# Extreme Classification for Answer Type Prediction in Question Answering

Vinay Setty

Department of Electrical Engineering and Computer Science
University of Stavanger
Stavanger, Norway
vsetty@acm.org

## ABSTRACT

Semantic answer type prediction (SMART) is known to be a useful step towards designing effective question answering (QA) systems. The SMART task involves predicting the top-$k$ knowledge graph (KG) types for a given natural language question. This is challenging due to the large number of types in KGs. In this paper, we propose use of extreme multi-label classification using Transformer models (XBERT) by clustering KG types using structural and semantic features based on question text. We specifically improve the clustering stage of the XBERT pipeline using the features derived from KGs. We show that these features can improve end-to-end performance for the SMART task, and yield state-of-the-art results.

## CCS CONCEPTS

• **Information systems** → **Information retrieval**;

## KEYWORDS

Question answering; Answer type prediction

## 1 INTRODUCTION

Question answering (QA) systems provide concise answers to natural language questions. Predicting the semantic answer type of question is an important component of QA systems [6, 10, 17, 25]. Answer type prediction in QA is in fact not a recent problem. Early works have focused on identifying coarse-grained types, e.g., wh-types of questions (who, when, what and where) [3, 9, 11, 12, 23, 31]. Mapping a given question to a semantic type, from type systems of large knowledge graphs (KGs), are known to benefit both open domain QA [28] and factoid questions in KGQA [19, 27, 32].

Both fine and coarse-grained type prediction are complementary to each other. Recognizing this, the *semantic answer type prediction* (SMART) challenge [15] was organized as part of the International Semantic Web Conference in 2020,[1] introducing the following task: given a natural language question, predict both the high-level answer category and a list of fine-grained types from an underlying KG (DBpedia or Wikidata). Coarse-grained categories are boolean, literal, or resource; fine-grained types only apply for the resource category.

The most effective approaches to the SMART challenge are based on Transformer models [18, 26]. Vanilla Transformer models work well out-of-the-box for coarse-grained category prediction, where there are only a handful of possible classes. However, they cannot be used for classifying the large number of semantic resource types found in KGs; for example, DBpedia has over 760 types and

---

[1]https://smart-task.github.io

Wikidata has over 50k types [5]. To overcome this, we propose a BERT-based extreme multi-label classification technique (XBERT) has been proposed that solves this problem using a three-stage pipeline: (1) **Clustering** reduces the number of target classes, (2) **Classifying** trains the Transformer model to predict type clusters to which the question belongs to, and (3) **Ranking** selects the top types within the cluster using a linear ranker [26].

In this three stage pipeline, we show that just by leveraging the *clustering* step, we can gain significant performance improvement in answer type classification. Since clustering is the first step in the three-stage pipeline, any errors occurring during this stage will propagate downstream and thus hinder the entire pipeline's performance. Instead of naively clustering the KG types based on associated question text as in [26], we aim to leverage additional information from the KG, like type descriptions and entity-type assignments. We show that using such additional KG features can significantly improve performance, without modifying other stages in the pipeline.

The type systems of KGs can vary a lot in terms of scale and depth of hierarchy [8], which represents another challenge. We therefore propose a KG-agnostic type clustering technique and perform experiments on two of the most popular KGs, DBpedia and Wikidata, with very different type systems.

Our study is driven by the following research questions:

- **RQ1**: Can we improve the clustering stage of the XBERT pipeline for the SMART task using signals from the KG?
- **RQ2**: How well do these findings generalize across KGs, which differ in the characteristics of their type systems?

The main technical contribution of this paper is the use of signals from a KG to improve the BERT-based extreme multi-label classification approach (XBERT) for the SMART task. Our experiments show that these features can improve results for both DBpedia and Wikidata and yield state-of-the-art performance over vanilla BERT.

## 2 RELATED WORK

The correct prediction of the expected answer type is shown to be one of the most important factors to a QA system's overall performance [17]. Types may be coarse grained entity classes [19, 27, 32] or fine-grained semantic types from the type systems of large knowledge graphs (like Wikidata, DBpedia and Freebase) [21, 28].

