



PNEL: Pointer Network Based End-To-End Entity Linking over Knowledge Graphs

Debayan Banerjee¹(✉), Debanjan Chaudhuri², Mohnish Dubey²,
and Jens Lehmann^{2,3}

¹ Language Technology Group, Universität Hamburg, Hamburg, Germany
banerjee@informatik.uni-hamburg.de

² Fraunhofer IAIS, Bonn/Dresden, Germany

{debanjan.chaudhuri,mohnish.dubey}@iais.fraunhofer.de

³ Smart Data Analytics Group, University of Bonn, Bonn, Germany
jens.lehmann@cs.uni-bonn.de

Abstract. Question Answering systems are generally modelled as a pipeline consisting of a sequence of steps. In such a pipeline, Entity Linking (EL) is often the first step. Several EL models first perform span detection and then entity disambiguation. In such models errors from the span detection phase cascade to later steps and result in a drop of overall accuracy. Moreover, lack of gold entity spans in training data is a limiting factor for span detector training. Hence the movement towards end-to-end EL models began where no separate span detection step is involved. In this work we present a novel approach to end-to-end EL by applying the popular Pointer Network model, which achieves competitive performance. We demonstrate this in our evaluation over three datasets on the Wikidata Knowledge Graph.

Keywords: Entity Linking · Question Answering · Knowledge Graphs · Wikidata

1 Introduction

Knowledge Graph based Question Answering (KGQA) systems use a background Knowledge Graph to answer queries posed by a user. Let us take the following question as an example (Fig. 1): *Who founded Tesla?*. The standard sequence of steps for a traditional Entity Linking system is as follows: The system tries to identify *Tesla* as a span of interest. This task is called Mention Detection (MD) or Span Detection. Then an attempt is made to link it to the appropriate entity in the Knowledge Base. In this work we focus on Knowledge Bases in the

D. Chaudhuri and M. Dubey—Equal contribution.

© Springer Nature Switzerland AG 2020

J. Z. Pan et al. (Eds.): ISWC 2020, LNCS 12506, pp. 21–38, 2020.

https://doi.org/10.1007/978-3-030-62419-4_2

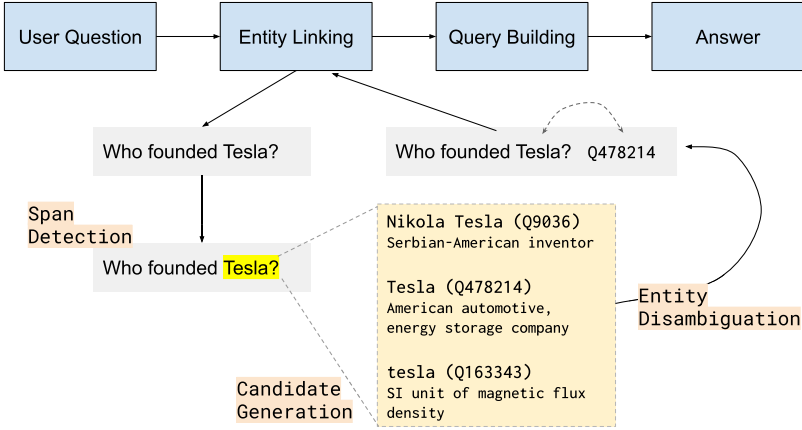


Fig. 1. Illustrating the use of Entity Linking in KGQA system.

form of graphs, hence the entity linker in this case tries to link *Tesla* to the appropriate node in the graph. For a human, it is evident that the question is looking for a person’s name who created an organisation named *Tesla*, since the text contains the *relation* *founded*. Hence, it is important that the entity linker understands the same nuance and ignores other entity nodes in the Knowledge Graph which also contain *Tesla* in their labels, e.g., *Nikola Tesla* (Q9036, *Serbian-American inventor*), *tesla* (Q163343, *SI unit*) when considering the example of the Wikidata knowledge graph. The task of ignoring the wrong candidate nodes, and identifying the right candidate node instead, is called *Entity Disambiguation (ED)*. The cumulative process involving Mention Detection and Entity Disambiguation is called *Entity Linking (EL)*.

Typically, the MD and ED stages are implemented by different machine learning models which require separate training. Especially for the MD part, sentences with marked entity spans are a requirement. In practice, such data is not easily available. Moreover, errors introduced by the MD phase cascade on to the ED phase. Hence, a movement towards end-to-end Entity Linkers began [11, 26]. Such systems do not require labelled entity spans during training. In spite of the benefits of end-to-end models some challenges remain: Due to the lack of a span detector at the initial phase, each word of the sentence needs to be considered as an entity candidate for the disambiguation which leads to the generation of a much larger number of entity candidates. To re-rank these candidates a large amount of time is consumed, not just in processing the features of the candidates, but also in compiling their features.

In this work, we remain cognizant of these challenges and design a system that completely avoids querying the Knowledge Graph during runtime. PNEL (Pointer Network based Entity Linker) instead relies on pre-computed and pre-indexed TransE embeddings and pre-indexed entity label and description text as the only set of features for a given candidate entity. We demonstrate that this

produces competitive performance while maintaining lower response times when compared to another end-to-end EL system, VCG [26].

