

# Description Boosting for Zero-Shot Entity and Relation Classification

Anonymous ACL submission

## Abstract

001 Zero-shot entity and relation classification mod- 040  
002 els leverage available external information of 041  
003 unseen classes – e.g., textual descriptions – to 042  
004 annotate input text data. Thanks to the min- 043  
005 imum data requirement, Zero-Shot Learning 044  
006 (ZSL) methods have high value in practice, es- 045  
007 pecially in applications where labeled data is 046  
008 scarce. Even though recent research in ZSL 047  
009 has demonstrated significant results, our anal- 048  
010 ysis reveals that those methods are sensitive 049  
011 to provided textual descriptions of entities (or 050  
012 relations). Even a minor modification of de- 051  
013 scriptions can lead to a change in the decision 052  
014 boundary between entity (or relation) classes. 053  
015 In this paper we formally define the problem of 054  
016 identifying effective descriptions for zero shot 055  
017 inference, we propose a strategy for generating 056  
018 variations of an initial description, a heuristic 057  
019 for ranking them and an ensemble method cap- 058  
020 able of boosting the predictions of zero-shot 059  
021 models through description enhancement. Em- 060  
022 pirical results on four different entity and rela- 061  
023 tion classification datasets show that our pro- 062  
024 posed method outperform existing approaches 063  
025 and achieve new SOTA results on these datasets 064  
026 under the ZSL settings. The source code of the 065  
027 proposed solutions and the evaluation frame- 066  
028 work are open-sourced.<sup>1</sup>

## 029 1 Introduction

030 Zero-shot learning (ZSL) is a classification task 067  
031 in machine learning where – at inference time – 068  
032 samples are classified into one of several classes 069  
033 which were not observed during training. Having a 070  
034 classifier that can generalize to new unseen classes 071  
035 is important for a variety of practical reasons. First, 072  
036 ZSL methods can be used to learn models that are 073  
037 more robust to labeled data shortages and distribu- 074  
038 tional shifts. Moreover, they can be used to extend 075  
039 the reach of models to new domains. 076

ZSL approaches in the Natural Language Pro- 040  
cessing (NLP) domain have seen significant im- 041  
provements in recent years thanks to the availability 042  
of large pre-trained Language Models (LMs). For 043  
example, it has been shown that models such as 044  
GPT-3 (Brown et al., 2020), OPT (Zhang et al., 045  
2022) and FLAN (DBL) achieve strong perfor- 046  
mances on many NLP tasks, including translation, 047  
question-answering, and cloze tests without any 048  
gradient updates or fine-tuning. 049

For entity recognition – including classification 050  
and linking – and relation classification problems, 051  
recent ZSL methods (Aly et al., 2021; Ledell Wu, 052  
2020; Chen and Li, 2021a) rely on textual descrip- 053  
tions of entities or relations. Descriptions provide 054  
the required information about the semantics of en- 055  
tities (or relations), which help the models to iden- 056  
tify entity mentions in texts without observing them 057  
during training. Works such as (Ledell Wu, 2020; 058  
De Cao et al., 2021) and (Aly et al., 2021) show 059  
how effective it is to use textual descriptions to per- 060  
form entity recognition tasks in the zero-shot con- 061  
text. The same mechanism can also be applied in 062  
other contexts such as relation classification (Chen 063  
and Li, 2021b). 064

An example of named entity classification with 065  
ZSL is demonstrated in Figure 1. At inference 066  
time, a zero-shot model is given short textual de- 067  
scriptions of new entity classes such as *Company* or 068  
*Fruits*, it then identifies and annotates mentions of 069  
those entity classes in an input sentence. Although 070  
state-of-the-art ZSL methods such as SMXM (Aly 071  
et al., 2021) have demonstrated significant results 072  
in recent research works, this toy example shows 073  
how the quality of the provided descriptions in- 074  
fluences the accuracy of these models. For exam- 075  
ple, in Figure 1 even with a small modification of 076  
the *Company* entity class description, the SMXM 077  
model changes its entity prediction. In practice, the 078  
sensitivity to entity descriptions is problematic be- 079  
cause, for non-expert users, it is not a trivial task to 080