In [28], they use probabilistic models to verify if the candidate answer types match the expected answer types to the question. Answer type prediction is also related to the task of inferring semantic types of queries, referred to as *target entity type identification* [1], which has been studied in the context of entity-oriented
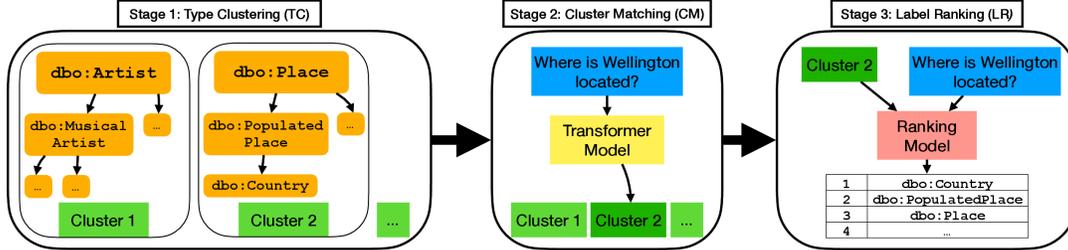
**Figure 1: Three phases of semantic answer type classification.**

search [2, 7, 29]. There, it is approached as a ranking task, using different ways of aggregating entity descriptions [2], which can be combined with additional taxonomic and embedding-based features [7]. Unlike in ad hoc retrieval, the evaluation measure considers the hierarchical relationships between types [2]—the same methodology has been followed in the SMART challenge.

Most participants at the SMART challenge employed classification methods. Common themes include data augmentation [20, 21], the use of word embeddings to represent the questions and types [4, 30], while the top performing approaches employed BERT-based classifiers [18, 26]. Closest to our approach is [18] in second place, performing two stage (coarse and fine-grained) type classification. For fine-grained type prediction, they train a BERT-based classifier using the most frequent types in the dataset. In their follow-up paper [19], they show how to use answer type prediction to improve a KGQA system. However, since this approach can predict only top-level types, it does not scale to type systems with a large number of types with a relatively flat hierarchy such as Wikidata. Moreover, the performance of this method on DBpedia is also not comparable to the XBERT approach. We extend and improve the top performing solution in [26], who introduce the use of XBERT for the SMART challenge.

Several extreme classification approaches have been proposed in the literature [13, 22]. We use XBERT, which provides a good trade-off between performance and efficiency [5]. A more efficient version of XBERT is also available [33].

## 3 PROBLEM STATEMENT

The SMART task is defined as follows: given a natural language input question $q$, (1) predict the coarse category $l \in \{\texttt{boolean}, \texttt{number}, \texttt{string}, \texttt{date}, \texttt{resource}\}$, and (2) if the category is $\texttt{resource}$, then return a list of top-$k$ most relevant types $t \in T$ from the type system $T$ of an underlying knowledge graph (KG).

The first sub-task, *coarse category classification*, may be regarded as a solved problem, as vanilla BERT models can perform it with over 97% accuracy [26]. The second sub-task, *type prediction*, proves to be more challenging due to the large number of types in KGs, ranging from several hundreds (e.g., DBpedia) to tens of thousands (e.g., Wikidata). The type prediction can be viewed as learning a multi-label classifier, which assigns a score to a given question and type $(q, t)$ pair.

The main challenge lies within the resource category, like *Who are some players of the England national football team?*, where the task is not only about identifying matching types, but also ranking them in order of relevance, from most specific to more generic: ["dbo:SoccerPlayer", "dbo:Athlete", "dbo:Person", "dbo:Agent"].

### 3.1 Solution Overview

To overcome the limitation of vanilla transformers for large number of classes, we use the XBERT approach by [5], which uses a three-stage pipeline for the type prediction task (see Fig. 1).