While there is a wide variety of KG embeddings to choose from, we confine our experiments to pre-computed TransE over Wikidata supplied by PyTorch-BigGraph [13]. Our choice was based on the popularity and ease of availability of these embeddings.

Traditionally, the Knowledge Graphs of choice for Question Answering research have been DBpedia [12], Freebase [2] and YAGO [27]. However, in recent times Wikidata [30] has received significant attention owing to the fact that it covers a large number of entities (DBpedia 6M¹, Yago 10M², Freebase 39M³, Wikidata 71M⁴). DBpedia, YAGO and Wikidata source their information from Wikipedia, however DBpedia and YAGO filter out a large percentage of the original entities, while Wikidata does not. While Wikidata has a larger number of entities it also adds to noise which is a challenge to any EL system. Wikidata also allows direct edits leading to up-to-date information, while DBpedia depends on edits performed on Wikipedia. Freebase has been discontinued and a portion of it is merged into Wikidata [19]. Moreover DBpedia now extracts data directly from Wikidata, apart from Wikipedia⁵ [8]. Hence, we decide to base this work on the Wikidata knowledge graph and the datasets we evaluate on are all based on Wikidata.

In this work our **contributions** are as follows:

1. PNEL is the first approach that uses the Pointer Network model for solving the End-to-End Entity Linking problem over Knowledge Graphs, inspired by the recent success of pointer networks for convex hull and generalised travelling salesman problems.
2. We are the first work to present baseline results for the entire LC-QuAD 2.0 [5] test set.
3. Our approach produces state-of-the-art results on the LC-QuAD 2.0 and SimpleQuestions datasets.

The paper is organised into the following sections: (2) Related Work, outlining some of the major contributions in entity linking used in question answering; (3) PNEL, where we discuss the pointer networks and the architecture of PNEL (4) Dataset used in the paper (5) Evaluation, with various evaluation criteria, results and ablation test (6) Error Analysis (7) Discussion and future direction.

¹ <https://wiki.dbpedia.org/develop/datasets/latest-core-dataset-releases>.

² <https://www.mpi-inf.mpg.de/departments/databases-and-information-systems/research/yago-naga/yago/>.

³ https://developers.google.com/freebase/guide/basic_concepts#topics.

⁴ <https://www.wikidata.org/wiki/Wikidata:Statistics>.

⁵ <https://databus.dbpedia.org/dbpedia/wikidata>.

2 Related Work

DBpedia Spotlight [16] is one of the early works for entity linking over DBpedia. It first identifies a list of surface forms and then generates entity candidates. It then disambiguates the entity based on the surrounding context. In spite of being an early solution, it still remains one of the strongest candidates in our own evaluations, at the same time it has low response times. Compared to PNEL it lags behind in precision significantly. S-MART [31] generates multiple regression trees and then applies sophisticated structured prediction techniques to link entities to resources. S-MART performs especially well in recall on WebQSP in our evaluations and the reason seems to be that they perform more complex information extraction related tasks during entity linking, e.g., “Peter Parker” span fetches “Stan Lee”⁶. However compared to PNEL it has low precision.

The journey towards end-to-end models which combine MD and ED in one model started with attempts to build feedback mechanisms from one step to the other so that errors in one stage can be recovered by the next stage. One of the first attempts, Sil et al. [25], use a popular NER model to generate extra number of spans and let the linking step take the final decisions. Their method however depends on a good mention spotter and the use of hand engineered features. It is also unclear how linking can improve their MD phase. Later, Luo et al. [15] developed competitive joint MD and ED models using semi-Conditional Random Fields (semi-CRF). However, the basis for dependency was not robust, using only type-category correlation features. The other engineered features used in their model are either NER or ED specific. Although their probabilistic graphical model allows for low complexity learning and inference, it suffers from high computational complexity caused by the usage of the cross product of all possible document spans, NER categories and entity assignments. Another solution J-NERD [18] addresses the end-to-end task using engineered features and a probabilistic graphical model on top of sentence parse trees. EARL [6] makes some rudimentary attempts towards a feedback mechanism by allowing the entity and relation span detector to make a different choice based on classifier score in the later entity linking stage, however it is not an End-to-End model.

Sorokin et al. [26] is possibly the earliest work on end-to-end EL. They use features of variable granularities of context and achieve strong results on Wikidata that we are yet unable to surpass on WebQSP dataset. More recently, Kolitsas et al. [11] worked on a truly end-to-end MD (Mention Detection) and ED (Entity Disambiguation) combined into a single EL (Entity Linking) model. They use context-aware mention embeddings, entity embeddings and a probabilistic mention - entity map, without demanding other engineered features. Additionally, there are a few recent works on entity linking for short text on Wikidata [30], which is also the area of focus of PNEL. OpenTapioca [4] works on a limited number of classes (humans, organisations and locations) when compared to PNEL, but is openly available both as a demo and as code and is

⁶ <https://github.com/UKPLab/starsem2018-entity-linking/issues/8#issuecomment-566469263>.

lightweight. Falcon 2.0 [22] is a rule based EL solution on Wikidata which is openly available and fast, but it requires manual feature engineering for new datasets. Sevgili et al. [24] performs ED using KG entity embeddings (DeepWalk [20]) on Wikidata, but they rely on an external MD solution. PNEL and Sorokin et al. both use TransE entity embeddings and also perform MD and ED end-to-end in a single model. Sorokin et al. has a more complex architecture when compared to PNEL. Apart from using TransE embeddings, they fetch neighbouring entities and relations on the fly during EL, which is a process PNEL intentionally avoids to maintain lower response times. KBPearl [14] is a recent work on KG population which also targets entity linking as a task for Wikidata. It uses dense sub-graphs formed across the document text to link entities. It is not an end-to-end model but is the most recent work which presents elaborate evaluation on Wikidata based datasets, hence we include it in evaluations.