<sup>1</sup>Anonymized for double-blind review

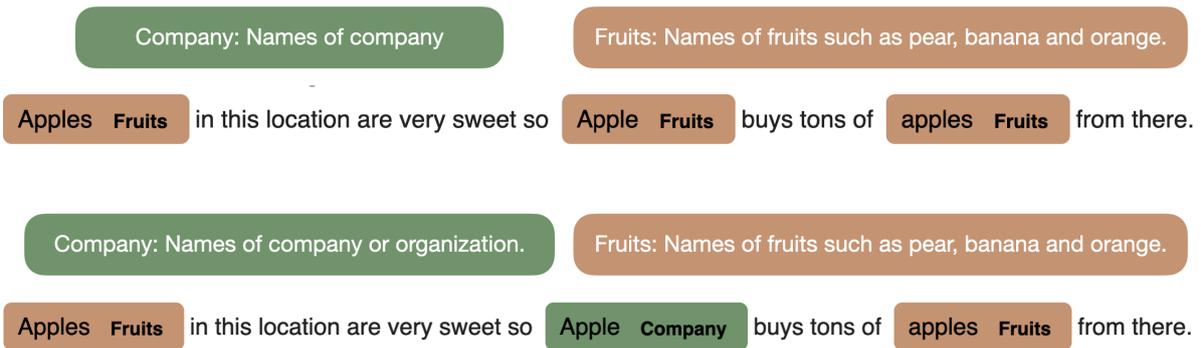


Figure 1: A small modification of the *Company* class description results in different entity predictions.

081 choose a proper description for black-box zero-shot  
082 models, in particular in an unfamiliar domain.

083 In this paper, we study different methods for  
084 boosting model performance with automatic de-  
085 scription enhancement. Specifically, we propose  
086 UDEBO (for Unsupervised DDescription BOosting),  
087 the first unsupervised method capable of automati-  
088 cally modifying/generating description to improve  
089 entity (or relation) predictions in the zero-shot set-  
090 tings. We present several strategies to alter descrip-  
091 tions, such as using a generative model, paraphras-  
092 ing, and summarization combined with description  
093 ranking/ensemble methods to reduce model uncer-  
094 tainty and increase overall performance. We em-  
095 pirically evaluate the performance of UDEBO on 4  
096 existing standard zero-shot datasets, spanning two  
097 tasks: (i) name entity classification and (ii) relation  
098 classification.

099 Our results show that for the zero-shot entity  
100 classification tasks, UDEBO improved the results  
101 of state-of-the-art models by 7 and 1.3 percentage  
102 points in terms of Macro F1 Score in the OntoNotes  
103 and MedMentions datasets, respectively. For what  
104 concerns relation classification, we achieve a per-  
105 formance improvement of 6 and 3 percentage  
106 points (Macro F1 Score) on the FewRel and Wik-  
107 iZS datasets over our baseline models, respectively.

108 We organize the paper as follows. In Section 2  
109 we provide a description of the zero-shot setting  
110 for entity recognition and relation classification.  
111 We also formally define the problem we aim to  
112 solve in this paper, i.e. how to enhance entity or  
113 relation descriptions to improve the performance  
114 of zero-shot models. In Section 3 we describe  
115 the proposed approaches for description boosting  
116 while in Section 4 we describe our experimental  
117 setup and results. We provide a literature review  
118 and draw the conclusions of our work in Sections

5 and 6, respectively.

## 2 Preliminaries and problem definition

119 Entity and relation classification are key steps to  
120 extract or query knowledge from unstructured doc-  
121 uments. Zero-shot approaches can identify which  
122 tokens in a text refer to an entity (its mention) and  
123 determine its type (entity typing) without the need  
124 of observing other instances of the same entity dur-  
125 ing training.

126 In ZSL, the sets of training and test entity (or  
127 relation) classes are disjoint. Therefore, the strat-  
128 egy employed by zero-shot models is to rely on  
129 prior general knowledge that could be transferred  
130 to unseen instances at inference time. In particu-  
131 lar, novel zero-shot approaches leverage the fact  
132 that textual descriptions for entity classes are ei-  
133 ther available in existing datasets or can be easily  
134 provided by users.