**Type Clustering (*TC*):** First, all unique types $T' \subset T$ which are in the training data are clustered using type vectors. In [26], types are represented using TF-IDF vectors, computed from the question text of all the instances for a given type in the training data. In this paper, we use structural and textual features derived from the KGs to represent types. This stage is the main focus of this paper and it is elaborated further in Sect. 4.

**Cluster Matching (*CM*):** Next, we fine-tune a pre-trained Transformer model to match a given question $q$ into one of the clusters $c$ from the previous stage. The output of this stage is a cluster matching score, $m(q, c)$, which is computed for each cluster based on the model's confidence.

**Label Ranking (*LR*):** Finally, a one-vs-all linear classifier is trained for each type within a cluster $c$ and matched to the given input question $q$ to predict a relevance score $h(q, t)$ for each type in that cluster. The final relevance prediction combines the scores from the *CM* and *LR* stages as follows: $f(q, t) = \sigma(m(q, c), h(q, t))$, where $t \in c$ and $\sigma$ is a non-linear activation function that learns the weights to combine the cluster matching score ($m(q, c)$) and the relevance score within the cluster ($h(q, t)$). The top-$k$ types are then chosen based on $f(q, t)$.

## 4 TYPE CLUSTERING

This section discusses the type clustering (*TC*) step. The main purpose of the TC phase is to reduce the label space for the BERT model used in the second stage (*CM*) of the XBERT. More formally, the *TC* step groups the set of types $T'$ from the training data into clusters $C$, where, $|C| << |T|$. In the past, question text-based TF-IDF vectors were used for clustering the types [26]. This method does not use any external features for the type prediction. In this paper, we investigate how features derived from a KG, such as type similarity vectors and embeddings encoded using the type description text and KG structure, could improve the type prediction step and thereby end-to-end performance on the SMART task.

### 4.1 Type Representation

Embedding KG types allows us to use any clustering algorithm within the XBERT approach. More formally, each type is represented with a vector $\{z_t : t \in T'\}$, where $z_t \in \mathbb{R}^d$ is a d-dimensional vector such that two KG types with high semantic similarity are closer together in the embedding space. Within the XBERT approach, these vectors are obtained by aggregating the features from

the question text of the instances for a give type from the training data. For example, for the type "dbo:SoccerPlayer" features are aggregated from the related questions from the training data such as *What are some players of England national football team?, Which soccer players were born on Malta?* etc. In [26], TF-IDF weights are used to represent the types. In this paper, we instead represent each type individually using textual and structural features from the type taxonomy of the underlying KG.

*4.1.1 KG Structure-Based Representations.* Based on the rationale that each type is defined by its neighborhood in the KG (for example, entities with that type and other related types) we propose the following two type representation methods:

**Type similarity vector:** Each type is represented by a vector of pairwise similarities to other types, based on the set of entities they share. We compute the similarity between two types $t$ and $t'$ using the Jaccard similarity, denoted as $J(t, t')$. Each type $t$ then is represented as a vector of pairwise similarities between the type and all other types $t_1 \ldots t_n$ in the KG:

$$z_t^{jaccard} = [J(t, t_1), J(t, t_2), \cdots, J(t, t_n)] \in \mathbb{R}^{|T|}.$$

**RDF2Vec embeddings:** To leverage the neighborhood structure of KGs, we use network embedding representations. Specifically, we use RDF2Vec [24] since it is trained on both DBPedia and Wikidata, and it provides embeddings for KG types in addition to entities. In short, RDF2Vec uses random walks to construct the sequences of nodes in the KGs. These sequences are treated like sentences and distributed representation of the nodes in the RDF graphs are learned, similar to how Word2Vec models [16] are trained. We use the pre-trained 200-dimensional RDF2Vec SkipGram embeddings for both DBpedia and Wikidata.[2]

*4.1.2 BERT-based Representations.* Next, we use textual descriptions of the types from the KG to represent them. These are sparse and short in practice, therefore, we augment them with the descriptions of entities with that type. A similar strategy has been used in retrieval methods for entity-bearing queries [2]. Based on this we propose the following:

**BERT type embeddings:** Each type is represented by concatenating the type description and the descriptions of the entities that are of that type. This text is then encoded using the BERT model to obtain an embedding for the type. More formally, a type description is a sequence of words $t_w$ and each entity description is a sequence of words $e_w^i$, that are concatenated with a special $[SEP]$ token and encoded as below:

$$\begin{aligned} t_w &= [t_{w_1}, t_{w_2}, \ldots, t_{w_n}] \\ e_w^i &= [e_{w_1}^i, e_{w_2}^i, \ldots, e_{w_m}^i] \\ Desc_t &= [t_w[SEP]e_w^1[SEP]e_w^2[SEP]\ldots[SEP]e_w^l] \\ z_t^{BERT} &= BERT(Desc_t) \in \mathbb{R}^{d_{BERT}}, \end{aligned} \quad (1)$$

where $d_{BERT}$ is the BERT embedding dimension (1024).

**Fine-tuned BERT type embeddings:** The BERT model used above is pre-trained using task independent objectives such as masked language modeling and next sentence prediction. Alternatively, we can fine-tune the BERT model specifically for the cluster

matching stage to encode the type embeddings. Intuitively, this would provide improved embeddings for the types.

## 4.2 Clustering Algorithm

Once we have the type representations, any clustering algorithm can be used to group the types into $|C|$ clusters. In this paper, we specifically report the results using simple a K-Means clustering algorithm. We note that other clustering methods were also tried (including KD-Trees and Spherical K-Means clustering), but the choice of clustering algorithm did not have a significant impact on type clustering performance.

## 5 EXPERIMENTS

In this section, we evaluate the different type representations for the *TC* stage and end-to-end SMART task.

## 5.1 Experimental Setup

*5.1.1 Datasets.* The DBpedia and Wikidata datasets each have around 17-18k training instances and 4k test instances. In both datasets, the majority of questions belong to the Resource category and have one or more ground truth types from the respective KG. For more details, see [15].

*5.1.2 Evaluation Metrics.* We follow the evaluation method used in the official SMART challenge. Type prediction is cast as a ranking task and is evaluated using rank-based metrics. It, however, considers only those questions that fall into the literal (number, string, date) or resource answer categories. Furthermore, evaluation is performed differently for DBpedia and for Wikidata, given the nature of their respective type taxonomies. Types in the DBpedia Ontology are organized hierarchically, up to 7 levels deep. There, a graded evaluation metric, Normalized Discounted Cumulative Gain (NDCG@k), is used. Specifically, literal answer types are either correct or incorrect, while resource answer types receive gain values depending on their distance from the gold types in the type hierarchy [2]. In case of Wikidata, the type hierarchy is rather flat. Therefore, type prediction is evaluated using a binary notion of relevance, with Mean Reciprocal Rank (MRR) as the metric.

*5.1.3 Implementation and baseline.* We compute the KG type features using the DBpedia dump from October 2016 and Wikidata dump from May 19, 2021. We use the RoBERTa [14] (roberta-large[3]) for both encoding the textual descriptions and for fine-tuning the classifier in the second stage (*CM*). We ran all experiments on a Linux server with Nvidia Tesla V100 GPU with 32 GB memory. Note that the evaluation script provided by the SMART task organizers had a known bug[4], which we fixed for this paper. Therefore, the scores we report for the baseline method in Table 1 slightly are higher than those reported in [26].

We use the Question text (TF-IDF) method [26] as baseline in the experiments. We represent our KG derived features as (1) **KG-TypeSim**: pairwise KG-TypeSim vector, (2) **RDF2Vec**: RDF2Vec embeddings, (3) **BERT-TypeDesc**: Type description (TD) encoded with BERT and (4) **BERT-TypeDesc (fine-tuned)**: TD encoded with fine-tuned BERT.

