

Description Boosting for Zero-Shot Entity and Relation Classification

Gabriele Picco

IBM Research Europe

gabriele.picco@ibm.com

Leopold Fuchs

IBM Research Europe

leopold.fuchs@ibm.com

Marcos Martinez Galindo

IBM Research Europe

marcos.martinez.galindo@ibm.com

Alberto Purpura

IBM Research Europe

alp@ibm.com

Vanessa Lopez

IBM Research Europe

vanlopez@ie.ibm.com

Hoang Thanh Lam

IBM Research Europe

t.l.hoang@ie.ibm.com

Abstract

Zero-shot entity and relation classification models leverage available external information of unseen classes - e.g., textual descriptions - to annotate input text data. Thanks to the minimum data requirement, Zero-Shot Learning (ZSL) methods have high value in practice, especially in applications where labeled data is scarce. Even though recent research in ZSL has demonstrated significant results, our analysis reveals that those methods are sensitive to provided textual descriptions of entities (or relations). Even a minor modification of descriptions can lead to a change in the decision boundary between entity classes. In this paper we formally define the problem of identifying effective descriptions for zero shot inference, we propose a strategy for generating variations of an initial description, a heuristic for ranking them and an ensemble method capable of boosting the predictions of zero-shot models through description enhancement. Empirical results on four different entity and relation classification datasets show that our proposed method outperform existing approaches and achieve new SOTA results on these datasets under the ZSL settings. The source code of the proposed solutions and the evaluation framework are open-sourced.¹

1 Introduction

Named Entity Recognition (NER) and Relation Extraction (RE) allow for the extraction and categorization of structured data from unstructured text, which in turn enables not only more accurate entity recognition and relationship extraction, but also getting data from several unstructured sources, helping to build knowledge graphs and the semantic web. However, these methods usually rely on labeled data (usually human-annotated data) for a good performance, usually requiring domain experts for data acquisition and labeling, which may

incur in high costs. Thus, it is not surprisingly that there is often a lack of labeled data for new domains, limiting the performance of these methods.

Zero-shot learning (ZSL) is a classification task in machine learning where - at inference time - samples are classified into one of several classes which were not observed during training. In ZSL, the sets of training and test entity classes are disjoint. Therefore, the strategy employed by zero-shot models is to rely on prior general knowledge that could be transferred to unseen instances at inference time. Having a classifier that can generalize to new unseen classes is important for a variety of practical reasons. First, ZSL methods can be used to learn models that are more robust to labeled data shortages and distributional shifts. Moreover, they can be used to extend the reach of models to new domains.

ZSL approaches in the Natural Language Processing (NLP) domain have seen significant improvements in recent years thanks to the availability of large pre-trained Language Models (LMs). For example, it has been shown that models such as GPT-3 (Brown et al., 2020), OPT (Zhang et al., 2022) and FLAN (DBL) achieve strong performances on many NLP tasks, including translation, question-answering, and cloze tests without any gradient updates or fine-tuning.

For entity recognition - including classification and linking - and relation classification problems, recent ZSL methods (Aly et al., 2021; Ledell Wu, 2020; Chen and Li, 2021) rely on textual descriptions of entities or relations. Descriptions provide the required information about the semantics of entities (or relations), which help the models to identify entity mentions in texts without observing them during training. Works such as (Ledell Wu, 2020; De Cao et al., 2021) and (Aly et al., 2021) show how effective it is to use textual descriptions to perform entity recognition tasks in the zero-shot context. However, the quality of the descriptions has

¹<https://github.com/IBM/zshot>

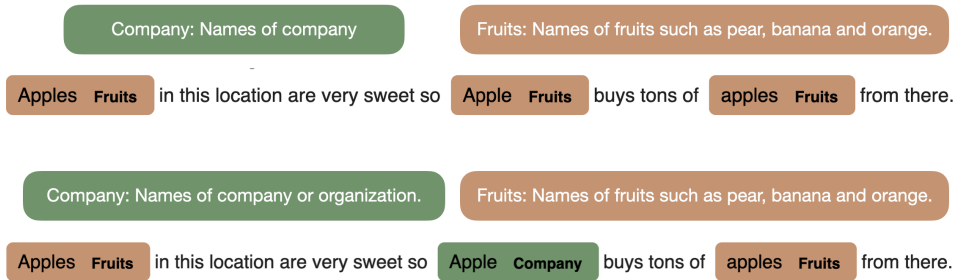


Figure 1: A small modification of the *Company* class description results in different entity predictions.

an impact on how effective the transfer of knowledge from observed to unseen entities (Aly et al., 2021). The same mechanism can also be applied in other contexts, such as relation classification (Chen and Li, 2021). From now on, we refer to entity as both named entities and named relations.

An example of named entity classification with ZSL is demonstrated in Figure 1. At inference time, a zero-shot model is given short textual descriptions of new entity classes such as *Company* or *Fruits*, it then identifies and annotates mentions of those entity classes in an input sentence. Although state-of-the-art ZSL methods such as SMXM (Aly et al., 2021) have demonstrated significant results in recent research works, this toy example shows how the quality of the provided descriptions influences the accuracy of these models. For example, in Figure 1 even with a small modification of the *Company* entity class description, the SMXM model changes its entity prediction. In practice, the sensitivity to entity descriptions is problematic because, for non-expert users, it is not a trivial task to choose a proper description for black-box zero-shot models, in particular in an unfamiliar domain.

In this paper, we study different methods for improving the descriptions in an unsupervised way. Specifically, we propose UDEBO (for Unsupervised Description Boosting), the first unsupervised method capable of automatically modifying/generating descriptions to improve entity predictions in the zero-shot settings. We present several strategies to alter descriptions, such as using a generative model, paraphrasing, and summarization combined with description ranking/ensemble methods to reduce model uncertainty and increase overall performance. We empirically evaluate the performance of UDEBO on 4 existing standard zero-shot datasets, spanning two tasks: (i) named entity classification and (ii) relation classification.

Our results show that for the zero-shot entity

classification tasks, UDEBO improved the results of state-of-the-art models by 7 and 1.3 percentage points in terms of Macro F1 Score in the OntoNotes and MedMentions datasets, respectively. For what concerns relation classification, we achieve a performance improvement of 6 and 3 percentage points (Macro F1 Score) on the FewRel and WikiZS datasets over our baseline models, respectively.

We organize the paper as follows. In Section 2 we formally define the problem we aim to solve in this paper, i.e. how to enhance entity or relation descriptions to improve the performance of zero-shot models. In Section 3 we describe the proposed approaches for description boosting while in Section 4 we describe our experimental setup and results. We provide a literature review, a discussion about Large Language Models (LLMs) and draw the conclusions of our work in Sections 5, 6 and 7, respectively.

2 Problem Definition

Given a set of entity classes E of interest with their textual descriptions D and a corpus of sentences S to annotate as input, we can define the problem of description enhancement as follows:

Problem 1 (Description enhancement) Denote $\phi(D, S)$ as the function estimating the accuracy of ZSL models when using a given entity description D for annotating an input text corpus S . Our goal is to generate a set of descriptions D^* such that:

$$D^* = \arg \max_D \phi(D, S) \quad (1)$$

In Section 3.1 we describe different strategies to generate new entity descriptions D' for the input set E , intending to improve the accuracy of the predictions by the ZSL models over that corpus. If the labeled data is known, it is possible to select the

best descriptions via a brute force search across different description reformulations by measuring the accuracy as a function of D and S . However, given the absence of labeled data in the zero-shot context, an unsupervised approach is needed for ranking the descriptions D that yield the highest classification accuracy. In Section 3.2 and Section 3.3, we will discuss methods for ranking or combining predictions from different description variations to achieve better results.

3 Methods

The UDEBO approach comprises 2 steps. First, the descriptions are generated or improved. Finally, the descriptions are ranked in order to select the best ones. As an optional step, we analyze the ensembling of descriptions for boosting performance.