135 Given a textual description of an entity class  
136 (or relation) of interest, zero-shot models recog-  
137 nize mentions in a text and predict whether the  
138 given mentions belong or not to the entity class  
139 (or relation) with a certain probability. One classic  
140 paradigm is to embed all entities with their tex-  
141 tual description and the input sentence with each  
142 mention into one common space and measure the  
143 probability of each entity by assessing their dis-  
144 tance. Descriptions for model pre-training are typi-  
145 cally sourced from Wikipedia by joining an entity  
146 page title or label with the first 10 sentences in the  
147 respective Wikipedia page (Wu et al., 2020). How-  
148 ever, the quality of the descriptions has an impact  
149 on how effective the transfer of knowledge from  
150 observed to unseen entities (Aly et al., 2021).

151 Given a set of entity classes  $E$  (or relations) of  
152 interest with their textual descriptions  $D$  and a cor-  
153 pus of sentences  $S$  to annotate as input, we describe  
154  
155

in Section 3.1 different strategies to generate new entity (or relation) descriptions  $D'$  for the input set  $E$ , intending to improve the accuracy of the predictions by the ZSL models over that corpus. 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 (or relation) description  $D$  for annotating an input text corpus  $S$ . Our goal is to generate a set of descriptions  $D^*$  such that:

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

As exemplified in Figure 1, 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

We begin our discussion with methods for generating description variations before providing details about description ranking and ensemble strategies in the following subsections.

#### 3.1 Generating description variations

Improving the completeness or clarity of entity (or relation) 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 (or relation) 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 (or relations) 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 (or relation), we propose to use a zero-shot model to first compute the probabilities of classes for each mention (or relation) 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 (or relation) descriptions. The main idea

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

Table 1: Number of sentences and entities (or relations) for each split of the considered datasets.

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 predict the label of  $\textit{London}$  as  $e_2 = \textit{Location}$  and the rest of the pipelines predict  $\textit{Bridge}$  as  $\textit{Facility}$ . Therefore, the accumulated number of votes for the span  $\textit{London Bridge}$  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  $\textit{London Bridge}$  is  $\textit{Facility}$ . Once the span  $\textit{London Bridge}$  has been assigned a label, all of its sub-spans become redundant and thus are removed from consideration.

**Relation description ensemble.** For each set of descriptions generated using the strategies discussed in Subsection 3.1, we run a pipeline to obtain the predicted relation for each provided pair of entities. The votes of all the relations are aggregated across different pipelines. We use the majority voting rule to select the relation with the highest aggregated number of votes from different pipelines. That relation is considered as the output relation for the given pair of entities.

## 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.<sup>2</sup> 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 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-BERT<sup>3</sup> (Chen and Li, 2021b), 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.<sup>3</sup> 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.

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

<sup>3</sup><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	<b>51.40</b>	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	<b>31.14</b>	46.51	<b>36.78</b>	<b>36.15</b>	<b>88.29</b>
MedMentionsZS	SMXM	16.79	<b>40.55</b>	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	<b>19.51</b>	32.73	<b>23.86</b>	<b>22.97</b>	<b>85.70</b>

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	<b>28.38</b>	<b>25.68</b>	<b>25.68</b>	<b>22.12</b>	<b>25.68</b>
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	<b>34.79</b>	<b>37.11</b>	<b>40.17</b>	<b>34.25</b>	<b>40.17</b>

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

## 4.2 Empirical results

This section discusses the results of entity (or relation) 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 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.

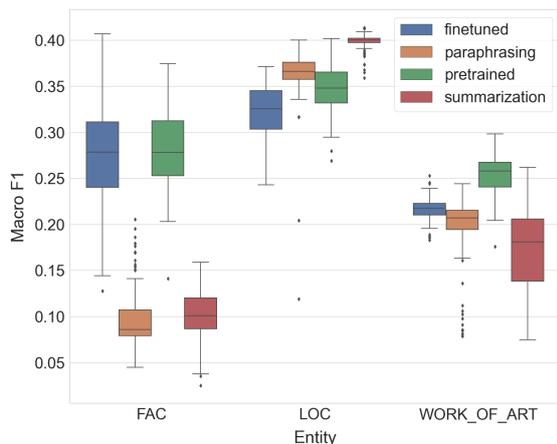


Figure 2: The figure shows the distributions of Macro F1 Score values on the test split of the OntoNotesZS dataset for each class, using the strategies described in Section 3.1 to generate 100 description variations for each class.

