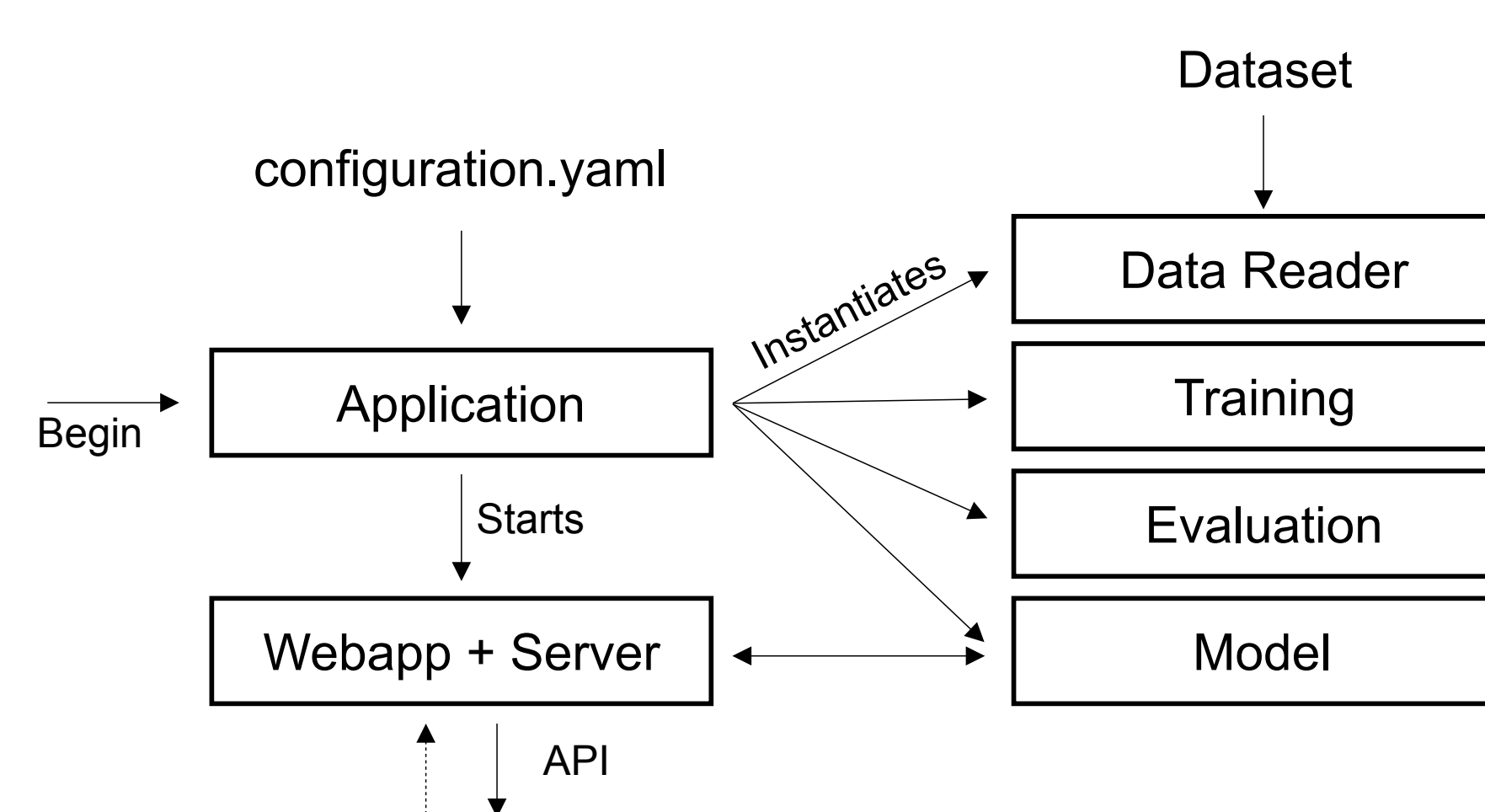
Extensibility

- The service implementation includes an extensible answer selection framework based on TensorFlow
- All components can freely be swapped and combined (dataset, model, training, evaluation)
- The framework is highly configurable through external YAML files (see the example on the right)

YAML Configuration:

```

1 data-module: data.insuranceqa.v2
2 model-module: model.ap_lstm
3 training-module: training.dynamic
4 evaluation-module: evaluation.default
5
6 data:
7   embeddings: data/glove.6B.100d.txt
8   insuranceqa: data/insuranceqa
9   lowercased: true
10  :
11
12 model:
13   lstm_cell_size: 141
14   margin: 0.2
15   trainable_embeddings: true
16  :
17
18 training:
19   batchsize: 20
20   epochs: 100
21   save_folder: checkpoints/ap_lstm
22   dropout: 0.3
23   optimizer: adam
24   initial_learning_rate: 0.001
25   scorer: accuracy
26  :
27
28 evaluation:
29   skip: true
    
```



Software Release



- The source code of our framework is **publicly available on GitHub**
- Includes a short documentation and an API reference

<https://github.com/UKPLab/acl2017-non-factoid-qa>

5 QA-Frontend

Ask arbitrary questions

Interactive visualization of Q/A attention

Choose between models

QA with Attention Visualization

is it possible to deduct health insurance from tax ?

Transparency Sensitivity (0.5): Threshold (0.5):

WEIGHT COMPARISON

LW-BiLSTM (Rücklé & Gurevych, IWCS 2017)

Attentive Pooling (Dos Santos et al., ACL 2016)

Importance: 0.12973786890506744