3.1 Generating description variations

Improving the completeness or clarity of entity descriptions is a complicated problem without a formal definition of an objective function, as there is a large space of candidates to explore. To enhance entity descriptions, in a more controlled way, we propose the following strategies.

Extension with pre-trained LMs. We propose to use large pre-trained LMs for generating text using the given description as context. Large LMs, as shown in (Petroni et al., 2019), capture linguistic and relational knowledge that can be extracted through generation to extend a given description. In Section 4 we analyse the use of GPT-2 (Radford et al., 2019) for generating descriptions variations.

Extension with a fine-tuned LM. We fine-tune a LM for description generation and expansion. The LM is fine-tuned on a large dataset containing about 5.3 million Wikidata instances, including the name and the first few sentences of the respective articles. The model is fine-tuned on extending a truncated sub-string of the textual description, using a sequence to sequence objective. In Section 4 we analyse the use of a T5 large (Raffel et al., 2020) fine-tuned model for generating descriptions variations.

Summarization. Text summarization can be used to generate a concise description with less noise compared to the original one. In the experimental results we analyse the effect of using a BERT2BERT (small) (Turc et al., 2019) model fine-

tuned on CNN/Dailymail for text summarization to enhance entities descriptions.

Paraphrasing. Paraphrasing a description can simplify its linguistic form, using more common and general terms. In the experimental results we analyse the effect of using a Pegasus (Zhang et al., 2019) model fine-tuned for paraphrasing.

3.2 Description ranking via entropy

To rank a description for an entity, we propose to use a zero-shot model to first compute the probabilities of classes for each mention in the input text with a candidate description. We then compute the information entropy H from this input. In information theory, entropy is the average level of "information" or "uncertainty" inherent to a variable's possible outcomes. Our assumption is that the lower the entropy is, the higher the confidence of the prediction will be, so Problem 1 can be reformulated as:

$$D^* = \arg \min_D H(D, S) \quad (2)$$

Where H is the entropy of a zero-shot model for a corpus S , using the description D to accomplish a certain classification task. This way we can rank different candidate descriptions and choose the best one without requiring any labeled data, which is ideal for the zero-shot setting.

3.3 Boosting performances with descriptions variations ensembling

Besides description ranking via entropy, we propose an ensemble method that combines predictions from multiple pipelines executed with different entity descriptions. The main idea behind this approach is to leverage the complementary information provided by the different definitions to make a more accurate prediction, reducing the variance and bias of an individual pipeline. Furthermore, using the methods described in section 3.1, the descriptions variations can provide additional information useful for correctly discriminating between unseen classes.

Entity description ensemble. Given a sentence, for each span s and an entity label $e \in E$, denote $v(s, e)$ as the number of pipelines that predict s or a sub-sequence of s with entity label e . For instance, given a span $s = \textit{London Bridge}$, assume that among ten pipelines, four pipelines predict the label of s as $e_1 = \textit{Facility}$, the other four pipelines

Dataset	Split	Instances	Entities / Relations
MedMentionsZS	train	26770	11
	val	1289	5
	test	1048	5
OntoNotesZS	train	41475	4
	val	1358	4
	test	426	3
Fewrel	train	44800	64
	test	11200	16
WikiZS	train	70952	83
	val	12982	15
	test	9494	15

Table 1: Number of sentences and entities for each split of the considered datasets.

predict the label of *London* as $e_2 = \text{Location}$ and the rest of the pipelines predict *Bridge* as Facility. Therefore, the accumulated number of votes for the span *London Bridge* are $v(s, e_1) = 6$ and $v(s, e_2) = 4$. Considering the majority of the votes, the final predicted label for the span *London Bridge* is Facility. Once the span *London Bridge* has been assigned a label, all of its sub-spans become redundant and thus are removed from consideration.

4 Experiments and Results

This section discusses experimental settings, baseline methods, and empirical results for both entity and relation classification tasks.

4.1 Datasets and experimental settings

We use two different settings: one for the Entity Classification (EC) task and one for the Relation Classification (RC) one.

Entity Classification setting. We use the pre-trained SMXM model (Aly et al., 2021) with the checkpoints available in the official GitHub repository.² We refer the reader to the original paper (Aly et al., 2021) to see the details of the implementation, the training parameters, and the datasets used for fine-tuning the model. There are two different checkpoints, one for each one of the datasets used, OntoNotes (Pradhan et al., 2013) and MedMentions (Mohan and Li, 2019). Both datasets have been processed as in the respective official GitHub repositories. Table 1 shows the number of rows and the entities of each dataset. Note that the number of rows reported in Table 1 refers to the zero-shot version of the dataset, containing only sentences with entities. See Appendix A for more

²<https://github.com/Raldir/Zero-shot-NERC/>

information on this process and the datasets. The results reported are all based on the *test* split of the datasets.

Relation Classification setting. For RC, we use ZS-BERT³ (Chen and Li, 2021), a multitask learning model, based on BERT, to directly predict unseen relations. We trained our checkpoint using the official implementation of the model and following the steps of the official repository.³ The datasets we use are FewRel (Han et al., 2018) and WikiZS (Sorokin and Gurevych, 2017). The results reported are all based on the *test* split of the datasets.

Description alteration settings. The language models used for the description alteration strategies: summarization, paraphrasing and pre-trained were obtained from the checkpoints available on Huggingface, while for the latter strategy we have fine-tuned a pre-trained T5-large model. We report detailed hyper-parameters of description alteration methods in section B of the appendix.

4.2 Empirical results

This section discusses the results of entity classification using methods for description enhancement.

4.2.1 Entity classification

Table 2 shows the results of the ensemble method (UDEBO) with ten descriptions generated by each of the description enhancing strategies, including pre-trained, finetuning, summarization and paraphrasing. For each enhancing strategy, we report the results when the descriptions with the lowest entropy are chosen for each class. The *Combined* strategy shows the results with the lowest entropy among all description-enhancing strategies.

We can see that the ensemble method (UDEBO) outperforms the SMXM baseline using the original descriptions provided on the OntoNotesZS dataset with a significant margin of 7 percentage points in terms of Macro F1 Score. On the MedMentionZS dataset, the improvement is 1.3 percentage points on the same reference performance measure (Macro F1 Score). Description ranking based on entropy works well with the pre-trained strategy on OntoNotesZS. However, the entropy does not seem to be a reliable score of model uncertainty on the MedMentionsZS dataset. Finding an alternative uncertainty score to entropy could be considered

³<https://github.com/dinobby/ZS-BERT>

Datasets	Methods	Precision	Recall	Micro F1	Macro F1	Accuracy
OntoNotesZS	SMXM	20.96	48.15	30.76	29.12	86.36
	SMXM (Pre-trained)	24.05	51.40	32.77	32.78	87.69
	SMXM (Finetuned)	17.97	42.21	25.21	23.90	85.76
	SMXM (Summarization)	18.93	35.45	24.68	19.47	85.93
	SMXM (Paraphrased)	18.49	40.90	25.46	23.41	85.14
	SMXM (Combined)	18.86	42.58	26.15	23.74	84.83
	UDEBO	31.14	46.51	36.78	36.15	88.29
MedMentionsZS	SMXM	16.79	40.55	20.38	21.70	83.05
	SMXM (Pre-trained)	13.25	37.98	19.64	18.26	81.88
	SMXM (Finetuned)	13.67	36.05	19.82	19.13	83.18
	SMXM (Summarization)	10.96	26.68	15.37	17.92	83.02
	SMXM (Paraphrased)	14.77	26.51	18.97	19.41	86.74
	SMXM (Combined)	12.80	37.15	19.04	17.92	81.63
	UDEBO	19.51	32.73	23.86	22.97	85.70

Table 2: UDEBO, i.e. the ensemble of predictions with description variations, compared to the SMXM baseline.