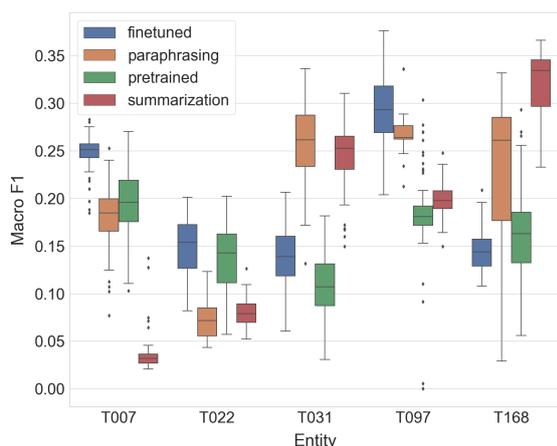


Figure 3: The figure shows the distributions of Macro F1 Score values on the test split of the MedMentionsZS dataset for each class, using the strategies described in Section 3.1 to generate 100 description variations for each class.

#### 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 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, indirectly unveiling the mechanism behind zero-shot predictions.

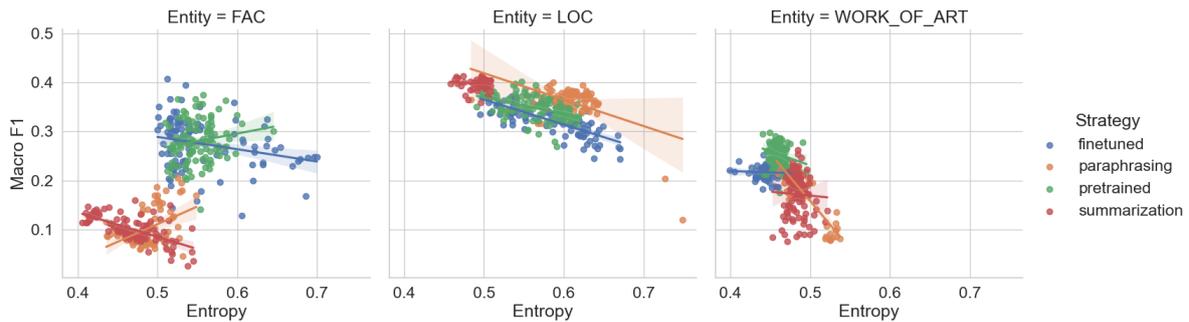


Figure 4: Analysis of the correlation between entropy and Macro F1 Score on unseen classes on the OntoNotesZS test split. Entropy can be calculated without the need for labeled data, therefore, if a correlation exists it can be used as an unsupervised heuristic to select descriptions that improve model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

## 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 (Akbik 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 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 *Med-Mentions* (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. Mention 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

493 to scale the approach from Wikipedia (5.9M enti- 544  
494 ties) to Wikidata (90M entities). 545

495 The analyses in (Aly et al., 2021) show that while 546  
496 Wikipedia descriptions work well on general entity 547  
497 types, they perform poorly on domain specific data, 548  
498 e.g. *MedMentions*. They also show the impact 549  
499 of using annotation guidelines for descriptions to 550  
500 improve the transfer of knowledge from observed 551  
501 to unseen entities. The adoption of this approach 552  
502 led to a better performance compared to using a 553  
503 class name itself or Wikipedia passages. In par- 554  
504 ticular, description vagueness, noise and negations 555  
505 had a negative effect, while annotation guidelines, 556  
506 including explicit examples and syntactic and mor- 557  
507 phological cues, improved the performance. 558

