

## Background

- Community forums where user post questions and answers are becoming increasingly popular
- Arabic community question answering (cQA) received little attention in the past
- The SemEval-2016 Task 3 offered a cQA Arabic dataset in the medical domain, where:
  - Given a question and 30 related question-answer pairs, rank the pairs with respect to the original question
  - QA pairs annotated as *Direct*, *Relevant*, or *Irrelevant*
- Challenges: long texts, user-generated content, medical terminology, mixed standard and colloquial language
- Example:
  - Q': ما هي اعراض الاصابة بمتلازمة القولون المتهيج العصبي
  - Q: ما هي أعراض القولون العصبي
  - A: تتباين علامات وأعراض الإصابة بمتلازمة القولون العصبي بين الأشخاص [...] كالتالي : - ألم أو تشنج البطن . - الشعور بالانتفاخ . - امتلاء [...] البطن بالغازات . - الاسهال أو الإمساك

## Preprocessing with Keyword Extraction

- Text preprocessing** in order to deal with several challenges:
  - Long texts: average question/answer length is 50/120 words
  - Rich morphology: multiple surface forms per lemma
  - Latin terminology in the medical domain
- Keyword extraction** with TextRank
  - Treat every thread as a document
  - Form a graph where nodes are word types and edges represent co-occurrence in  $N$ -sized window
  - Compute importance weight iteratively and keep top  $P\%$  of words
- Lemmatization**: we apply MADA for finding lemmas and part-of-speech tags
- Stop-word removal**: we keep only content words, Latin words, and words with no morphological analysis.
- Sentence Splitting**: We use the splitter from Stanford's core NLP toolkit.
- Syntactic parsing**: We use the Berkeley parser to get the syntactic tree for each each sentence.

## Experiments

- Preprocessing settings**
  - No preprocessing
  - Only keeping content lemmas
  - Only content lemmas and keyword extraction with TextRank params  $N=3, P=5$
  - Same, with TextRank params  $N=4, P=1$
- Tree kernels settings**
  - ConvKN-contrastive1: only basic features
  - ConvKN-contrastive2: MT features
- ConvKN-primary**: basic features + tree kernels

	MAP	AvgRec	MRR	P	R	F1	ACC
Rand	29.79	31.00	33.71	19.53	20.66	20.08	68.35
a	38.33	42.09	43.75	20.38	<b>96.95</b>	33.68	26.58
b	39.98	43.68	46.41	26.26	68.39	37.95	57.00
c	<b>45.50</b>	<b>50.13</b>	<b>52.55</b>	<b>28.55</b>	64.53	<b>39.58</b>	<b>62.10</b>
i	44.94	49.72	51.58	<b>62.96</b>	2.40	4.62	<b>80.95</b>
iii	42.95	47.61	49.55	27.20	<b>74.40</b>	39.84	56.76
i-iv	<b>45.83</b>	<b>51.01</b>	<b>53.66</b>	34.45	52.33	<b>41.55</b>	71.67

## Feature Representation

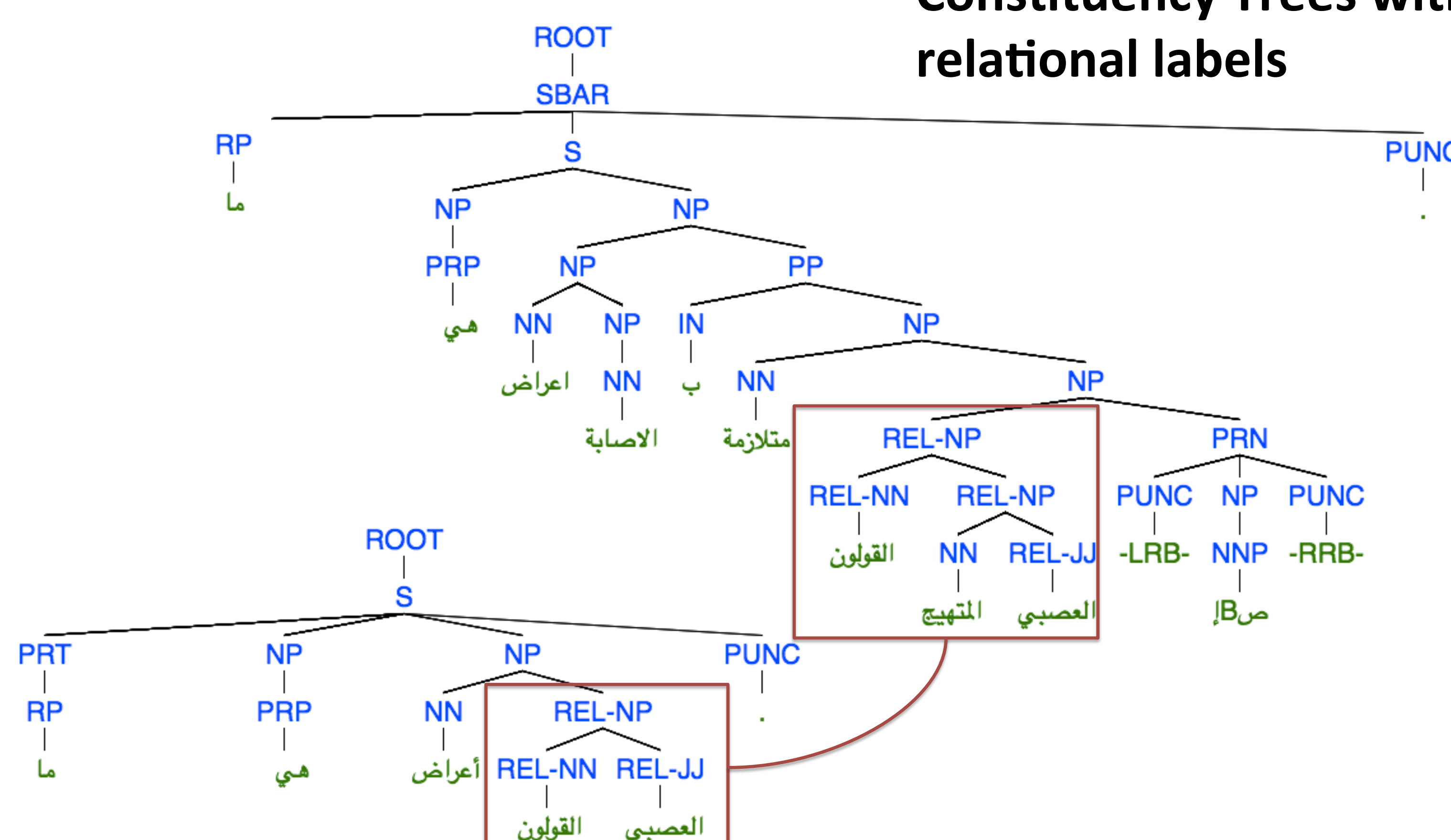
- Given a question  $q'$  and a related question-answer pair  $q-a$ , compute features between the pairs  $q'-q$  and  $q'-a$
- Text-based features**
  - Various text-similarity metrics such as Longest Common Substring, Longest Common Subsequence, Greedy String Tiling, etc. (Belinkov et al. 2015)
- Vector-based features**
  - Vector representations of closest pairs of words or sentences in  $q'-q$  and  $q'-a$
  - Word vectors computed from Arabic Gigaword and medical domain raw data using *Word2Vec*
  - Sentence representation is average of word vectors
- Machine translation evaluation features**
  - BIEU, TER, Meteor

## Tree Kernels

### Syntactic Tree Kernels

$$K((t_1, t_2), (u_1, u_2)) = TK(t_1, u_1) + TK(t_2, u_2)$$

### Constituency Trees with relational labels



## Future Work

- Combine keyword extraction with tree kernels
- How do deal with grammatical structure after keyword extraction?
- Automatically detecting the most important sentences to be matched with the tree kernels

## References

- Belinkov et al. 2015. VectorSLU: A Continuous Word Vector Approach to Answer Selection in Community Question Answering Systems.
- Mohtarami et al. 2016. SLS at SemEval-2016 Task 3: Neural-based Approaches for Ranking in Community Question Answering.
- ConvKN at SemEval-2016 Task 3: Answer and Question Selection for Question Answering on Arabic and English Fora