Datasets	Methods	Precision	Recall	Micro F1	Macro F1	Accuracy
Fewrel	ZS-BERT	25.08	21.59	21.59	17.89	21.59
	ZS-BERT (Pre-trained)	18.25	25.29	25.29	19.10	25.29
	ZS-BERT (Finetuned)	19.39	16.09	16.09	14.59	16.09
	ZS-BERT (Summarization)	19.83	19.81	19.81	15.21	19.81
	ZS-BERT (Paraphrased)	25.89	21.76	21.76	19.90	21.76
	ZS-BERT (Combined)	17.09	16.53	16.53	16.53	16.53
	UDEBO	28.38	25.68	25.68	22.12	25.68
WikiZS	ZS-BERT	34.18	33.90	37.14	30.97	37.14
	ZS-BERT (Pre-trained)	14.73	15.80	14.29	11.72	14.29
	ZS-BERT (Finetuned)	16.23	16.26	16.62	13.65	16.62
	ZS-BERT (Summarization)	19.07	19.57	19.62	16.87	19.62
	ZS-BERT (Paraphrased)	25.50	27.60	27.60	24.56	27.60
	ZS-BERT (Combined)	17.34	19.62	18.43	16.27	18.43
	UDEBO	34.79	37.11	40.17	34.25	40.17

Table 3: UDEBO, i.e. the ensemble of predictions with description variations, compared to the ZS-BERT baseline.

as future work. Overall, these results confirm our hypothesis - discussed in Section 1 - that zero-shot methods are sensitive to provided descriptions and that an ensemble of description enhancement methods is needed to obtain more robust results.

4.2.2 Relation classification

In Table 3, we report our evaluation of the proposed approaches on the RC task. The results we observe here are similar to what we described for entity classification where the proposed ensembling method (UDEBO) achieves a higher performance across different measures compared to the baseline ZS-BERT model that does not rely on any relation description reformulation approach. We also observe on the FewRel dataset a higher Macro F1

Score associated with most of the description enhancement variants when employed independently from each other. These results further validate the strength of the proposed approach to enhance relation descriptions employed by ZSL models to improve their performance.

4.2.3 Descriptions enhancement strategies comparison and limitations

Generating variations of descriptions is relatively simple, as described in Subsection 3.1, several strategies allow to generate plausible extensions or variations of a text. Considering the results of ranking the descriptions using entropy in Section 4, we analyze and discuss here the correlation between Macro F1 Score and entropy measures and

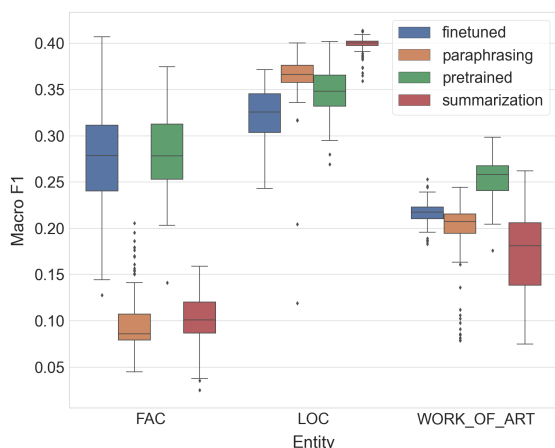


Figure 2: The figure shows the distributions of Macro F1 Score values on the test split of the OntoNotesZS dataset for each class, using the strategies described in Section 3.1 to generate 100 description variations for each class.

the limitations of the proposed approach.

Figure 2 and Figure 3 show the distributions of the Macro F1 Score on the test split of the OntoNotesZS and the MedmentionsZS dataset for each class, using the strategies described in Section 3.1 to generate 100 description variations for each class. None of the strategies is a clear champion over all the classes. The high variance of the performance explains the fact that the ensemble method makes a better prediction as observed in Table 2 and Table 3 thanks to successfully combining the strength of individual description alteration strategies. Figure 4 shows the correlations between Macro F1 Score and entropy for each unseen class on the OntoNotesZS test split with 100 description variations. Although there appears to be a significant statistical correlation using a sign test with ($p\text{-value} = 0.03$) between Macro F1 Score and entropy measures on the OntoNotesZS test set, the correlation does not appear to be statistically significant in the MedMentionsZS dataset. Also, as evidenced by the results in Table 2 and 3, using the descriptions with minimum entropy does not seem like a good strategy for selecting descriptions.

This phenomenon may be due to several factors like the change in the style of generated descriptions compared to the ones observed during training. Although a new description might seem more relevant, it could make the model more uncertain. See an example in Appendix C.2. The importance of this problem motivates the future study of alternative heuristics with more significant correlations,

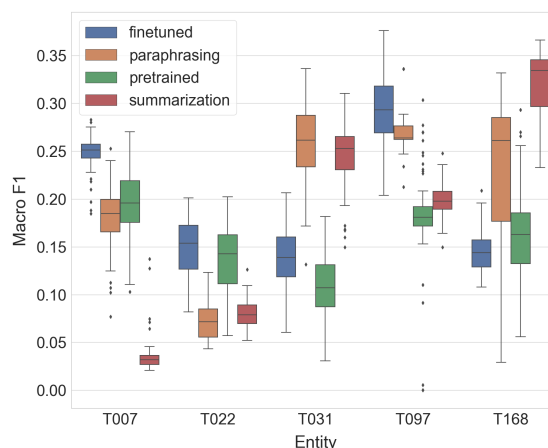


Figure 3: The figure shows the distributions of Macro F1 Score values on the test split of the MedMentionsZS dataset for each class, using the strategies described in Section 3.1 to generate 100 description variations for each class.

indirectly unveiling the mechanism behind zero-shot predictions.

5 Related work

Zero-shot entity recognition and linking. Zero-shot end-to-end entity linking refers to the task of detecting and disambiguating entity mentions by linking them to an entity in a Knowledge Base (KB), without requiring new labeled data. KBs are inherently incomplete and evolve over time with the addition of new entities and relations. Zero-shot entity linking usually relies on available textual information, or other set of relations in the KB, to generalise to entity sets unseen in the training data.

BLINK (Wu et al., 2020) is a BERT-based solution for Zero-shot linking of textual mentions - extracted for example using FLAIR (Akbi et al., 2018) - to entities in Wikipedia. It follows a bi-encoder architecture, each mention is encoded in a dense space, together with its context (left and right part of the input sentence). Independently, each entity in the KB is encoded in the same dense space together with its context e.g., entity description. Mentions are linked to entities in the dense space using a nearest neighbour search. To improve accuracy, candidate entities are ranked by passing each concatenated mention, its context and entity description to a more expensive cross-encoder.

GENRE (De Cao et al., 2021) is a BART based model fine-tuned using a sequence to sequence objective, which claims to outperform BLINK. It is an autoregressive end-to-end entity linker, it detects

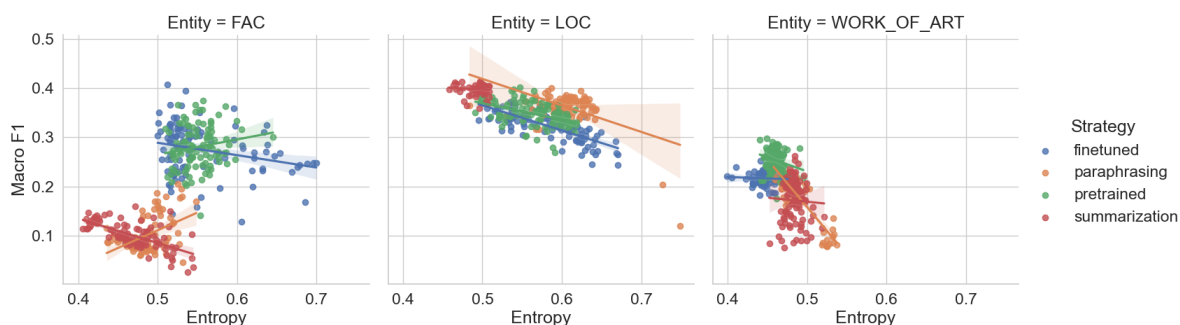


Figure 4: Analysis of the correlation between entropy and Macro F1 Score on unseen classes on the OntoNotesZS test split. Entropy can be calculated without the need for labeled data, therefore, if a correlation exists it can be used as an unsupervised heuristic to select descriptions that improve model performance.

and retrieves mentions and the respective entities in a KB by generating their unique textual name - left to right, token-by-token. To do so, it uses a constrained decoding strategy that forces the generated name to be in a predefined candidate set. Compared to multi-class classification models such as BLINK, GENRE has a lower memory footprint to store dense vectors for large KBs, scaling linearly with vocabulary size, not entity count, and does not need to subsample negative data during training.