508 **Zero-shot relation classification.** Textual de- 549  
509 scriptions have also been employed in the rela- 550  
510 tion classification task to predict new relations that 551  
511 could not be observed at training time. For exam- 552  
512 ple, ZS-BERT (Chen and Li, 2021c) learns two 553  
513 functions – one to project sentences and the other 554  
514 to project relation descriptions into an embedding 555  
515 space. The objective is first to jointly minimise 556  
516 the distance between the embedding vectors for 557  
517 an input sentence and the relation description for 558  
518 positive entity pairs and then to classify the relation 559  
519 (using a softmax layer to produce a classification 560  
520 probability). At inference time, the prediction of 561  
521 unseen relation classes can be achieved through 562  
522 nearest neighbor search. Overall, using descrip- 563  
523 tions seems to improve existent zero-shot methods 564  
524 and expand their domains of application. Still, de- 565  
525 scriptions are not always good enough to get good 566  
526 predictions. Improving the accuracy of these ap- 567  
527 proaches remains an open challenge. The better 568  
528 the separation between embedding of different re- 569  
529 lations, the more accurate the model predictions, 570  
530 however, as the number of unseen relations in- 571  
531 creases, it becomes more difficult to predict the 572  
532 right one (Chen and Li, 2021c). 573

533 Existent ZSL methods usually rely on external 574  
534 knowledge from KGs, ranging from textual in- 575  
535 formation, class attributes, hierarchy, domain and 576  
536 range constrains and relations to logic rules. There 577  
537 are relatively few studies evaluating their perfor- 578  
538 mance for unseen relations, a comparison using dif- 579  
539 ferent external knowledge settings for zero-shot re- 580  
540 lation classification and KG completion can be seen 581  
541 in (Geng et al., 2021). To the best of our knowl- 582  
542 edge, we present the first approach to automatically 583  
543 predict and generate entity (or relation) descrip-

544 tions to improve the accuracy of entity recognition 545  
and relation classification models.

**Query auto completion in information retrieval systems.** 546  
in relation to our work, query auto com- 547  
pletion is the problem where a computer extends 548  
the initial parts of user queries to a search engine 549  
to save users time and enhance search performance 550  
(Cai et al., 2016). Most query auto completion ap- 551  
proaches are based on mining query logs (Whiting 552  
and Jose, 2014). The most related approach to our 553  
work is based on personalised LMs fine-tuned on 554  
users’ historical data (Jaech and Ostendorf, 2018). 555  
The key difference between our work and the query 556  
auto completion setting is that in the context of 557  
named entity recognition, we don’t have histori- 558  
cal data to learn from. Moreover, the objective of 559  
query extension is to maximise the retrieval doc- 560  
uments accuracy while named entity recognition 561  
looks at the descriptions that maximise the entity 562  
annotation accuracy. 563

## 6 Conclusion and future work 564

565 In this paper, we formally defined the problem of 566  
567 selecting descriptions to make predictions about 568  
569 unseen classes in the ZSL context. We empiri- 570  
571 cally evaluated the sensitivity of two ZSL methods 572  
573 to description changes, and proposed 4 different 574  
575 strategies to enhance them using the implicit knowl- 576  
577 edge of pre-trained language models. We also stud- 578  
579 ied in detail the efficacy of the proposed entropy- 580  
581 based heuristic to rank different description formu- 582  
583 lations, analyzing its correlation with the perfor- 584  
585 mance (in terms of Macro F1 Score) of the model. 586  
587 We observed a negative correlation between the 588  
589 proposed heuristic and Macro F1 Score on two out 590  
591 of four of the considered datasets (OntoNotesZS 592  
and FewRel). The same assumption however was 593  
not valid for the other datasets (MedMentionsZS 594  
and WikiZS), thus motivating the need to develop 595  
more effective heuristics in future. Finally, we de- 596  
scribed the UDEBO method, which combines the 597  
predictions obtained by the same model using dif- 598  
ferent automatically generated variants of entity 599  
and relation descriptions. Our experimental results, 600  
on 4 different datasets, spanning across two differ- 601  
ent NLP tasks (Entity Classification and Relation 602  
Classification) showed how UDEBO outperforms 603  
the baselines by a significant margin and achieves 604  
new state-of-the-art results on these benchmarks 605  
under the zero-shot setting. 606

## 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.
- 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.
- Fei Cai, Maarten De Rijke, et al. 2016. A survey of query auto completion in information retrieval. *Foundations and Trends® in Information Retrieval*, 10(4):273–363.
- Chih-Yao Chen and Cheng-Te Li. 2021a. Zsbert: Towards zero-shot relation extraction with attribute representation learning. *arXiv preprint arXiv:2104.04697*.
- Chih-Yao Chen and Cheng-Te Li. 2021b. ZS-BERT: towards zero-shot relation extraction with attribute representation learning. *CoRR*, abs/2104.04697.
- Chih-Yao Chen and Cheng-Te Li. 2021c. ZS-BERT: towards zero-shot relation extraction with attribute representation learning. *CoRR*, abs/2104.04697.
- 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.
- 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.
- Aaron Jaech and Mari Ostendorf. 2018. Personalized language model for query auto-completion. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 700–705, Melbourne, Australia. Association for Computational Linguistics.
- 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.
- 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*.