---

**Table 1: Results for the type prediction and end-to-end SMART task. Best scores in each metric are boldfaced. † indicates statistically significant result with $p < 0.05$ and Bonferroni correction when compared with the baseline solution from [26].**

| Method | DBpedia | | | | | | Wikidata | |
| | Type prediction | | | End-to-end | | | Type prediction | End-to-end |
| | NDCG@3 | NDCG@5 | NDCG@10 | NDCG@3 | NDCG@5 | NDCG@10 | MRR | MRR |
|---|---|---|---|---|---|---|---|---|
| Question text (TF-IDF) | 0.717 | 0.693 | 0.650 | 0.824 | 0.811 | 0.787 | 0.66 | 0.76 |
| KG-TypeSim | 0.725 | 0.697 | 0.662 | 0.828 | 0.813 | 0.793 | 0.67 | 0.77 |
| KG-RDF2Vec | 0.729 | 0.701 | 0.656 | 0.831 | 0.815 | 0.791 | 0.67 | 0.78 |
| BERT-TypeDesc | 0.727 | 0.706 | 0.665 | 0.830 | 0.818 | 0.795 | 0.67 | 0.78 |
| BERT-TypeDesc (fine-tuned) | **0.734**† | **0.712**† | **0.678**† | **0.834** | **0.822** | **0.802** | **0.68** | **0.79** |

*5.1.4 Parameter settings.* The number of clusters ($|C|$) is a hyperparameter of the *TC* stage, which we set experimentally, using cross-validation on the training set. We consider the values {32, 64, 128, 256, 512}. We observe the best performance for DBpedia with 64 clusters and for Wikidata with 128 clusters, which will be the values used in the remaining of the experiments. We also note that this parameter has hardly any effect in case of Wikidata.

## 5.2 Results

Table 1 shows the results for type classification (i.e., only for the `resource` category) as well as on the end-to-end task (which also includes coarse category classification).

**RQ1** Can we improve the clustering stage for the SMART task by using additional signals from the KG?

We observe that the TF-IDF-based question text baseline performs reasonably well. The KG features "KG-TypeSim" and "KG-RDF2Vec" yield slightly higher NDCG scores. The best performance is obtained with the BERT models that encode textual descriptions of types as well as of corresponding entities from the KG. We also find that fine-tuning can further enhance performance.

We also observe a similar pattern regarding end-to-end performance improvements. However, the results are not statistically significant—this is because the end-to-end pipeline also includes coarse category classification. Specifically, end-to-end evaluation also includes the questions which belong to the `boolean` and various `literal` categories, and hence downplay the overall impact of fine-grained type classification.

**RQ2** How well do these findings generalize across KGs, which differ in the characteristics of their type systems?

When we compare the performance of DBpedia and Wikidata, we see notable (and in some cases statistically significant) improvements for DBpedia, while Wikidata shows only marginal differences. This is also apparent from the fact that the results are statistically significant for DBpedia but not for the Wikidata. We attribute this to the fact the DBpedia has a rich type system, which can by leveraged for type prediction, while Wikidata has a relatively flat type hierarchy.

## 5.3 Analysis

In this section, we analyze anecdotal examples to highlight the strengths and weaknesses of our approach.

For types with good descriptions from KG, **BERT-TypeDesc** methods perform well. For example, "Which are exclusions of Asperger syndrome?", the ground-truth is 'dbo:Disease' and our method predicts ['dbo:Biomolecule', 'dbo:Gene', 'dbo:Disease'] as top-3 types. While, the baseline method predicts ['dbo:Media', 'dbo:EthnicGroup', 'dbo:Religious'], since the question text alone is insufficient. Our methods fail when the type descriptions and questions have little in common. For example, for "Name a participant of the American Revolutionary War." our method predicts ['dbo:Agent', 'dbo:Person', 'dbo:AmericanFootballTeam'] while the ground-truth is ['dbo:Country', 'dbo:State'].

In some cases method is able to predict more appropriate types, while the ground-truth is incorrect. For example, "Which is the vessel class of the galleon" has the ground-truth 'dbo:Work' and our method predicts ['dbo:MeanOfTransportation', 'dbo:Ship', 'dbo:Spacecraft'].

## 6 CONCLUSION AND FUTURE WORK

In this paper, we have proposed different ways to represent KG types to improve the clustering stage of the XBERT pipeline for the SMART task. We have shown that features derived from KGs to represent types are beneficial, especially for KGs with a rich type system, such as DBpedia. Our error analysis suggests that the BERT-based type embeddings are effective in capturing question context, and our method is able to deal with noisy training data.

In addition to improving type prediction, cluster-based features could also be used for improving the third stage of the XBERT pipeline for SMART in the future. It would also be worth exploring if KGs with a richer structure, like DBpedia, can be used to improve prediction effectiveness on KGs with relatively flat type systems, like Wikidata. Moreover, improved type clustering can potentially be useful for other tasks as well that involve KGs. Specifically, in future work, we would like to explore how to use answer type prediction in a conversational setting. For instance, this method can be used to decide if a question can be answered with a KG.

## 7 ACKNOWLEDGEMENTS

# REFERENCES

[1] K. Balog. *Entity-Oriented Search*, volume 39 of *The Information Retrieval Series*. Springer, 2018.

[2] K. Balog and R. Neumayer. Hierarchical target type identification for entity-oriented queries. In *Proceedings of the 21st ACM international conference on Information and knowledge management*, CIKM '12, pages 2391–2394, 2012.

[3] H. Bast and E. Haussmann. More accurate question answering on freebase. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, CIKM '15, pages 1431–1440, 2015.

[4] E. Bill and E. Jiménez-Ruiz. Question embeddings for semantic answer type prediction. In *Proceedings of the SeMantic AnsweR Type prediction task (SMART) at ISWC 2020*, ISWC SMART '20, pages 71–80, 2020.

[5] W.-C. Chang, H.-F. Yu, K. Zhong, Y. Yang, and I. S. Dhillon. Taming pretrained transformers for extreme multi-label text classification. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '20, pages 3163–3171, 2020.

[6] D. A. Ferrucci, E. W. Brown, J. Chu-Carroll, J. Fan, D. Gondek, A. Kalyanpur, A. Lally, J. W. Murdock, E. Nyberg, J. M. Prager, N. Schlaefer, and C. A. Welty. Building Watson: An overview of the DeepQA project. *AI Magazine*, 31(3):59–79, 2010.

[7] D. Garigliotti, F. Hasibi, and K. Balog. Target type identification for entity-bearing queries. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '17, pages 845–848, 2017.

[8] D. Garigliotti, F. Hasibi, and K. Balog. Identifying and exploiting target entity type information for ad hoc entity retrieval. *Information Retrieval Journal*, 22(3):285–323, 2019.

[9] Z. Huang, M. Thint, and Z. Qin. Question classification using head words and their hypernyms. In *Proceedings of the 2008 Conference on empirical methods in natural language processing*, EMNLP '08, pages 927–936, 2008.

[10] S. Kamath, B. Grau, and Y. Ma. Predicting and integrating expected answer types into a simple recurrent neural network model for answer sentence selection. *Computación y Sistemas*, 23(3):665–673, 2019.

[11] C. C. Kwok, O. Etzioni, and D. S. Weld. Scaling question answering to the web. In *Proceedings of the 10th international conference on World Wide Web*, WWW '01, pages 150–161, 2001.

[12] X. Li and D. Roth. Learning question classifiers: the role of semantic information. *Natural Language Engineering*, 12(3):229–249, 2006.

[13] W. Liu, H. Wang, X. Shen, and I. Tsang. The emerging trends of multi-label learning. 2020.

[14] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. Roberta: A robustly optimized bert pretraining approach. 2019.

[15] N. Mihindukulasooriya, M. Dubey, A. Gliozzo, J. Lehmann, A. N. Ngomo, and R. Usbeck, editors. *Proceedings of the SeMantic AnsweR Type prediction task (SMART) at ISWC 2020 Semantic Web Challenge co-located with the 19th International Semantic Web Conference (ISWC 2020), Virtual Conference, November 5th, 2020*, volume 2774 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2020.

[16] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26:3111–3119, 2013.

[17] D. Moldovan, M. Pașca, S. Harabagiu, and M. Surdeanu. Performance issues and error analysis in an open-domain question answering system. *ACM Trans. Inf. Syst.*, 21(2):133–154, apr 2003.

[18] C. Nikas, P. Fafalios, and Y. Tzitzikas. Two-stage semantic answer type prediction for question answering using BERT and class-specificity rewarding. In *Proceedings of the SeMantic AnsweR Type prediction task (SMART) at ISWC 2020*, SMART '20, pages 19–28, 2020.

[19] C. Nikas, P. Fafalios, and Y. Tzitzikas. Open domain question answering over knowledge graphs using keyword search, answer type prediction, sparql and pre-trained neural models. In *International Semantic Web Conference*, pages 235–251, 2021.

[20] A. Perevalov and A. Both. Augmentation-based answer type classification of the smart dataset. In *Proceedings of the SeMantic AnsweR Type prediction task (SMART) at ISWC 2020*, ISWC SMART '20, pages 1–9, 2020.

[21] A. Perevalov and A. Both. Improving answer type classification quality through combined question answering datasets. In *International Conference on Knowledge Science, Engineering and Management*, ISWC '21, pages 191–204, 2021.

[22] Y. Prabhu, A. Kag, S. Harsola, R. Agrawal, and M. Varma. Parabel: Partitioned label trees for extreme classification with application to dynamic search advertising. In *Proceedings of the 2018 World Wide Web Conference*, WWW '18, pages 993–1002, 2018.

[23] E. Riloff and M. Thelen. A rule-based question answering system for reading comprehension tests. In *ANLP-NAACL 2000 Workshop: Reading Comprehension Tests as Evaluation for Computer-Based Language Understanding Systems*, ANLP-NAACL '00, pages 13–19, 2000.

[24] P. Ristoski and H. Paulheim. Rdf2vec: RDF graph embeddings for data mining. In *The Semantic Web - ISWC 2016 - 15th International Semantic Web Conference, Kobe, Japan, October 17-21, 2016, Proceedings, Part I*, volume 9981 of *ISWC '16*, pages 498–514, 2016.

[25] R. S. Roy and A. Anand. *Question Answering for the Curated Web: Tasks and Methods in QA over Knowledge Bases and Text Collections*. Morgan & Claypool Publishers, 2021.

[26] V. Setty and K. Balog. Semantic answer type prediction using BERT IAI at the ISWC SMART task 2020. In *Proceedings of the SeMantic AnsweR Type prediction task (SMART) at ISWC 2020*, ISWC SMART '20, pages 10–18, 2020.

[27] T. Shen, X. Geng, Q. Tao, D. Guo, D. Tang, N. Duan, G. Long, and D. Jiang. Multi-task learning for conversational question answering over a large-scale knowledge base. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, EMNLP '19, pages 2442–2451, 2019.

[28] H. Sun, H. Ma, W.-t. Yih, C.-T. Tsai, J. Liu, and M.-W. Chang. Open domain question answering via semantic enrichment. In *Proceedings of the 24th International Conference on World Wide Web*, WWW '15, pages 1045–1055, 2015.

[29] D. Vallet and H. Zaragoza. Inferring the most important types of a query: A semantic approach. In *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '08, pages 857–858, 2008.

[30] S. Vallurupalli, J. Sleeman, and T. Finin. Fine and ultra-fine entity type embeddings for question answering. In *Proceedings of the SeMantic AnsweR Type prediction task (SMART) at ISWC 2020*, ISWC SMART '20, pages 57–70, 2020.

[31] E. M. Voorhees. The trec question answering track. *Nat. Lang. Eng.*, 7(4):361–378, Dec. 2001.

[32] S. Yavuz, I. Gür, Y. Su, M. Srivatsa, and X. Yan. Improving semantic parsing via answer type inference. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, EMNLP '16, pages 149–159, 2016.

[33] H.-F. Yu, K. Zhong, and I. S. Dhillon. Pecos: Prediction for enormous and correlated output spaces. 2020.