Zero-shot entity classification. Entity classification consists in predicting a probability for each semantic type of an entity mention, given a set of types (e.g. organisation, organic compound). The most straightforward feature used to generalise to unseen types is the textual descriptions. For example, SMXM (Aly et al., 2021) uses a cross-attention encoder to generate a vector representation for each type description and token in the input sentence and recognizes as entity types those representations that are closer to each other, including rarer classes unseen in training. It is evaluated using zero-shot adaptations of *OntoNotes* (Pradhan et al., 2013) and the domain specific biomedical dataset *MedMentions* (Mohan and Li, 2019), it also considers *out-of-KB* predictions i.e., *nil* predictions for mentions that do not have a valid gold entity.

ReFinED (Ayoola et al., 2022) is an end-to-end entity linking model optimised to perform mention detection, fine-grained entity typing (classification), and entity disambiguation in a single pass. Similar to BLINK, ReFinED uses a bi-encoder architecture modified to encode all mentions in a document simultaneously, which improves efficiency relatively to zero-shot models such as (Wu et al., 2020) that requires a forward-pass for each mention. Men-

tion embeddings and entity description embeddings are projected into a shared vector space to calculate their dot product as the entity score. A fast bi-encoder combined with a score for unseen entities, computed based on the scores for entity types and description, is enough for ReFinED to obtain state-of-the-art performance on entity linking and to scale the approach from Wikipedia (5.9M entities) to Wikidata (90M entities).

The analyses in (Aly et al., 2021) show that while Wikipedia descriptions work well on general entity types, they perform poorly on domain specific data, e.g. *MedMentions*. They also show the impact of using annotation guidelines for descriptions to improve the transfer of knowledge from observed to unseen entities. The adoption of this approach led to a better performance compared to using a class name itself or Wikipedia passages. In particular, description vagueness, noise and negations had a negative effect, while annotation guidelines, including explicit examples and syntactic and morphological cues, improved the performance.

Zero-shot relation classification. Textual descriptions have also been employed in the relation classification task to predict new relations that could not be observed at training time. For example, ZS-BERT (Chen and Li, 2021) learns two functions – one to project sentences and the other to project relation descriptions into an embedding space. The objective is first to jointly minimise the distance between the embedding vectors for an input sentence and the relation description for positive entity pairs and then to classify the relation (using a softmax layer to produce a classification probability). At inference time, the prediction of unseen relation classes can be achieved through

nearest neighbor search. Overall, using descriptions seems to improve existent zero-shot methods and expand their domains of application. Still, descriptions are not always good enough to get good predictions. Improving the accuracy of these approaches remains an open challenge. The better the separation between embedding of different relations, the more accurate the model predictions, however, as the number of unseen relations increases, it becomes more difficult to predict the right one (Chen and Li, 2021).

Existent ZSL methods usually rely on external knowledge from KGs, ranging from textual information, class attributes, hierarchy, domain and range constrains and relations to logic rules. There are relatively few studies evaluating their performance for unseen relations, a comparison using different external knowledge settings for zero-shot relation classification and KG completion can be seen in (Geng et al., 2021). To the best of our knowledge, we present the first approach to automatically predict and generate entity descriptions to improve the accuracy of entity recognition and relation classification models.

6 Discussion about Large Language Models

In the era of Large Language Models, all kinds of problems are being solved with LLMs, that achieve outstanding results in different tasks. However, LLMs also raise some concerns, being one of the most important, the green footprint of these models. Serving a single 175 billion LLM requires at least 350 GB GPU memory using specialized infrastructure, (Zheng et al., 2022). This makes it unfeasible for a lot of users to use LLMs, and even if it's possible to use, there is a lot of concern about using them, specially for tasks that could be solved with smaller models. With UDEBO we try to push the research in a direction that improves the performance of small LMs to achieve results comparable to LLMs. However, to compare the performance of UDEBO against LLMs, we select 3 open-source LLMs available to the community and evaluate them. The results, discussion, and settings can be found in Appendix D.

7 Conclusion and future work

In this paper, we formally defined the problem of selecting descriptions to make predictions about unseen classes in the ZSL context. To the best of our

knowledge, this is the first time for entity/relation ZSL problems in which the impact of description variations on prediction performance is studied, and different methods for automatic creation of descriptions are considered. We empirically evaluated the sensitivity of two ZSL methods to description changes, and proposed 4 different strategies to enhance them using the implicit knowledge of pre-trained language models. We also studied in detail the efficacy of the proposed entropy-based heuristic to rank different description formulations, analyzing its correlation with the performance (in terms of Macro F1 Score) of the model. We observed a negative correlation between the proposed heuristic and Macro F1 Score on two out of four of the considered datasets (OntoNotesZS and FewRel). The same assumption however was not valid for the other datasets (MedMentionsZS and WikiZS), thus motivating the need to develop more effective heuristics in the future. Finally, we described the UDEBO method, which combines the predictions obtained by the same model using different automatically generated variants of entity and relation descriptions. Our experimental results, on 4 different datasets, spanning across two different NLP tasks (Entity Classification and Relation Classification) showed how UDEBO outperforms the baselines by a significant margin and achieves new state-of-the-art results on these benchmarks under the zero-shot setting. Existing ensemble methods focus on ensembling different models trained on different data or models with different structures. Our work is orthogonal to these approaches, we proposed methods that consider entity/relation description variations as the hyper-parameters that need to vary. Most importantly, the description variations are not provided by the users but were generated from the initially provided descriptions. Therefore, this is a new way of creating ensembles, at least in the context of Zero-shot Entity/Relation extraction this is the first time a method for descriptions generation to diversify the pipelines and make an ensemble that improve the quality of the results is proposed.

Limitations

While our proposed method of boosting for zero-shot entity and relation classification shows promising results, there are several limitations that need to be acknowledged. Firstly, ensembling methods can be computationally intensive, which can limit their applicability to large-scale datasets. Our current implementation combining multiple model predictions, by varying the descriptions, which requires a significant amount of computational resources. Therefore, future research should explore alternative methods for ensembling that are more computationally efficient. Secondly, while we used entropy as a metric to identify helpful descriptions, it may not always be the most effective metric. Entropy measures the uncertainty or randomness of a distribution, but it may not necessarily capture the semantic relevance of a description. Therefore, there is a need for further research to develop better metrics or heuristics for identifying helpful descriptions. Finally, we have only experimented with a limited set of generation techniques. Future research should explore ways to improve the quality and coverage of the descriptions.

Supplemental Material Statement and Reproducibility

Source code availability the source code including the frameworks for pipeline ensemble and evaluation on the standard benchmark datasets are available as an open-source github repository: <https://github.com/IBM/zshot>.

Dataset availability Both datasets used in our evaluation are publicly available, also included in our open-source evaluation framework. We also release the enhanced descriptions generated by the generative models used to create the ensemble pipelines.

References

Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. Contextual string embeddings for sequence labeling. In *COLING 2018, 27th International Conference on Computational Linguistics*, pages 1638–1649.