703 Stewart Whiting and Joemon M Jose. 2014. Recent and  
704 robust query auto-completion. In *Proceedings of the*  
705 *23rd international conference on World wide web*,  
706 pages 971–982.

707 Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian  
708 Riedel, and Luke Zettlemoyer. 2020. [Scalable zero-](#)  
709 [shot entity linking with dense entity retrieval](#). In  
710 *Proceedings of the 2020 Conference on Empirical*  
711 *Methods in Natural Language Processing (EMNLP)*,  
712 pages 6397–6407, Online. Association for Computa-  
713 tional Linguistics.

714 Yongqin Xian, Christoph H. Lampert, Bernt Schiele,  
715 and Zeynep Akata. 2019. [Zero-shot learning—a com-](#)  
716 [prehensive evaluation of the good, the bad and the](#)  
717 [ugly](#). *IEEE Transactions on Pattern Analysis and*  
718 *Machine Intelligence*, 41(9):2251–2265.

719 Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Pe-  
720 ter J. Liu. 2019. [Pegasus: Pre-training with extracted](#)  
721 [gap-sentences for abstractive summarization](#).

722 Susan Zhang, Stephen Roller, Naman Goyal, Mikel  
723 Artetxe, Moya Chen, Shuohui Chen, Christopher De-  
724 wan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022.  
725 Opt: Open pre-trained transformer language models.  
726 *arXiv preprint arXiv:2205.01068*.

## A 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 Huggingface<sup>4</sup> 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).<sup>2</sup> 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.
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.

<sup>4</sup>[https://huggingface.co/datasets/conll2012\\_ontonotesv5](https://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.

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 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

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.

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.<sup>5</sup> 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 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.

<sup>5</sup>[https://huggingface.co/datasets/few\\_rel](https://huggingface.co/datasets/few_rel)

**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.<sup>6</sup> 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 Hub<sup>7</sup> 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.<sup>8</sup> 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.<sup>9</sup> We used *min\_length* = 80, *max\_length* = 512, *num\_beams* = 8, *temperature* = 1 and *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.<sup>10</sup> We used *min\_length* =

<sup>6</sup><https://huggingface.co/gpt2>

<sup>7</sup><https://huggingface.co/t5-large>

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

<sup>9</sup>[https://huggingface.co/mrm8488/bert-small2bert-small-finetuned-cnn\\_daily\\_mail-summarization](https://huggingface.co/mrm8488/bert-small2bert-small-finetuned-cnn_daily_mail-summarization)

<sup>10</sup>[https://huggingface.co/tuner007/pegasus\\_paraphrase](https://huggingface.co/tuner007/pegasus_paraphrase)

Dataset	Split	Mean #Tokens	Max #Tokens	Min #Tokens	Mean #Entities	Max #Entities	Min #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.

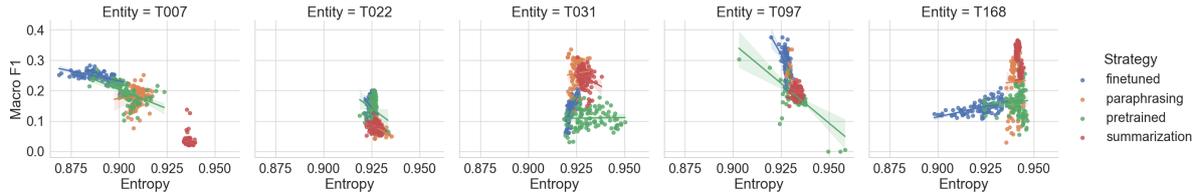


Figure 5: Analysis of the correlation between entropy and Macro F1 Score on unseen classes on the MedmentionsZS test split.

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

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, stage-play, 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 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.