We also include QKBFly [17] and TagME [7] in our evaluations because KBPearl includes results for these systems on a common dataset (LC-QuAD 2.0). QKBFly performs on-the-fly knowledge base construction for ad-hoc text. It uses a semantic-graph representation of sentences that captures per-sentence clauses, noun phrases, pronouns, as well as their syntactic and semantic dependencies. It retrieves relevant source documents for entity centric text from multiple sources like Wikipedia and other news websites. TagME is an older system that spots entity spans in short text using a Lucene index built out of anchor text in Wikipedia. It then performs a mutual-voting based disambiguation process among the candidates and finishes with a pruning step.

3 PNEL

PNEL stands for Pointer Network based Entity Linker. Inspired by the use case of Pointer Networks [29] in solving the convex hull and the generalised travelling salesman problems, this work adapts the approach to solving entity linking. *Conceptually, each candidate entity is a point in an euclidean space, and the pointer network finds the correct set of points for the given problem.*

3.1 Encoding for Input

The first step is to take the input sentence and vectorise it for feeding into the pointer network. We take varying length of n-grams, also called n-gram tiling and vectorise each such n-gram.

Given an input sentence $S = \{s_1, s_2 \dots s_n\}$ where s_k is a token (word) in the given sentence, we vectorise s_k to v_k , which is done in the following manner:

1. Take the following 4 n-grams: $[s_k], [s_{k-1}, s_k], [s_k, s_{k+1}], [s_{k-1}, s_k, s_{k+1}]$.
2. For each such n-gram find the top L text matches in the entity label database. We use the OKAPI BM25 algorithm for label search.
3. For each such candidate form a candidate vector comprising of the concatenation of the following features

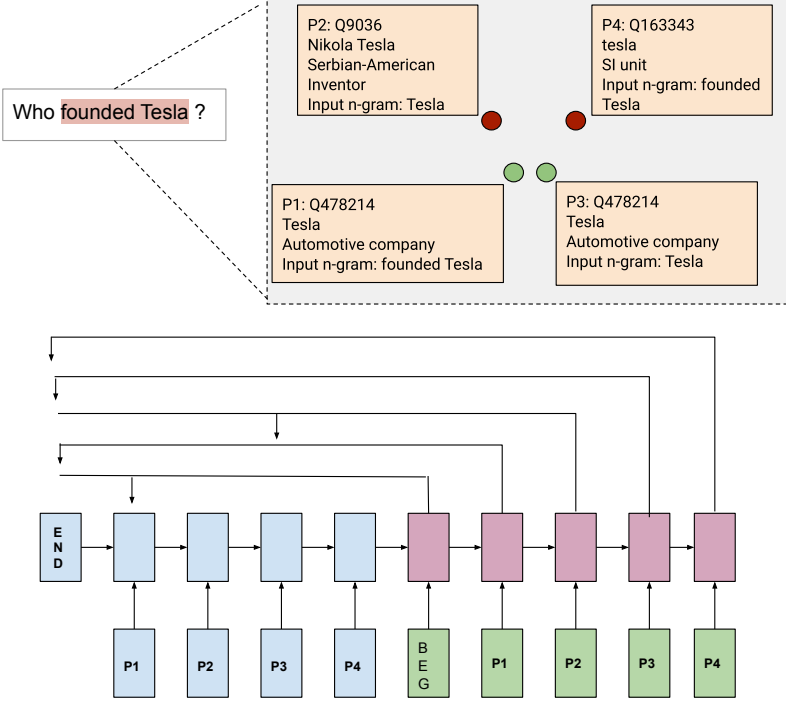


Fig. 2. The red and green dots represent entity candidate vectors for the given question. The green vectors are the correct entity vectors. Although they belong to the same entity they are not the same dots because they come from different n-grams. At each time step the Pointer Network points to one of the input candidate entities as the linked entity, or to the END symbol to indicate no choice. (Color figure online)

- (a) R_{kl} = Rank of entity candidate in text search (length 1), where $1 \leq l \leq L$.
- (b) $ngramlen$ = The number of words in the current n-gram under consideration where $1 \leq ngramlen \leq 3$ (length 1).
- (c) k = The index of the token s_k (length 1).
- (d) pos_k = A one-hot vector of length 36 denoting the PoS tag of the word under consideration. The 36 different tags are as declared in the Penn Treebank Project [23] (length 36).
- (e) $EntEmbed_{kl}$ = TransE Entity Embedding (length 200).
- (f) $SentFTEEmbed$ = fastText embedding of sentence S (length 300), which is a mean of the embeddings of the tokens of S . In some sense this carries within it the problem statement.
- (g) $TokFTEEmbed_k$ = fastText embedding of token s_k (length 300). Addition of this feature might seem wasteful considering we have already added the sentence vector above, but as shown in the ablation experiment in Table 6, it results in an improvement.
- (h) $DescriptionEmbed_{kl}$ = fastText embedding of the Wikidata description for entity candidate kl (length 300).

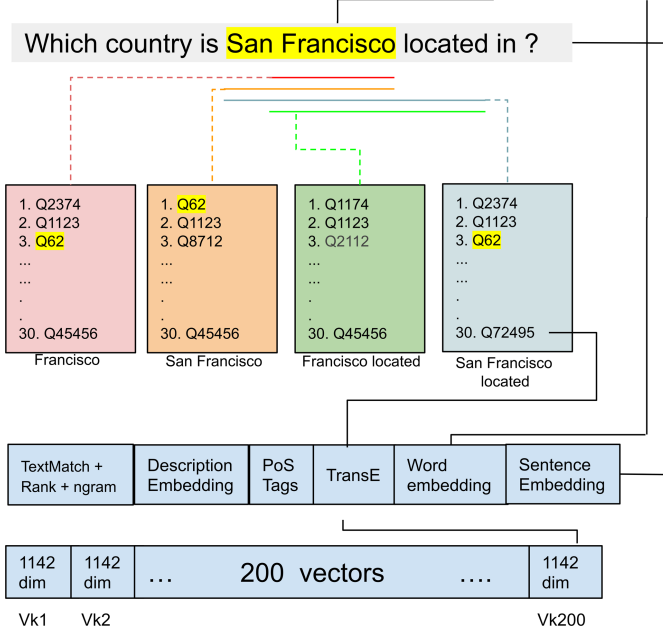


Fig. 3. The word “Francisco” is vectorised in the following manner: 4 n-grams represented by the underlines are considered and searched against an entity label database. The top 50 search results are depicted for each of the n-grams resulting in 200 candidates. For the entity Q72495, for example, we fetch its TransE embedding, add its text search rank, n-gram length, word index position as features. Additionally we also append the fastText embedding for “Francisco” and the entire fastText embedding for the sentence (average of word vectors) to the feature. We then append the fast-Text embeddings for the label and description for this entity. Hence we get a 1142 dimensional vector V_{k120} **corresponding to entity candidate Q72495**. For all 200 candidate entities for “Francisco”, we have a sequence of two hundred 1142 dimensional vectors as input to the pointer network. For the sentence above which has 7 words, this results in a final sequence of $7 \times 200 = 1400$ vectors each of length 1142 as input to our pointer network. Any one or more of these vectors could be the correct entities.

- (i) $TextMatchMetric_{kl}$ = This is a triple of values, each ranging from 0 to 100, that measures the degree of text match between the token under consideration s_k and the label of the entity candidate kl . The three similarity matches are *simple ratio*, *partial ratio*, and *token sort ratio*. In case of *simple ratio* the following pair of text corresponds to perfect match: "Elon Musk" and "Elon Musk". In case of *partial ratio* the following pair of text corresponds to a perfect match: "Elon Musk" and "Musk". In case of *token sort ratio* the following pair of text corresponds to a perfect match: "Elon Musk" and "Musk Elon" (length 3).

For each token s_k we have an expanded sequence of token vectors, comprising of 4 n-grams, upto 50 candidates per n-gram, where each vector is of length

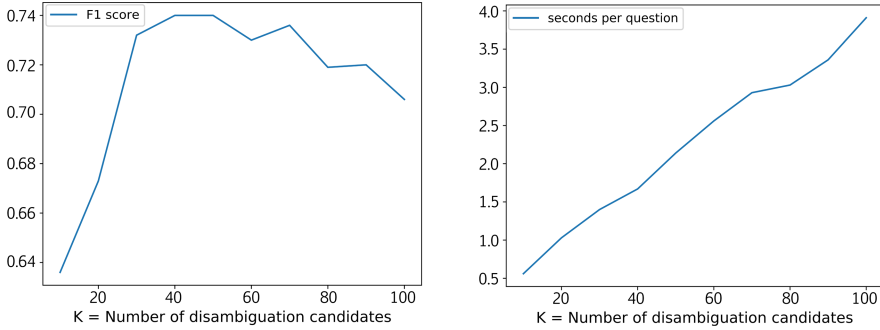


Fig. 4. K = Number of search candidates per n-gram. On the left: K vs F1 score on a set of 100 WebQSP test questions, with average word length of 6.68. F1 is maximum for $K=40$ and 50. On the right: K vs time taken for PNEL to return a response. The relationship appears to be close to linear.

1142. Hence each token s_k is transformed into $4 \times 50 = 200$ vectors, each a 1142 length vector (see Fig. 3). We may denote this transformation as $s_k \rightarrow \{v_{k1}, v_{k2} \dots v_{k200}\}$. Note that there may be less than 50 search results for a given token so there may be less than 200 entity candidates in the final vectorisation. Each of these v_k vectors is an entity candidate (Fig. 3).

3.2 Training

For the entire sentence, a sequence of such vectors is provided as input to the pointer network. During training the labels for the given input sequence are the index numbers of the correct entities in the input sequence. Note that the same entity appears multiple times because of n-gram tiling. During each decoding time step the decoder produces a softmax distribution over the input sequence (see Fig. 2), which in our implementation has a maximum sequence length of 3000. Additionally the **BEGIN**, **END**, **PAD** symbols add to a total of 3003 symbols to softmax over. The cross entropy loss function is averaged over the entire output sequence of labels and is considered the final loss for the entire input sequence.

3.3 Network Configuration

We use a single layer bi-LSTM [9] pointer network with 512 hidden units in a layer and an attention size of 128. Addition of an extra layer to the network did not result in an improvement. The Adam optimizer [10] was used with an initial learning rate of 0.001. A maximum input sequence length of 3000 and a maximum output length of 100 were enforced.

4 Datasets

For reasons explained in Sect. 1 we evaluate on Wikidata based datasets. We use the following:

- **WebQSP:** We use the dataset released by Sorokin et al. [26] where the original WebQSP dataset by Yih et al. [32], which was based on Freebase, has been adapted and all Freebase IDs converted to their respective Wikidata IDs. WebQSP contains questions that were originally collected for the WebQuestions dataset from web search logs (Berant et al. [1]). WebQSP is a relatively small dataset consisting of 3098 train 1639 test questions which cover 3794 and 2002 entities respectively. The dataset has a mixture of simple and complex questions. We found some questions in the test set that had failed Freebase to Wikidata entity ID conversions. We skipped such questions during PNEL’s evaluation.
- **SimpleQuestions:** To test the performance of PNEL on simple questions, we choose SimpleQuestions [3], which as the name suggests, consists only of Simple Questions. The training set has more than 30,000 questions while the test set has close to 10,000 questions. This dataset was also originally based on Freebase and later the entity IDs were converted to corresponding Wikidata IDs. However out of the 10,000 test questions only about half are answerable on the current Wikidata.
- **LC-QuAD 2.0:** Unlike the first two datasets, LC-QuAD 2.0 [5] is based on Wikidata since its inception and is also the most recent dataset of the three. It carries a mixture of simple and complex questions which were verbalised by human workers on Amazon Mechanical Turk. It is a large and varied dataset comprising of 24180 train questions and 6046 test questions which cover 33609 and 8417 entities respectively.

5 Evaluation

In this section we evaluate our proposed model(s) against different state-of-the-art methods for KGQA. As notations, PNEL-L stands for PNEL trained on LC-QuAD 2.0. PNEL-W and PNEL-S stand for PNEL trained on WebQSP and SimpleQuestions respectively.

5.1 Experiment 1: EL over KBPearl Split of LC-QuAD 2.0 Test Set

Objective: The purpose of this experiment is to benchmark PNEL against a large number of EL systems, not just over Wikidata, but also other KBs.

Method: The results are largely taken from KBPearl. PNEL is trained on the LC-QuAD 2.0 training set. For a fair comparison, the systems are tested on the 1294 questions split of test set provided by KBPearl. We train PNEL for 2 epochs.

Remarks: Results for Falcon 2.0 and OpenTapioca were obtained by accessing their live API. The original Falcon 2.0 paper provides an F1 of 0.69 on 15% of randomly selected questions from a combination of the train and test splits of the dataset. Several systems in the table below do not originally produce Wikidata entity IDs, hence the authors of KBpearl have converted the IDs to corresponding Wikidata IDs.

Analysis: As observed from the results in Table 1, PNEL outperforms all other systems on this particular split of LC-QuAD 2.0 dataset.

Table 1. Evaluation on KBPearl split of LC-QuAD 2.0 test set

Entity linker	Precision	Recall	F1
Falcon [21]	0.533	0.598	0.564
EARL [6]	0.403	0.498	0.445
Spotlight [16]	0.585	0.657	0.619
TagMe [7]	0.352	0.864	0.500
OpenTapioca [4]	0.237	0.411	0.301
QKBfly [17]	0.518	0.479	0.498
Falcon 2.0	0.395	0.268	0.320
KBPearl-NN	0.561	0.647	0.601
PNEL-L	0.803	0.517	0.629

5.2 Experiment 2: EL over Full LC-QuAD 2.0 Test Set

Objective: The objective of this experiment is to compare systems that return Wikidata IDs for the EL task.

Method: We train PNEL on LC-QuAD 2.0 train set and test on all 6046 questions in test set. PNEL was trained for 2 epochs.

Remarks: Results for competing systems were obtained by accessing their live APIs. We choose systems that return Wikidata IDs.

Analysis: As seen in Table 2, similar to the previous experiment, PNEL performs the best on the LC-QuAD 2.0 test set.

5.3 Experiment 3: EL over WebQSP Test Set

Objective: Benchmark against an end-to-end model that returns Wikidata IDs.

Table 2. Evaluation on LC-QuAD 2.0 test set

Entity linker	Precision	Recall	F1
VCG [26]	0.516	0.432	0.470
OpenTapioca [4]	0.237	0.411	0.301
Falcon 2.0	0.418	0.476	0.445
PNEL-L	0.688	0.516	0.589

Method: Train and test PNEL on WebQSP train and test sets respectively. PNEL is trained for 10 epochs.

Remarks: Results for the competing systems were taken from Sorokin et al. [26].

Table 3. Evaluation on WebQSP

Entity linker	Precision	Recall	F1
Spotlight	0.704	0.514	0.595
S-MART [31]	0.666	0.772	0.715
VCG [26]	0.826	0.653	0.730
PNEL-L	0.636	0.480	0.547
PNEL-W	0.886	0.596	0.712

Analysis: As seen in Table 3 PNEL comes in third best in this experiment, beaten by VCG and S-MART. S-MART has high recall because it performs semantic information retrieval apart from lexical matching for candidate generation, as explained in Sect. 2. VCG is more similar to PNEL in that it is also an end-to-end system. It has higher recall but lower precision than PNEL.

5.4 Experiment 4: EL over SimpleQuestions Test Set

Objective: Benchmark systems on the SimpleQuestions Dataset.

Method: Train and test PNEL on SimpleQuestions train and test sets respectively. PNEL is trained for 2 epochs.

Remarks: We extended the results from Falcon 2.0 [22].

Table 4. Evaluation on SimpleQuestions

Entity linker	Precision	Recall	F1
OpenTapioca [4]	0.16	0.28	0.20
Falcon 2.0	0.38	0.44	0.41
PNEL-L	0.31	0.25	0.28
PNEL-S	0.74	0.63	0.68

Analysis: As seen in Table 4, PNEL outperforms the competing systems both in precision and recall for SimpleQuestions dataset. As observed, PNEL has the best precision across all datasets, however, recall seems to be PNEL’s weakness.

5.5 Experiment 5: Candidate Generation Accuracy

Objective: The purpose of this experiment is to see what percentage of correct entity candidates were made available to PNEL after the text search phase. This sets a limit on the maximum performance that can be expected from PNEL.

Remarks: PNEL considers each token a possible correct entity, but since it only considers top-K text search matches for each token, it also loses potentially correct entity candidates before the disambiguation phase. The results in Table 5 are for K = 30.

Table 5. Entity Candidates available post label search

Dataset	PNEL (%)
WebQSP	73
LC-QuAD 2.0	82
SimpleQuestions	90

5.6 Experiment 6: Ablation of Features Affecting Accuracy

Objective: We present an ablation study on the WebQSP dataset to understand the importance of different feature vectors used in the model.

Analysis: As seen in Table 6 it appears that the most important feature is the TransE entity embedding, which implicitly contains the entire KG structure information. On removing this feature there is drop in F1 score from 0.712 to 0.221. On the other hand the least important feature seem to be the description embedding. Removal of this feature merely leads to a drop in F1 from 0.712 to 0.700. A possible reason is that the Text Search Rank potentially encodes

Table 6. Ablation test for PNEL on WebQSP test set features

Sentence embed.	Word embed.	Descript. embed.	TransE	PoS tags	Text rank	n-gram length	Text match metric	F1 score
✓	✓	✓	✓	✓	✓	✓	✓	0.712
	✓	✓	✓	✓	✓	✓	✓	0.554
✓		✓	✓	✓	✓	✓	✓	0.666
✓	✓		✓	✓	✓	✓	✓	0.700
✓	✓	✓		✓	✓	✓	✓	0.221
✓	✓	✓	✓		✓	✓	✓	0.685
✓	✓	✓	✓	✓		✓	✓	0.399
✓	✓	✓	✓	✓	✓		✓	0.554
✓	✓	✓	✓	✓	✓	✓		0.698

significant text similarity information, and TransE potentially encodes other type and category related information that description often adds. Removal of the Text Search Rank also results in a large drop in F1 reaching to 0.399 from 0.712.

5.7 Experiment 7: Run Time Evaluation

Objective: We look at a comparison of run times across the systems we have evaluated on

Table 7. Time taken per question on the WebQSP dataset of 1639 questions

System	Seconds	Target KG
VCG	8.62	Wikidata
PNEL	3.14	Wikidata
Falcon 2.0	1.08	Wikidata
EARL	0.79	DBpedia
TagME	0.29	Wikipedia
Spotlight	0.16	DBpedia
Falcon	0.16	DBpedia
OpenTapioca	0.07	Wikidata

Analysis: QKBFly and KBPearl are off-line systems, requiring separate steps for entity candidate population and entity linking, hence they are not evaluated in Table 7. VCG and PNEL are end-to-end systems while the others are modular

systems. VCG and PNEL were installed locally on a machine with the following configuration: 256 GB RAM, 42 core E5-2650 Intel Xeon v4@2.2 GHz. No GPU was present on the system during run time. For VCG and PNEL, the times taken for first runs were recorded, where the corresponding databases such as Virtuoso and Elasticsearch, were started just before the evaluation. This was done so that the times were not affected by caching from previous runs. For systems except PNEL and VCG, the times mentioned in the table were collected from API calls to their hosted services. It must be considered that, due to network latency, and other unknown setup related configurations at the service end, the times may not be comparably directly. PNEL performs faster than VCG since it avoids querying the KG during runtime, and instead relies on pre-computed KG embeddings. PNEL also uses lesser number of features than VCG. A visible trend is that the more accurate system is slower, however Spotlight is an exception, which performs well in both speed and accuracy.

6 Error Analysis

A prominent feature of PNEL is high precision and low recall. We focus on loss in recall in this section. For LC-QuAD 2.0 test set consisting of 6046 questions, the precision, recall and F-score are 0.688, 0.516 and 0.589 respectively. We categorise the phases of loss in recall in two sections 1) Failure in the candidate generation phase 2) Failure in re-ranking/disambiguation phase. When considering the top 50 search candidates during text label search, it was found that 75.3% of the correct entities were recovered from the entity label index. This meant that before re-ranking we had already lost 24.7% recall accuracy. During re-ranking phase, a further 23.7% in absolute accuracy was lost, leading to our recall of 0.516. We drill down into the 23.7% absolute loss in accuracy during re-ranking, attempting to find the reasons for such loss, since this would expose the weaknesses of the model. In the plots below, we consider all those questions which contained the right candidate entity in the candidate generation phase. Hence, we discard those questions for our analysis, which already failed in the candidate generation phase.

Table 8. Comparison of PNEL’s performance with respect to number of entities in a question.

Entity count	Questions count	Precision	Recall	F1
1	3311	0.687	0.636	0.656
2	1981	0.774	0.498	0.602
3	88	0.666	0.431	0.518

It is observed in Table 8 that recall falls as the number of entities per question rises. It must not be concluded however, that PNEL fails to recognise more

than an entity per question. There were 375 questions with multiple entities where PNEL was able to link all the entities correctly. In Fig. 5 we observe that the recall does not exhibit significant co-relation with either the length of the question, or the length of entity label. The recall remains stable. There seems to be some co-relation between the amount of data available for a given length of question, and the recall on it. It appears that the model performs better on question lengths it has seen more often during training.

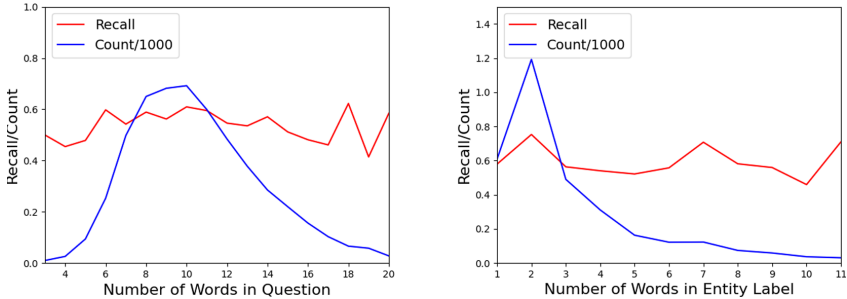


Fig. 5. Plots of recall variation versus 1) Length of Question 2) Length of entity span 3) Frequency of questions with the given lengths in the dataset (scaled down by a factor of 1000)

7 Discussion and Future Work

In this work we have proposed PNEL, an end-to-end Entity Linking system based on the Pointer Network model. We make no modifications to the original Pointer Network model, but identify its utility for the problem statement of EL, and successfully model the problem so the Pointer Network is able to find the right set of entities. We evaluate our approach on three datasets of varying complexity and report state of the art results on two of them. On the third dataset, WebQSP, we perform best in precision but lag behind in recall. We select such features that require no real time KG queries during inference. This demonstrates that the Pointer Network model, and the choice of features presented in this work, result in a practical and deployable EL solution for the largest Knowledge Graph publicly available - Wikidata.

For future work: PNEL being based on the LSTM cell inevitably processes tokens sequentially increasing the response times. This limitation could be overcome by using some variant of the Transformer model [28] instead, which is not only a powerful model but also able to process tokens in parallel. As a future work we would also like to explore different entity embedding techniques and investigate which characteristics make an embedding suitable for the entity linking task.

Acknowledgement. We would like to thank Prof. Dr. Chris Biemann of the Language Technology Group, University of Hamburg, for his valuable suggestions towards improving this work.

References

1. Berant, J., Chou, A., Frostig, R., Liang, P.: Freebase from question-answer pairs. In: Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics (2013)
2. Bollacker, K., Evans, C., Paritosh, P., Sturge, T., Taylor, J.: Freebase: a collaboratively created graph database for structuring human knowledge. In: Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data. ACM (2008)
3. Bordes, A., Usunier, N., Chopra, S., Weston, J.: Large-scale Simple Question Answering with Memory Networks (2015)
4. Delpeuch, A.: OpenTapioca: lightweight entity linking for Wikidata. arXiv preprint [arXiv:1904.09131](https://arxiv.org/abs/1904.09131) (2019)
5. Dubey, M., Banerjee, D., Abdelkawi, A., Lehmann, J.: LC-QuAD 2.0: a large dataset for complex question answering over Wikidata and DBpedia. In: Ghidini, C., et al. (eds.) ISWC 2019. LNCS, vol. 11779, pp. 69–78. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-30796-7_5
6. Dubey, M., Banerjee, D., Chaudhuri, D., Lehmann, J.: EARL: joint entity and relation linking for question answering over knowledge graphs. In: Vrandečić, D., et al. (eds.) ISWC 2018. LNCS, vol. 11136, pp. 108–126. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-00671-6_7
7. Ferragina, P., Scaiella, U.: TAGME: on-the-fly annotation of short text fragments (by Wikipedia entities). In: Proceedings of the 19th ACM International Conference on Information and Knowledge Management (2010)
8. Frey, J., Hofer, M., Obraczka, D., Lehmann, J., Hellmann, S.: DBpedia FlexiFusion the best of Wikipedia > Wikidata > your data. In: Ghidini, C., et al. (eds.) ISWC 2019. LNCS, vol. 11779, pp. 96–112. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-30796-7_7
9. Graves, A., Fernández, S., Schmidhuber, J.: Bidirectional LSTM networks for improved phoneme classification and recognition. In: Duch, W., Kacprzyk, J., Oja, E., Zadrozny, S. (eds.) ICANN 2005, Part II. LNCS, vol. 3697, pp. 799–804. Springer, Heidelberg (2005). https://doi.org/10.1007/11550907_126
10. Kingma, D.P., Ba, J.: Adam: a method for stochastic optimization. arXiv preprint [arXiv:1412.6980](https://arxiv.org/abs/1412.6980) (2014)
11. Kolitsas, N., Ganea, O.-E., Hofmann, T.: End-to-end neural entity linking. In: Proceedings of the 22nd Conference on Computational Natural Language Learning (2018)
12. Lehmann, J., et al.: DBpedia – A Large-Scale, Multilingual Knowledge Base Extracted from Wikipedia. Semantic Web (2015)
13. Lerer, A., et al.: PyTorch-BigGraph: a large-scale graph embedding system. In: Proceedings of the 2nd SysML Conference (2019)
14. Lin, X., Li, H., Xin, H., Li, Z., Chen, L.: KBPearl: a knowledge base population system supported by joint entity and relation linking. Proc. VLDB Endow. **13**, 1035–1049 (2020)

15. Luo, G., Huang, X., Lin, C.-Y., Nie, Z.: Joint entity recognition and disambiguation. In: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics (2015)
16. Mendes, P.N., Jakob, M., García-Silva, A., Bizer, C.: DBpedia spotlight: shedding light on the web of documents. In: Proceedings the 7th International Conference on Semantic Systems (2011)
17. Nguyen, D.B., Abujabal, A., Tran, K., Theobald, M., Weikum, G.: Query-driven on-the-fly knowledge base construction. *Proc. VLDB Endow.* **11**, 66–79 (2017)
18. Nguyen, D.B., Theobald, M., Weikum, G.: NERD: joint named entity recognition and disambiguation with rich linguistic features. *Trans. Assoc. Comput. Ling.* **4**, 215–229 (2016)
19. Pellissier Tanon, T., Vrandečić, D., Schaffert, S., Steiner, T., Pintscher, L.: From freebase to Wikidata: the great migration. In: Proceedings of the 25th International Conference on World Wide Web. International World Wide Web Conferences Steering Committee (2016)
20. Perozzi, B., Al-Rfou, R., Skiena, S.: DeepWalk. In: Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD 2014 (2014)
21. Sakor, A., et al.: Old is gold: linguistic driven approach for entity and relation linking of short text. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). Association for Computational Linguistics (2019)
22. Sakor, A., Singh, K., Patel, A., Vidal, M.-E.: FALCON 2.0: an entity and relation linking tool over Wikidata (2019)
23. Santorini, B.: Part-of-speech tagging guidelines for the Penn Treebank Project (1990)
24. Sevgili, Ö., Panchenko, A., Biemann, C.: Improving neural entity disambiguation with graph embeddings. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop. Association for Computational Linguistics (2019)
25. Sil, A., Yates, A.: Re-ranking for joint named-entity recognition and linking. In: Proceedings of the 22nd ACM International Conference on Information and Knowledge Management. Association for Computing Machinery (2013)
26. Sorokin, D., Gurevych, I.: Mixing context granularities for improved entity linking on question answering data across entity categories. In: Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics. Association for Computational Linguistics (2018)
27. Suchanek, F.M., Kasneci, G., Weikum, G.: Yago: a core of semantic knowledge. In: Proceedings of the 16th International Conference on World Wide Web. ACM (2007)
28. Vaswani, A., et al.: Attention is all you need. In: Guyon, I., et al. (eds.) *Advances in Neural Information Processing Systems*, vol. 30, pp. 5998–6008. Curran Associates Inc. (2017)
29. Vinyals, O., Fortunato, M., Jaitly, N.: Pointer networks. In: Cortes, C., Lawrence, N.D., Lee, D.D., Sugiyama, M., Garnett, R. (eds.) *Advances in Neural Information Processing Systems*, vol. 28. Curran Associates Inc. (2015)
30. Vrandečić, D., Krötzsch, M.: Wikidata: a free collaborative knowledge base. *Commun. ACM* **57**, 78–85 (2014)

31. Yang, Y., Chang, M.-W.: S-MART: novel tree-based structured learning algorithms applied to tweet entity linking. In: ACL 2015 (2015)
32. Yih, W.-T., Richardson, M., Meek, C., Chang, M.-W., Suh, J.: The value of semantic parse labeling for knowledge base question answering. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers). Association for Computational Linguistics (2016)