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, et al. 2023. Falcon-40b: an open large language model with state-of-the-art performance.

Rami Aly, Andreas Vlachos, and Ryan McDonald. 2021. Leveraging type descriptions for zero-shot named entity recognition and classification. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1516–1528, Online. Association for Computational Linguistics.

Tom Ayoola, Shubhi Tyagi, Joseph Fisher, Christos Christodoulopoulos, and Andrea Pierleoni. 2022. ReFinED: An efficient zero-shot-capable approach to end-to-end entity linking. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Track*, pages 209–220, Hybrid: Seattle, Washington + Online. Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Chih-Yao Chen and Cheng-Te Li. 2021. ZS-BERT: towards zero-shot relation extraction with attribute representation learning. *CoRR*, abs/2104.04697.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways.

Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. 2021. Autoregressive entity retrieval.

- In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *CoRR*, abs/2305.14314.
- Ning Ding, Shengding Hu, Weilin Zhao, Yulin Chen, Zhiyuan Liu, Hai-Tao Zheng, and Maosong Sun. 2021. Openprompt: An open-source framework for prompt-learning. *arXiv preprint arXiv:2111.01998*.
- Yuxia Geng, Jiaoyan Chen, Xiang Zhuang, Zhuo Chen, Jeff Z. Pan, Juan Li, Zonggang Yuan, and Huajun Chen. 2021. Benchmarking knowledge-driven zero-shot learning. *Journal of Web Semantics*.
- Xu Han, Hao Zhu, Pengfei Yu, Ziyun Wang, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2018. [FewRel: A large-scale supervised few-shot relation classification dataset with state-of-the-art evaluation](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4803–4809, Brussels, Belgium. Association for Computational Linguistics.
- Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alexander Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. [Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes](#).
- Martin Josifoski Sebastian Riedel Luke Zettlemoyer Ledell Wu, Fabio Petroni. 2020. Zero-shot entity linking with dense entity retrieval. In *EMNLP*.
- Sunil Mohan and Donghui Li. 2019. [Medmentions: A large biomedical corpus annotated with {umls} concepts](#). In *Automated Knowledge Base Construction (AKBC)*.
- Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. 2019. Language models as knowledge bases? *arXiv preprint arXiv:1909.01066*.
- Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Hwee Tou Ng, Anders Björkelund, Olga Uryupina, Yuchen Zhang, and Zhi Zhong. 2013. [Towards robust linguistic analysis using OntoNotes](#). In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*, pages 143–152, Sofia, Bulgaria. Association for Computational Linguistics.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(140):1–67.
- Daniil Sorokin and Iryna Gurevych. 2017. Context-aware representations for knowledge base relation extraction. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1784–1789.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and finetuned chat models](#).
- Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Well-read students learn better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962v2*.
- Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, and Guoyin Wang. 2023. [Gpt-ner: Named entity recognition via large language models](#).
- BigScience Workshop, :, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Lucic, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, Dragomir Radev, Eduardo González Ponferrada, Efrat Levkovich, Ethan Kim, Eyal Bar Natan, Francesco De Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady Elsahar, Hamza Benyamina, Hieu Tran, Ian Yu, Idris Abdulmumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jörg Froberg,

- Joseph Tobing, Joydeep Bhattacharjee, Khalid Al-mubarak, Kimbo Chen, Kyle Lo, Leandro Von Werra, Leon Weber, Long Phan, Loubna Ben allal, Ludovic Tanguy, Manan Dey, Manuel Romero Muñoz, Maraim Masoud, María Grandury, Mario Šaško, Max Huang, Maximin Coavoux, Mayank Singh, Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad A. Jauhar, Mustafa Ghaleb, Nishant Subramani, Nora Kassner, Nurulaqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, Paulo Villegas, Peter Henderson, Pierre Colombo, Priscilla Amuok, Quentin Lhoest, Rheza Harliman, Rishi Bommasani, Roberto Luis López, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, Shayne Longpre, Somaieh Nikpoor, Stanislav Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan Thrush, Valentin Danchev, Vassilina Nikoulina, Veronika Laippala, Violette Lepercq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heinzerling, Chenglei Si, Davut Emre Taşar, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, Andrea Santilli, Antoine Chaffin, Arnaud Stiegler, Debajyoti Datta, Eliza Szczechla, Gunjan Chhablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, Lintang Sutawika, M Saiful Bari, Maged S. Al-shaibani, Matteo Manica, Nihal Nayak, Ryan Teehan, Samuel Albanie, Sheng Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Evry, Trishala Neeraj, Urmish Thakker, Vikas Raunak, Xiangru Tang, Zhengxin Yong, Zhiqing Sun, Shaked Brody, Yallow Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, Jaesung Tae, Jason Phang, Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared Casper, Jeff Rasley, Max Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myriam Peyrounette, Nicolas Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre Cornette, Pierre François Lavallée, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Requena, Suraj Patil, Tim Dettmers, Ahmed Baruwa, Amanpreet Singh, Anastasia Cheveleva, Anne-Laure Ligozat, Arjun Subramonian, Aurélie Névéol, Charles Lovering, Dan Garrette, Deepak Tunuguntla, Ehud Reiter, Ekaterina Taktasheva, Ekaterina Voloshina, Eli Bogdanov, Genta Indra Winata, Hailey Schoelkopf, Jan-Christoph Kalo, Jekaterina Novikova, Jessica Zosa Forde, Jordan Clive, Jungo Kasai, Ken Kawamura, Liam Hazan, Marine Carpuat, Miruna Clinciu, Najeon Kim, Newton Cheng, Oleg Serikov, Omer Antverg, Oskar van der Wal, Rui Zhang, Ruochen Zhang, Sebastian Gehrmann, Shachar Mirkin, Shani Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun, Tomasz Limisiewicz, Verena Rieser, Vitaly Protasov, Vladislav Mikhailov, Yada Pruksachatkun, Yonatan Belinkov, Zachary Bamberger, Zdeněk Kasner, Alice Rueda, Amanda Pestana, Amir Feizpour, Ammar Khan, Amy Faranak, Ana Santos, Anthony Hevia, Antigona Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh HajiHosseini, Bahareh Behroozi, Benjamin Ajbade, Bharat Saxena, Carlos Muñoz Ferrandis, Danish Contractor, David Lansky, Davis David, Douwe Kiela, Duong A. Nguyen, Edward Tan, Emi Baylor, Ezinwanne Ozoani, Fatima Mirza, Frankline Ononiwu, Habib Rezanejad, Hessie Jones, Indrani Bhat-tacharya, Irene Solaiman, Irina Sedenko, Isar Nejadgholi, Jesse Passmore, Josh Seltzer, Julio Bonis Sanz, Livia Dutra, Mairon Samagaio, Maraim Elbadri, Margot Mieskes, Marissa Gerchick, Martha Akinlolu, Michael McKenna, Mike Qiu, Muhammed Ghauri, Mykola Burynok, Nafis Abrar, Nazneen Rajani, Nour Elkott, Nour Fahmy, Olanrewaju Samuel, Ran An, Rasmus Kromann, Ryan Hao, Samira Alizadeh, Sarmad Shubber, Silas Wang, Sourav Roy, Sylvain Viguier, Thanh Le, Tobi Oyebade, Trieu Le, Yoyo Yang, Zach Nguyen, Abhinav Ramesh Kashyap, Alfredo Palasciano, Alison Callahan, Anima Shukla, Antonio Miranda-Escalada, Ayush Singh, Benjamin Beilharz, Bo Wang, Caio Brito, Chenxi Zhou, Chirag Jain, Chuxin Xu, Clémentine Fourrier, Daniel León Perrián, Daniel Molano, Dian Yu, Enrique Manjavacas, Fabio Barth, Florian Fuhrmann, Gabriel Altay, Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec, Imane Bello, Ishani Dash, Jihyun Kang, John Giorgi, Jonas Golde, Jose David Posada, Karthik Rangasai Sivaraman, Lokesh Bulchandani, Lu Liu, Luisa Shinzato, Madeleine Hahn de Bykhovetz, Maiko Takeuchi, Marc Pàmies, Maria A Castillo, Marianna Nezhurina, Mario Sängler, Matthias Samwald, Michael Cullan, Michael Weinberg, Michiel De Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank, Myungsun Kang, Natasha Seelam, Nathan Dahlberg, Nicholas Michio Broad, Nikolaus Muellner, Pascale Fung, Patrick Haller, Ramya Chandrasekhar, Renata Eisenberg, Robert Martin, Rodrigo Canalli, Rosaline Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda, Shlok S Deshmukh, Shubhanshu Mishra, Sid Kiblawi, Simon Ott, Sinee Sang-aroonisiri, Srishti Kumar, Stefan Schweter, Sushil Bharati, Tanmay Laud, Théo Gigant, Tomoya Kainuma, Wojciech Kusa, Yanis Labrak, Yash Shailesh Bajaj, Yash Venkatraman, Yifan Xu, Yingxin Xu, Yu Xu, Zhe Tan, Zhongli Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and Thomas Wolf. 2023. [Bloom: A 176b-parameter open-access multilingual language model](#).
- Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian Riedel, and Luke Zettlemoyer. 2020. [Scalable zero-shot entity linking with dense entity retrieval](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6397–6407, Online. Association for Computational Linguistics.
- Yongqin Xian, Christoph H. Lampert, Bernt Schiele, and Zeynep Akata. 2019. [Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly](#). *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(9):2251–2265.
- Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. 2019. [Pegasus: Pre-training with extracted gap-sentences for abstractive summarization](#).
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel

Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.

Lianmin Zheng, Zhuohan Li, Hao Zhang, Yonghao Zhuang, Zhifeng Chen, Yanping Huang, Yida Wang, Yuanzhong Xu, Danyang Zhuo, Joseph E. Gonzalez, and Ion Stoica. 2022. *Alpa: Automating inter- and intra-operator parallelism for distributed deep learning*. *CoRR*, abs/2201.12023.

8 Appendix

A Datasets

As mentioned in Section 4.1, we evaluate our approach on four different datasets, two for EC and two for RC. For EC, we use OntoNotes (Pradhan et al., 2013) and MedMentions (Mohan and Li, 2019). OntoNotes is a dataset that comprises various genres of text (news, conversational telephone speech, weblogs, usenet newsgroups, broadcast, talk shows). We use the version available in Huggingface⁴ and adapt it to perform zero-shot as explained in (Aly et al., 2021), removing all the entities that are out of the split - i.e., each split has a unique set of entities, so all the entities labeled with entities out of that set are removed - removing sentences without any entity labelled and using the same train/test/dev splits, so the pre-trained model has not seen the entities in the test set neither. The entity descriptions used for OntoNotesZS (the zero-shot version of OntoNotes) were provided by the authors of (Aly et al., 2021).

MedMentions is a corpus of Biomedical papers annotated with mentions of UMLS entities. We apply the same preprocessing steps we used for the MedMentions dataset, with the descriptions available in the official GitHub repository of (Aly et al., 2021).² The version of the MedMentionsZS dataset we use is also available on Huggingface. Both of them in their zero-shot version, as proposed in (Aly et al., 2021). To convert them to the zero-shot version, we follow the following steps:

1. Get the train/test/dev splits of the datasets;
2. Collect the entities in each split;
3. Remove entities out of the split i.e., if one entity e belongs to the train split, all mentions labelled as e in the test and dev splits will be replaced with the O label.

⁴https://huggingface.co/datasets/conll2012_ontonotesv5

Split	Entity	Count
Train	O	515420
	T103	22360
	T038	25007
	T033	9824
	T062	5445
	T098	3574
	T017	12575
	T074	1165
	T082	7511
	T058	14779
	T170	5996
T204	4922	
Test	O	27433
	T031	212
	T097	360
	T007	448
	T168	321
T022	89	
Validation	O	34400
	T201	404
	T091	196
	T037	434
	T005	224
T092	452	

Table 4: Number of entities labelled in each split in MedMentionsZS.

4. Remove sentences without labels. As the previous processing step (3) may remove all the entities of one sentence, the result dataset will have a lot of empty sentences. These sentences are removed in the final dataset.

Table 4 and Table 5 report the entities for each split in the dataset and the number of entities for MedMentionsZS and OntoNotesZS, respectively. As we can observe, both datasets are highly imbalanced, with some entities appearing 25007 times and some others only 89 in the case of MedMentionsZS, and 24163 and 65 times for OntoNotesZS. However, the most common entities are used only for training and the ones with fewer examples are used for validation and testing. As pointed in (Xian et al., 2019), real-world scenarios annotated data is likely to be available for the more common ones.

In Table 6 we report some statistics concerning the length of sentences on both MedMentionsZS and OntoNotesZS. In both datasets, there are sentences containing only 1 token and 1 entity. The maximum number of tokens also varies across

Split	Entity	Count
Train	O	909142
	ORG	24163
	GPE	21938
	DATE	18791
	PERSON	22035
Test	O	11299
	FAC	149
	LOC	215
	WORK_OF_ART	169
Validation	O	36790
	NORP	1277
	LAW	65
	EVENT	179
	PRODUCT	214

Table 5: Number of entities labelled in each split in OntoNotesZS.

datasets and splits, with a maximum of 179 for MedMentionsZS and 210 for OntoNotesZS.

For RC, we use the FewRel (Han et al., 2018) and WikiZS (Sorokin and Gurevych, 2017) datasets. FewRel is a dataset for RC compiled by collecting entity-relation triplets with sentences from Wikipedia articles, and manually filtered to ensure the data quality and class balance. We use different relations for the train and the test split to ensure the zero-shot version of the dataset. The dataset is available in the Huggingface hub.⁵ We use the *train_wiki* split in Huggingface as training split for the ZS-BERT model and the *wiki_val* as test split. Table 1 shows the total number of sentences in FewRel, and the number of different relations for each split. There are 700 samples for each relation in each split, thus the number of sentences reported in Table 1 is equal to the number of relations times the number of samples for each of them (e.g. train split: $44800 = 64 * 700$). Differently from FewRel, WikiZS was constructed using the Wikidata knowledge base. The dataset contains a total of 93431 sentences, each with an entity pair and a labelled relation between them. In this case, the number of instances per relation class is not balanced and we employ our own random splits containing different distinct sets of relations for the training (83 relations), validation (15 relations) and

⁵https://huggingface.co/datasets/few_rel

testing (15 relations) of the ZS-BERT model. More information on the dataset is contained in Table 1.

B Additional details on the models used for generating description variations

In this section, we report additional details on the methods used to generate description variations described in Section 3.1.

Extension with pre-trained LMs. An off-the-shelf GPT-2 pre-trained model was used for generating the variations, using the checkpoint from the Huggingface Hub.⁶ We used *min_length* = 80, *max_length* = 120, *num_beams* = 8, *temperature* = 1 and *no_repeat_ngram_size* = 2 for the generation.

Extension with a fine-tuned LM. A model based on T5 large (Raffel et al., 2020) and fine-tuned on the task of description generation and extension was used for generating the variations. As a starting point for the fine-tuning, the checkpoint from Huggingface Hub⁷ was used. The Wikidata dataset, containing the name and the first few sentences of included Wikipedia articles where the model was fine-tuned on, was taken from Facebook Research’s BLINK project.⁸ After cleaning the data i.e., removing instances with no or too short (less than 10 words) descriptions, about 5,310,000 samples were available for training the model to perform a new sequence to sequence task using *learning_rate* = $3e - 05$ and *epochs* = 1. The objective was to complete the input description, starting from a sub-string containing the first ten words of it. For the generation task, just the name of the description was used. In the latter case, we set *min_length* = 80, *max_length* = 120, *num_beams* = 8, *temperature* = 1 and *no_repeat_ngram_size* = 2.

Summarization. A warm-started BERT2BERT (small) model fine-tuned on the CNN/Dailymail for document summarization was used for generating the descriptions variations, using the checkpoint from the Huggingface Hub.⁹ We used *min_length* = 80, *max_length* = 512, *num_beams* = 8, *temperature* = 1 and

⁶<https://huggingface.co/gpt2>

⁷<https://huggingface.co/t5-large>

⁸<http://dl.fbaipublicfiles.com/BLINK/entity.jsonl>

⁹https://huggingface.co/mrm8488/bert-small2bert-small-finetuned-cnn-daily_mail-summarization

Dataset	Split	Mean	Max	Min	Mean	Max	Min
		#Tokens	#Tokens	#Tokens	#Entities	#Entities	#Entities
MedMentionsZS	train	26	179	1	6	78	1
	test	28	102	2	2	33	1
	validation	28	119	4	2	12	1
OntoNotesZS	train	25	210	1	3	99	1
	test	29	108	2	3	39	1
	validation	28	186	3	1	27	1

Table 6: Entity classification datasets details.

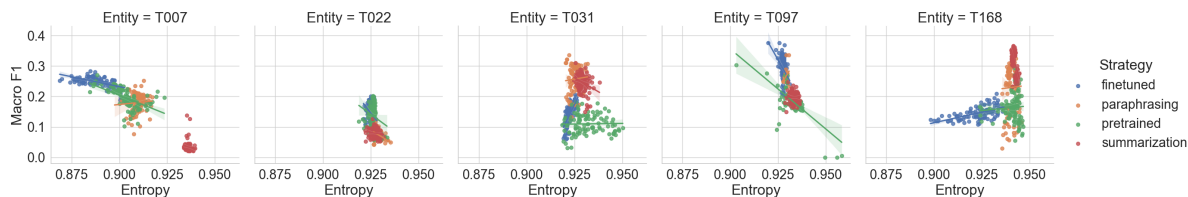


Figure 5: Analysis of the correlation between entropy and Macro F1 Score on unseen classes on the MedmentionsZS test split.

$no_repeat_ngram_size = 2$ for this set of experiments.

Paraphrasing. A PEGASUS model fine-tuned for paraphrasing was used for generating the description variations, using the checkpoint from the Huggingface Hub.¹⁰ We used $min_length = 10$, $max_length = 60$, $num_beams = 8$, $temperature = 1$ and $no_repeat_ngram_size = 2$ for the generation of text.

C Additional experiments on Entropy and Macro F1 Score correlations

In this section we report additional insights on the correlation analysis discussed in the paper.

C.1 Correlations analysis of Macro F1 Score and entropy on MedmentionsZS

Figure 5 reports the correlations between Macro F1 Score and entropy on MedmentionsZS test-set. As discussed in the paper we did not observe any statistically significant correlations, with $p\text{-value} = 0.50$.

C.2 Example of generated descriptions and entropy values

We analyze the variations for both entity and relation descriptions.

¹⁰https://huggingface.co/tuner007/pegasus_paraphrase

Entity Tables 7, 8 and 9 show some examples of the description variations generated with the different methods. The original descriptions used are shown in Table 10. It can be seen that the models used for generation might generate wrong descriptions, e.g. Summarization variations in Table 7, in which the model, that has been fine-tuned over CNN/Dailymail, a news dataset, replicates the style of writing as if it was a reporter. Also, some hallucinations lead the models to generate URLs, that have been removed. This means that the base models used for generating the descriptions are crucial, and they should be chosen and fine-tuned carefully for the description generation task. We leave this as a task for future work.

Relations Given the relation **Film Director** described as:

"director(s) of film, TV-series, stageplay, video game or similar".

The fine-tuned approach for generating variations produces the alternative description:

The director(s) of a film, TV-series, stageplay, video game or similar is the person who directs the production of the film or television series. The term "director" is also used to describe an individual or group of people who are responsible for the creation, production, and/or

Enhancement method	Description	Entropy
Summarization	"in reporting live from right outside hospital," says massachusetts general hospital. buildings are referred to using the name of the company or organization that uses them. names of man - made structures should be marked as fac when they refer to the structure, usually in a locative way. " he's reporting. here : " were looking for people who don't know the exact location of this type "	0.460
Summarization	" in reporting live from right outside hospital," says massachusetts general hospital. buildings are referred to as fac when they refer to the physical structure of the building itself. names of man - made structures include infrastructure (streets, bridges), buildings, monuments, etc. " were still trying to make a difference in our lives,says hospital's director of medical examiners.	0.538
Pretrained	Names of man-made structures: infrastructure (streets, bridges), buildings, monuments, etc. 1. The name of the building or structure is used to distinguish it from other buildings or structures. For example, a building might be called a museum or a church. A building can also be referred to as a "museum" or "church" if it is built on top of an existing structure (e.g., a university building). It is also possible to refer to a structure by its name. Examples of buildings that are known as museums include: a school building;	0.497
Pretrained	Names of man-made structures: infrastructure (streets, bridges), buildings, monuments, etc. The number of buildings in a city is based on the size of the city and its population. The larger the number, the more likely a building is to be built. For example, if you have 10,000 people, you would expect to have a 10% chance of having 10 buildings per square mile. However, this is not always the case. In fact, there are some buildings that are more than 10 times as large as the population of your city.	0.573
Paraphrasing	Buildings that are referred to using the name of a company or organization should be marked as FAC if they refer to the physical structure of the building itself.	0.438
Paraphrasing	The names of the man-made structures should be marked FAC when references are made to the physical structure.	0.558
Finetuned	Names of man-made structures: infrastructure (streets, bridges), buildings, monuments etc. are the names given to a number of different types of structures in the world, e.g. roads, railways and roads. The names of these structures are usually derived from the words "infrastructure" and "building" or "domestic".	0.512
Finetuned	Names of man-made structures: infrastructure (streets, bridges), buildings, monuments etc. are the names given to the physical structures of a country or region. These names are usually derived from the name of the municipality in which they are located, or the place where it is located (e.g., city or town). For example, in the United States, the city of San Francisco is known as "San Francisco International Airport", while in Mexico, Mexico City is called "Mexico City" or "California".	0.699

Table 7: Example descriptions variations for the class FAC of OntoNotes. The worst and best variations, according to the entropy, are selected for each method.

Enhancement method	Description	Entropy
Summarization	also included in this category are named regions such as the middle east, areas, neighborhoods, continents and regions of continents. these include mountain ranges, coasts, borders, planets, geo - coordinates, bodies of water. don't mark deictics or other non - proper nouns, but don't be marked when they are part of the location name itself. the list is based on the names of geographical locations other than gpes.	0.465
Summarization	these include mountain ranges, coasts, borders, planets, geo - coordinates, bodies of water. also included in this category are named regions such as the middle east, areas, neighborhoods, continents and regions of continents. don't mark deictics or other non - proper nouns, but don't mark the names. do you know a hero? nominations are open at cnn. com / heroes.	0.510
Pretrained	Names of geographical locations other than GPEs. These include mountain ranges, rivers, lakes, and streams. For more information, please visit:	0.492
Pretrained	Names of geographical locations other than GPEs. These include mountain ranges, lakes, rivers, streams, and oceans. The following table lists the geographic locations of each of the three types of geographic data used in this report. The geographic coordinates for each type of data are listed in the table below. For more information about the geospatial data, please refer to the Geographic Information System (GIS) Web site at	0.583
Paraphrasing	The geographical locations that aren't included in GPEs include mountain ranges, coasts, borders, planets, as well as the regions of continents and the Middle East.	0.582
Paraphrasing	There are geographical locations beyond GPEs, including mountain ranges, coasts, borders, planets, and bodies of water.	0.656
Finetuned	These are the names of geographical locations other than GPEs. These include mountain ranges, peaks, lakes, rivers, and other geographical features that are not part of a geographic area. Some of these names are also used in other places such as the United States, United Kingdom, Canada, Australia, New Zealand, France, Germany, Italy, Japan, Norway, Sweden, Switzerland and Switzerland.	0.521
Finetuned	These are the names of geographical locations other than GPEs. These include mountain ranges, mountain peaks, and other geographical features, such as rivers, streams, lakes, or other bodies of water. The following is a list of names used by the United States Geological Survey (USGS) to identify geological features that are not part of the U.S. Geographic Names Information System (GNIS).	0.618

Table 8: Example descriptions variations for the class LOC of OntoNotes. The worst and best variations, according to the entropy, are selected for each method.

directing of video games, films, television shows, or other forms of media..

Although the generated description seems more complete and containing relevant additional information, the entropy calculated with ZS-BERT is higher in this case than when using the original description. This means that the model is more uncertain of its prediction.

D Discussion about Large Language Models

Recently, several large language models (LLMs) have been released demonstrating high capabilities for diverse NLP tasks including, but not limited to, text generation, question answering, text summarization and also NER and RE, (Workshop et al., 2023), (Chowdhery et al., 2022). In the zero-shot setting, the LLMs perform exceptionally well. However, there are some problems that might limit the usage of LLMs for NER and RE, like the need of prompt engineering and result parsing, the token

Enhancement method	Description	Entropy
Summarization	the headline of the article being annotated should only be marked if they are referential. in other words, a reference to an article is markable as a work of art. paper headlines should be treated if it's referred to as an art work. but if in the body of a text, then it is marked as art by art, but it should not be used.	0.464
Summarization	paper headlines should only be marked if they are referential. in other words the headline of the article is markable as a work of art. there is a reference to an article being annotated if she is referred to. but in some words it should not be used as an art work or art, it's been used to make art more modern than it was used in the original version of this article.	0.492
Pretrained	Titles of books, songs, television programs and other creations. Also available on iTunes. Bookmark this page:	0.452
Pretrained	Titles of books, songs, television programs and other creations. Also available in English and Spanish. In addition to the titles of titles, there are a number of other titles available for download on the Internet. The following are examples of some of the most popular titles that have been released in the past few years. For more information, please see the full list of available titles on this site. You can also check out the list at the top of this page to see which titles are currently available and which ones are not. If you have any questions about any of these titles or would like	0.477
Paraphrasing	The titles of books, songs, TV programs and other creations should be marked if they are referential.	0.471
Paraphrasing	It is not necessary to include quotations in the article if the headline is referential.	0.531
Finetuned	The titles of books, songs, television programs and other creations. Also known as the title of a book, song, TV program, or television series, are the names of the creators of those works. Often, these titles are also referred to as "titles" or "pronunciations". In the United States, the term 'title' is used to indicate the author's title. The term is also used in other countries, such as Canada, Australia, New Zealand, and the U.S. of Azerbaijan.	0.426
Finetuned	Titles of books, songs, television programs and other creations. Also referred to as "titles", are the titles given to a book, song, TV program, or other work of art. The term is derived from the Latin "titus" ("title"). It can also be used to describe the title of the work itself, the author's name, etc.	0.454

Table 9: Example descriptions variations for the class WORK_OF_ART of OntoNotes. The worst and best variations, according to the entropy, are selected for each method.

Class	Description
FAC	Names of man-made structures: infrastructure (streets, bridges), buildings, monuments, etc. belong to this type. Buildings that are referred to using the name of the company or organization that uses them should be marked as FAC when they refer to the physical structure of the building itself, usually in a locative way: "I'm reporting live from right outside [Massachusetts General Hospital]"
LOC	Names of geographical locations other than GPEs. These include mountain ranges, coasts, borders, planets, geo-coordinates, bodies of water. Also included in this category are named regions such as the Middle East, areas, neighborhoods, continents and regions of continents. Do NOT mark deictics or other non-proper nouns: here, there, everywhere, etc. As with GPEs, directional modifiers such as "southern" are only marked when they are part of the location name itself.
WORK_OF_ART	Titles of books, songs, television programs and other creations. Also includes awards. These are usually surrounded by quotation marks in the article (though the quotations are not included in the annotation). Newspaper headlines should only be marked if they are referential. In other words the headline of the article being annotated should not be marked but if in the body of the text here is a reference to an article, then it is markable as a work of art.

Table 10: Original OntoNotes class descriptions (Source: annotation guidelines, <https://catalog.ldc.upenn.edu/docs/LDC2013T19/OntoNotes-Release-5.0.pdf>)

limitation, the hallucinations and out of context generation or the efficiency.

- Prompt Engineering and result parsing. The performance of these models is highly dependent on the prompt used, and also the output of the model (Ding et al., 2021). Thus, depending on the use case, one prompt might or might not be adequate. Also, the output has to be processed to extract the actual entities, and depending on the prompt this process can be different.
- Token Limitation. These models are all based on text generation, having a minimum and a maximum number of tokens to generate. Depending on these hyperparameters, the result might not be complete or might lead to hallucinations or false positives.
- Hallucinations and out of context generation. LLMs often suffer from hallucination and out of context generation, which in the case of NER and RE might result into entities and relations extracted that are not present in the

text. Some approaches add a self-verification strategy to alleviate the hallucination issue, which requires further executions of the LLM (Wang et al., 2023). Moreover, common NER and RE use cases focus only on some specific entities or relations to be extracted, but these models can extract other entities and relations that may not be of interest to the user.

- Efficiency. Serving a single 175 billion LLM requires at least 350 GB GPU memory using specialized infrastructure, (Zheng et al., 2022). This makes it unfeasible for a lot of users to use LLMs, and even if it's possible to use, there is a lot of concern of the green footprint of these models.

The efficiency problem is one of the most important problems of the LLMs, and thus some approaches try to reduce the size of the models or to train new models via knowledge distillation. Recent works approach the knowledge distillation process with human rationales to improve the performance of the distilled model (Hsieh et al., 2023). In this approach, the human rationale adds information to the

input, so the model can use it to perform the task. Similarly, the descriptions of the entities and relations are used to add information to the input. We leave to future work the usage of the descriptions for knowledge distillation. In either cases, UDEBO could be used to improve the descriptions or the human rationales to improve the performance of the models.

The focus of this work is the evaluation of the method UDEBO, and not the performance of the model itself, as a different size of the model, more pretraining, or even a different architecture could lead to changes in the results. However, we evaluate 3 LLMs, BLOOM (Workshop et al., 2023), FALCON (Almazrouei et al., 2023) and LLAMA 2 (Touvron et al., 2023). In Table 11 we evaluate the largest version of the models in zero-shot. In Table 12 we fine-tune a smaller version of the models (7B) using QLORA (Detmiers et al., 2023), with 0.06%, 0.03%, and 0.12% trained parameters for BLOOM, FALCON, LLAMA 2, respectively.

Model	Size	MedMentions	OntoNotes
BLOOM	176B	0.11	0.14
FALCON	40B	0.12	0.09
LLAMA 2	70B	0.25	0.10

Table 11: F1-Score for LLMs evaluated on MedMentions and OntoNotes.

Model	Size	MedMentions	OntoNotes
BLOOM	7B	0.25	0.00
FALCON	7B	0.20	0.00
LLAMA 2	7B	0.33	0.09

Table 12: F1-Score for QLORA fine-tuned LLMs evaluated on MedMentions and OntoNotes.

The fine-tuned version of the models benefit in the MedMentions dataset, which is specific, but they suffer in the generic domain (OntoNotes), as they extracted entities from the training set, and not the ones in the test set. We use the following prompt (example for OntoNotes):

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

Instruction:

From the input context below extract instances of the following labels: ['LOCATION', 'WORK_OF_ART', 'BUILDING_NAME']

Input: