ZERO-SHOT LEARNING OVER LARGE OUTPUT SPACES: UTILIZING INDIRECT KNOWLEDGE EXTRACTION FROM LARGE LANGUAGE MODELS

Jinbin Zhang, Nasib Ullah

Aalto University
Finland
{jinbin.zhang, nasibullah.nasibullah}@aalto.fi

Rohit Babbar

University of Bath United Kingdom rb2608@bath.ac.uk

ABSTRACT

Extreme Multi-label Learning (XMC) is a task that allocates the most relevant labels for an instance from a predefined label set. Extreme Zero-shot XMC (EZ-XMC) is a special setting of XMC wherein no supervision is provided; only the instances (raw text of the document) and the predetermined label set are given. The scenario is designed to address cold-start problems in categorization and recommendation. Traditional state-of-the-art methods extract pseudo labels from the document title or segments. These labels from the document are used to train a zero-shot bi-encoder model. The main issue with these generated labels is their misalignment with the tagging task. In this work, we propose a framework to train a small bi-encoder model via the feedback from the large language model (LLM), the bi-encoder model encodes the document and labels into embeddings for retrieval. Our approach leverages the zero-shot ability of LLM to assess the correlation between labels and the document instead of using the low-quality labels extracted from the document itself. Our method also guarantees fast inference without the involvement of LLM. The performance of our approach outperforms the SOTA methods on various datasets while retaining a similar training time for large datasets.

Keywords LLM · XMC · Zero-shot · Tagging

1 Introduction

Extreme multi-label text classification (XMC) involves the task of assigning relevant labels to documents from an extensive pool of all possible labels [1]. The size of label space in typical applications of XMC, such as product-to-product recommendations, product search[2], labeling Wikipedia pages [3], and categorizing Amazon products [4], often ranges from hundreds of thousands to millions. Despite its wide-spread application, current supervised learning models for XMC heavily rely on expert-annotated (Wikipedia) or user-annotated labels (Amazon purchase history) to train the models and the label set remains fixed throughout the training and prediction process. Furthermore, the framework of supervised XMC setting faces two additional challenges. Firstly, it is difficult to get annotations due to the vast number of labels involved. This makes it very challenging for annotators to select the relevant labels from such a large pool, which could lead to missing labels [5; 6; 7; 8; 9]. Secondly, the emergence of new labels is a common occurrence in XMC, especially in scenarios like cold start, where most of the labels are unseen (hence remain unannotated), and there is also a need for the new ones to be added to the label set. Most traditional XMC algorithms are not able to handle the unseen labels during the inference process, making them incapable to effectively address these complexities.

There are two distinct settings for zero-shot extreme classification: (i) Generalized Zero-Shot Extreme Multi-label Learning (GZXML) [10], which enables the model to make predictions for labels that have not been encountered before. Nevertheless, the training dataset necessitates annotation, even if certain labels are absent from the training data. Hence, it is not appropriate for situations when there is no annotated data available, such as cold start scenarios. (ii) Extreme

Zero-Shot Extreme Multi-Label Text Classification (EZ-XMC) [11; 12], which is designed to handle unseen labels even in the absence of an annotated training dataset. The latter setting of EZ-XMC is the one which we follow in this work.

The current methods in EZ-XMC primarily concentrate on training a resilient sentence embedding encoder by leveraging the pseudo positive labels generated from the document itself. This approach enables encoding label texts into embeddings through the sentence encoder, facilitating efficient retrieval aligned with the document embedding. This approach obviates the need for the training dataset to cover the entire spectrum of labels. For instance, as shown in Figure 2, MACLR [11] proposed constructing instance & pseudo-labels pairs using the (content, title) combinations from documents, while RTS [12] randomly split the document and choose two spans to form the instance & pseudo labels pairs. However, these methods often neglect the direct matching relationship between pairs of document and pseudo label. For example, the two segments in one document, as constructed by RTS [12], help the model grasp semantic similarity, but they fall short in determining whether the label truly represents the document (see Figure 1). Furthermore, the pseudo labels might not harmonize effectively with the

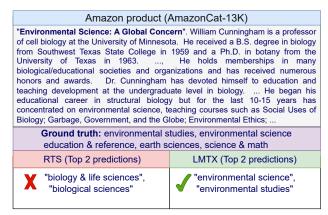


Figure 1: RTS [12] predicts labels that are semantically similar, whereas our LMTX can predict relevant labels that match the document.

domain of the predefined label set, resulting in a discrepancy between the target task and the available training pairs.

Recently, even though Large Language Models (LLMs) have exhibited impressive reasoning and zero-shot capabilities across various NLP tasks [13; 14; 15; 16; 17; 18; 19; 20; 21], with the exception of the only work [22; 23], these remain unexplored in the XMC setting. This is primarily due to the additional computational costs of deploying LLMs in inference stage. To mitigate this issue, we employ a relevance assessment strategy using an LLM to meticulously select pertinent pseudo labels from a short-listed label set for each document, thereby make it possible to train a lightweight model while inheriting the knowledge of the LLM. A pictorial description of the distinction in constructing pseudo positive labels between current EZ-XMC algorithms (MACLR [11] and RTS [12]) and our proposed method is illustrated in Figure 2. It may be noted that, in the EZ-XMC setting, there does not exist an instance-wise correspondence between documents and labels in the training set, and these methods mainly differ in the way the instance and label sets are used in their respective pipelines.

In this work, we present, LMTX, a novel approach that utilizes an LLM as a Teacher for eXtreme zero-shot classification. LMTX employs a bi-encoder, which acts as a retrieval framework, and consists of two encoders — one each for encoding the documents and labels into their corresponding embeddings. This framework is well-suited for the extreme zero-shot scenario, where label embeddings can be efficiently retrieved through maximum inner-product search (MIPS) [24]. The LLM is not utilized to directly generate labels for each document; instead, it focuses on evaluating the relevance of labels from a narrowed-down shortlist, rather than assessing the entire label set. Following this relevance assessment, the verified relevant labels are then provided to the bi-encoder, which uses these to fine-tune and optimize the training process. By adopting this strategy, we can benefit from the capabilities of the LLM model while ensuring swift inference via MIPS. Our contributions can be summarized as follows:

- LMTX introduces a novel training approach for bi-encoders, emphasizing a curriculum-based method that
 dynamically adjusts based on the relevance feedback from an LLM by leveraging its zero-shot learning
 abilities,
- The proposed LMTX requires less training data because there is a higher correlation between the pseudolabels and documents, resulting in higher-quality training pairs. Consequently, our approach achieves better performance while maintaining similar or reduced training time compared to traditional methods on some large datasets.
- LMTX significantly outperforms current state-of-the-art methods, demonstrating comprehensive advancements in performance metrics. This is particularly evident in the improvement of prediction performance over exising approaches, where it shows 18-38% improvement on a range of benchmark datasets with upto 500,000 labels.

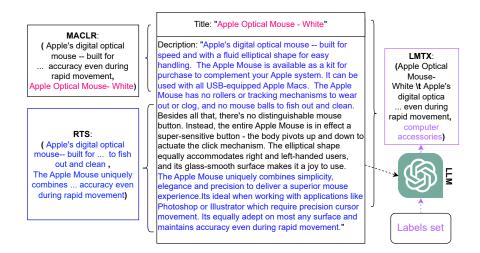


Figure 2: An example of how to construct training pairs using state-of-the-art methods MACLR [11] and RTS [12]. MACLR utilizes the 'Title' of a document to generate pseudo labels, while the 'Description' serves as the training document. Conversely, RTS forms its training pairs by selecting two random segments solely from the 'Description'. Differently, our proposed model, LMTX, adopts a more refined approach. It selects 'computer accessories' as a pseudo positive label from a predefined set, a choice validated by the LLM model.

2 Proposed Method

In this section, we start by providing a definition of the problem statement for extreme zero-shot extreme multi-label classification (EZ-XMC), and present our framework in detail.

2.1 Problem Definition

Let's denote $X_i \in \mathcal{X}$ as the text for an instance in a particular domain; i.e., X_i could be the textual description for a product on Amazon. Unlike the standard (supervised) extreme multi-label classification problem, the key characteristic of the EZ-XMC setting is that we do not have the corresponding well-annotated labels Y_i for each training instance X_i . However, besides having the original text of instances $\{X_i\}_{i=1}^N$, we also have access to the predetermined labels along with their texts, i.e., we have $\{l_k\}_{k=1}^L$. We refer to this collection of predetermined labels as the "labels set". The goal of EZ-XMC, which is the one that we consider in this paper, is to assign the document $X_i \in \mathcal{X}$ to set of labels $\{l_j\} \subseteq \{l_k\}_{k=1}^L$ that are relevant to the document. To achieve this objective, the task requires learning a mapping function from text to embedding, denoted as $\mathcal{E}_\theta: \mathcal{X} \to \mathbb{S}^{D-1}$, where θ represents the training parameters, \mathcal{E} represents the encoder for documents and labels, and \mathbb{S}^{D-1} is the D-dimensional unit sphere. The mapping function is typically implemented as a bi-encoder, where both the text of documents and labels are embedded within \mathbb{S}^{D-1} .

2.2 Bi-Encoder Model

To extract the respective embeddings for the document and labels text, we introduce the bi-encoder architecture \mathcal{E}_{θ} . It consists of two encoders: one for encoding the text of documents, and the other for encoding the label text. It must be noted that the weights between both the encoders are shared. The embeddings of document and label can be expressed as follows: $\mathcal{E}_{\theta}(X_i)$ for document X_i and $\mathcal{E}_{\theta}(l_k)$ for label l_k , where X_i represents the document and l_k represents the label text. The relevance score between the document X_i and the label l_k is determined by the cosine similarity of $\mathcal{E}_{\theta}(X_i)$ and $\mathcal{E}_{\theta}(l_k)$. A pictorial depiction of the bi-encoder part is shown in the Figure 3 and is built upon the Distill-Bert transformer [25] as the base model.

2.3 Training the Bi-encoder from the Feedback of LLM

Training process overview: Our training methodology adopts an iterative framework, encompassing three distinct stages within each cycle. In the first stage, we embed all the documents and labels, followed by constructing an Approximate Nearest Neighbor Search (ANNS [26]) over these label embeddings. This process facilitates the retrieval of a refined set of label candidates for each document. Subsequently, in the second stage, the LLM is deployed to

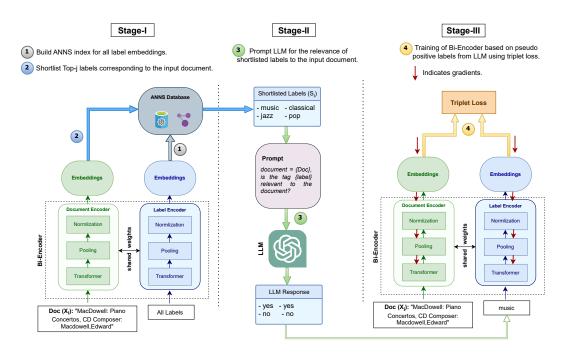


Figure 3: The process of getting feedback from LLM model for training the bi-encoder. First, for a given document, the (pre)trained bi-encoder and ANNS are employed to create a short-list of potential labels. Next, the LLM assesses the relevance between the labels in this shortlist and the document. Finally, the selected labels are utilized to further train the bi-encoder.

scrutinize these candidates, effectively identifying pseudo positive labels for each document. The final stage involves the training of the bi-encoder model, utilizing the labels identified in the previous stage. The subsequent section provides a detailed exposition of these three stages. Additionally, Figure 3 illustrates the mechanism through which the bi-encoder acquires feedback from the LLM and progresses through its training regimen.

Data embedding & shortlist generation (stage-I): The LLM model demonstrates zero-shot ability in determining relevance between two text segments [15]. However, this approach encounters practical challenges when applied to a vast array of labels, as in our context. Specifically, the computational complexity involved in assessing the relevance between each document and every label in a large set becomes formidable, being $\mathcal{O}(NL)$ in complexity. This can be quite prohibitive, even for a dataset with a moderate number $(\mathcal{O}(10^3))$ of instances and labels. To mitigate this, our strategy involves condensing the label space presented to the LLM. We utilize the (pre)trained bi-encoder to process the document and label text into embeddings. Once the embeddings are learnt, we utilize the ANNS to efficiently select the top-j most relevant labels for each document. These selected labels, denoted as $S_i = \{l_{i1}, l_{i2}, ..., l_{ij}\}$, constitute a focused subset for subsequent processing.

LLM model as a teacher (stage-II): Once we obtain the label shortlist S_i for the i-th document, we can employ the LLM as a teacher to determine the relevance between the document and the top-j labels in a shortlist. Let X_i denote a particular document and l_{ik} be its k-th label in the shortlist. To assess the relevance between X_i and l_{ik} , we instruct the LLM with the question, "document = $\{X_i\}$, is the tag $\{l_{ik}\}$ relevant to the document? answer yes or no". If the LLM outputs "Yes", we consider l_{ik} to be relevant to X_i . Conversely, if the model outputs "No", we consider l_{ik} as an unrelated label and discard it. We keep all the labels from the shortlist that received a positive feedback "yes") from the LLM. Then, we use these selected relevant labels to train the bi-encoder model. Our analysis in Section 2.10 shows that using labels which are rejected by the LLM, as hard negatives, leads to degradation in prediction performance. A detailed discussion of different prompts used for the LLM can also be found in Appendix 5.1.

Training bi-encoder with pseudo positive labels (stage-III): To train the bi-encoder, we follow the training procedure in [27]. Out of the labels identified by the LLM as the pseudo positives, we choose only one of the pseudo positive labels for each document during the training process. This is shown to help in achieving faster convergence in the earlier work [27]. Regarding the negatives, which we need to compute the instance-wise loss, we use in-batch negative sampling, in which the negatives for a document come from the pseudo positive labels of other documents in the same batch. For the label l_k , the predicted relevance score between document X_i and l_k is computed through the cosine

Algorithm 1 Training the bi-encoder with the feedback from LLM teacher (LMTX)

```
Input: Initial bi-encoder \mathcal{E}_{\theta}, LLM model \mathcal{M}_{LLM}, data instances \{X_i\}_{i=1}^N, labels set \{l_k\}_{i=1}^L, dev set instances \{X_j\}, and number
of cycles T
Output: Trained bi-encoder \mathcal{E}_{\theta}
 1: c = 0.
 2: while c < T do
        Compute \mathcal{E}_{\theta}(X_i), \mathcal{E}_{\theta}(l_k) for all \{X_i\}_{i=1}^N and \{l_k\}_{i=1}^L
        Retrieve top labels S_i = ANNS(\mathcal{E}_{\theta}(X_i), \mathcal{E}_{\theta}(l_k)_{k=1}^L) for each X_i
 4:
 5:
        Fetch pseudo positive labels P_i^+ = \mathcal{M}_{LLM}(X_i, S_i) for all X_i
        for i=0 to N_batches do
 6:
            Sample a mini_batch B_i = \{X_i, P_i^+\} where, |B_i| = m
 7:
 8:
            Update \mathcal{E}_{\theta} using mini-batch B_i, loss \mathcal{L} and AdamW optimizer.
 9:
        Evaluate \mathcal{E}_{\theta} with \mathcal{M}_{LLM} on the dev set \{X_j\} and obtain P@1 over pseudo labels.
10:
11:
         if P@1 does not improve on dev dataset then
12:
             Stop training cycle
13:
         end if
14:
         c = c + 1
15: end while
16: return model \mathcal{E}_{\theta}
```

similarity $\langle f_e(X_i), f_e(l_k) \rangle$, and triplet loss is used to train the bi-encoder [28; 29; 27]:

$$\mathcal{L} = \sum_{i=1}^{N} \sum_{k'} [\langle f_e(X_i), f_e(l_{k'}) \rangle - \langle f_e(X_i), f_e(l_p) \rangle + \gamma]_+ \tag{1}$$

where γ is the margin, the k' stands for the index of hard negative labels from the mini batch, $l_{k'}$ and l_p correspond to the text of the negative labels and the pseudo positive label.

As the training progresses, the bi-encoder gradually improves, leading to an enhancement in the quality of labels within the shortlist and an increased relevance to the corresponding document. During training, we evaluate the model on the development dataset and choose the best one based on performance evaluated by the LLM since under the EZ-XMC setting one does not have access to annotated ground-truth labels. If there is no performance improvement on the development dataset, the training is halted, so the number of cycles is actually determined by the performance on the development dataset. The pseudo code of the proposed algorithm LMTX, for training the bi-encoder model with feedback from LLM, is presented in Algorithm 1.

2.4 Inference

The model's inference procedure is analogous to the formation of the shortlist during training, as depicted in Stage-I of Figure 3. Initially, we extract text embeddings for the document and the labels from the pre-saved model. Subsequently, we build the MIPS [24] over these label embeddings. For each document, we employ its embedding as a query to retrieve the top-m labels, which ultimately serve as the predicted results. The use of MIPS in the inference process ensures a sublinear time complexity for each instance. The label embedding extraction and construction of MIPS index are performed just once, hence amortizing the cost of this step. To assess the model's performance, we carry out inference on the annotated test

Dataset	N	N_{test}	N_{label}
EURLex-4K	15,511	3,803	3,956
Wiki10-31K	14,146	6,616	30,938
AmazonCat-13K	1,186,239	306,782	13,330
LF-WikiSeeAlso-320K	693,082	177,515	312,330
LF-Wikipedia-500K	1,813,391	783,743	501,070

Table 1: The statistical information for the datasets: N, N_{test} and N_{label} are the number of training samples, test samples, and the total number of labels.

dataset and then compare the top-m labels obtained from MIPS with the annotated ground truth.

¹https://github.com/facebookresearch/faiss

2.5 Dataset

We used five datasets for evaluation. EURlex-4k, Wiki10-31k, and AmazonCat-13K were obtained from the XLNet-APLC repository², while the remaining datasets were downloaded from the extreme classification repository³. For a detailed statistical information regarding all these datasets, please refer to Table 1. Due to the extensive and time-consuming process of LLM judgement, we have decided to limit the training data for AmazonCat-13K, LF-WikiSeeAlso-320K and LF-Wikipedia-500K datasets to only 30,000 documents each. In contrast, the baseline models utilize the entire dataset rather than just a subset.

2.6 Evaluation Metrics

We employ the commonly used evaluation metrics [30; 31; 12] for the EZ-XMC setting : Precision@k(P@m) and Recall@m(R@m).

$$P@m = \frac{1}{m} \sum_{i \in rank_m(\hat{y})} y_i \tag{2}$$

$$R@m = \frac{1}{\sum_{l} y_{l}} \sum_{i \in rank_{m}(\hat{y})} y_{i}$$
(3)

where $\hat{y} \in \mathbb{R}^L$ represents a vector containing the predicted labels' score for each instance, while $y \in \{0,1\}^L$ corresponds to a vector representing the ground truth for each document. The term $rank_m(\hat{y})$ refers to a list of the predicted top-m label indices. The definition of the two metrics applies to a single instance; for multiple instances, the performance is the average across all instances.

2.7 Implementation Details

Bi-Encoder: In our bi-encoder framework, we adopt a siamese network architecture for sentence encoding. The core of this network is DistilBERT [25], comprising six transformer layers. For the generation of sentence embeddings, we apply mean pooling, yielding embeddings of 768 dimensions. The bi-encoder is initialized using the msmarco-distilbert-base-v4⁴, and ANNs is built via the HNSW package⁵. For optimization, we employ the AdamW optimizer [32] with a learning rate of 0.0002, setting the batch size to 128. All experiments for training the bi-encoder are conducted on a single A100 GPU. Following the supervised method in [27], γ is set to 0.3. To assess performance, a development set of 800 documents is randomly selected from the training dataset, with pseudo labels derived from the top-k labels as determined by the LLM model.

LLM inference: For our Large Language Model (LLM) component, we employ the WizardLM-13B-V1.0 model [33], an open-source LLM notable for achieving 89.1% of GPT-4's [34] performance with approximately 13 billion parameters. In addition, for the purposes of this study, we incorporate Llama2 [35] and vicuna-13b-v1.3 [36] models in our ablation experiments to serve as comparative benchmarks. The computational setup for these LLM models involves the utilization of 2 × A100 GPUs. Instances encoded into the LLM are truncated to 430 tokens.

2.8 Baselines

We have incorporated state-of-the-art extreme zero-shot extreme classification models as our baselines.

- GloVe [37]: The method takes the average of Glove word embeddings to represent the sentence and then retrieves the top-m labels by comparing their cosine similarity.
- Inverse Cloze Task (ICT) [38]: An unsupervised sentence embedding method which constructs the pseudo pairs with (title, document).
- SentBERT [39]: A siamese network structure used for sentence similarity. The model underwent fine-tuning on the sentence matching dataset.

²https://github.com/huiyegit/APLC XLNet

³http://manikvarma.org/downloads/XC/XMLRepository.html

⁴https://huggingface.co/sentence-transformers/msmarco-distilbert-base-v4

⁵https://github.com/kunaldahiya/pyxclib

Method	P@1	P@3	P@5	R@1	R@3	R@5	R@10	R@100
					Lex-4K			
GloVe	1.66	1.11	1.04	0.37	0.73	1.08	1.88	13.72
SentBERT	8.52	7.70	6.83	1.7	4.54	6.69	10.20	32.45
SimCSE	5.86	4.44	3.85	1.20	2.86	3.93	6.12	23.66
MPNet	10.81	8.65	7.21	2.27	5.28	7.28	10.85	35.00
Msmarco-distilbert	15.91	9.89	7.81	3.33	6.16	8.08	11.22	32.82
RTS§	30.58	21.54	17.73	6.19	13.01	17.72	25.34	59.61
LMTX	47.12 [†]	28.87 [†]	21.80 [†]	9.59 [†]	17.42 [†]	21.78 [†]	28.06 [†]	50.69
				Wiki	10-31K			
GloVe	3.87	3.11	2.87	0.24	0.57	0.89	1.48	6.63
SentBERT	9.39	6.93	5.81	0.60	1.31	1.81	2.70	8.89
SimCSE	23.55	17.21	14.01	1.42	3.07	4.13	6.01	15.10
MPNet	44.82	29.18	22.38	2.63	5.12	6.52	8.89	21.58
RTS	47.73	31.03	23.65	2.81	5.41	6.84	9.12	20.58
Msmarco-distilbert§	54.17	33.44	25.38	3.18	5.82	7.32	9.70	19.50
LMTX	57.89 [†]	38.00 [†]	29.09 [†]	3.41 [†]	6.68 [†]	8.46 [†]	11.14 [†]	21.79 [†]
				Amazor	Cat-13K			
GloVe	4.83	3.89	3.42	0.99	2.46	3.67	6.05	24.46
SentBERT	5.21	4.22	3.68	0.99	2.34	3.37	5.35	17.71
SimCSE	2.84	2.60	2.42	0.52	1.41	2.17	3.75	15.25
ICT	15.52	10.48	8.34	2.91	5.93	7.86	11.04	27.79
MPNet	18.01	12.84	10.51	3.63	7.68	10.48	15.73	43.00
MACLR	10.66	6.75	5.14	1.98	3.79	4.81	6.35	13.42
Msmarco-distilbert	16.36	10.96	8.68	3.29	6.62	8.73	12.23	29.08
RTS§	18.89	13.59	11.07	3.69	8.03	10.97	16.20	41.50
LMTX	26.18 [†]	17.45 [†]	13.39 [†]	5.52 [†]	10.95 [†]	13.83 [†]	17.95 [†]	36.10
				-WikiSe	eAlso-32			
GloVe*	3.86	2.76	2.21	2.12	4.11	5.22	6.95	15.33
SentBERT*	1.71	1.27	1.06	1.08	2.16	2.90	4.17	10.76
SimCSE*	9.03	6.64	5.22	4.99	9.89	12.34	15.93	30.11
ICT*	10.76	10.05	8.12	6.12	14.32	18.05	23.01	39.77
MPNet*	13.75	11.93	9.58	8.14	17.77	22.21	28.11	45.91
MACLR*	16.31	13.53	10.78	9.71	20.39	25.37	32.05	53.83
Msmarco-distilbert	14.93	12.65	10.08	8.99	19.25	23.99	30.19	50.25
RTS*	18.64	15.14	12.07	10.86	22.68	28.29	35.47	57.30
LMTX	19.09	14.19	11.12	11.49	21.76	26.59	32.92	54.15
					edia-500			
GloVe*	2.19	1.52	1.23	0.85	1.66	2.18	3.10	8.52
SentBERT*	0.17	0.15	0.13	0.05		0.18	0.30	1.29
SimCSE*	14.32	6.84	4.55	4.24	8.03	11.26	14.35	27.68
ICT*	17.74	9.67	7.06	7.35	11.60	13.84	17.19	31.08
MPNet*	22.46	12.87	9.49	8.74	14.07	16.76	20.64	34.72
MACLR*	28.44	17.75	13.53	10.40	18.16	22.38	28.52	50.09
RTS*	30.67	19.03	14.34	10.58	18.48	22.51	28.23	48.00
Msmarco-distilbert	21.62	12.75	9.52	8.27	13.81	16.68	20.89	37.60
Wisinarco-distribut								50.14

Table 2: Comparison of LMTX model with state-of-the-art EZ-XMC methods. The results with * are taken from [12], others are implemented with respective code. We only check the significance test for EURLex-4k, Wiki10-31K and AmazonCat-13k, the symbol \dagger indicates a statistically significant improvement over the best baseline model (paired t-test with $p \leq 0.01$) and the symbol \S represents the best baseline model.

Dataset	LLM Methods	P@1	P@3	P@5	R@1	R@3	R@5
EURLex-4K	ICXML (w/o reranking)	3.2	2.33	2.56	0.76	1.64	2.87
	LMTX	48.8	29.6	22.44	10.05	17.98	22.68
LF-WikiSeeAlso-320K	ICXML (w/o reranking)	12.00	9.07	6.44	7.82	16.28	18.05
LF-WIKISEEAISO-320K	LMTX	22.6	15.8	11.88	13.82	24.29	28.94

Table 3: Comparison with LLM baseline on 500 samples

- SimCSE [40]: An unsupervised sentence similarity method which adopt the dropout as augmentation for the positives.
- MPNet [41]: Pre-trained language model for paraphrase embedding.
- Msmarco-distilbert [42]: Another siamese network which is pre-trained on the MS MARCO dataset. Our
 models are initialized with this model.
- MACLR [11]: The multi-stage method attempts to create pseudo positive labels by utilizing clustering and (content, title) pairs. The model is applicable in both zero-shot and few-shot settings. We benchmark our results against the zero-shot MACLR models.
- RTS [12]: The approach creates positive labels by randomly dividing the documents and using segments from within the same documents.
- ICXML [22]: The method utilizes the LLM model to address the EZ-XMC problem during inference.

To assess the baseline performance of LF-WikiSeeAlso-320K and LF-Wikipedia-500k, we obtained the results from [12]. As for the other baselines, we acquired their performance by executing the respective baseline. The MACLR [11] and ICT [38] require documents to have titles, so we could not run these on the EURLex-4k and Wiki10-31K datasets. For a fair and feasible comparison with the baseline based on the LLM model, we replaced GPT-3.5 in ICXML [22] with WizardLM-13B-V1.0. Due to the high cost of LLM, we limited the evaluation to 500 samples. It may be noted that despite using an LLM in large label space setting, DDR [23] is applicable to few-shot learning problems, and hence is not directly comparable to the EZ-XMC methods above.

2.9 Results

Comparison with standard baselines: In Table 2, we present a comparative analysis of our model's performance against other models. Notably, our LMTX model demonstrates substantial improvements in both Precision@m & Recall@m, especially for datasets like EURLex-4k, Wiki10-31k, AmazonCat-13k, and LF-Wikipedia-500k. Particularly striking are the results in LF-Wikipedia-500k and AmazonCat-13K, where our model shows an increase of 31% and 38%, respectively, for P@1.

A strength of our approach is its ability to efficient employ LLM in the learning pipeline which can assess the relevance of a document to the short-listed label. This ability helps overcome situations where positive labels created using previous methods might not precisely match the intended task. Moreover, it tackles the potential issue of relevance in randomly generated pairs. For example, consider the labels in datasets like AmazonCat-13k and LF-Wikipedia-500k. These labels indicate the category of the Wikipedia page and the

Dataset	Models	Training	Inference
	MACLR	28.86	0.38
AmazonCat-13K	RTS	35.60	0.73
	LMTX 30k	23.00	0.29
	MACLR	28.88	0.39
LF-WikiSeeAlso-320K	RTS	26.66	1.09
	LMTX 30k	30.25	0.21

Table 4: The training and inference time (hours) of LMTX

product description, the recent methods [11; 12] create the pseudo labels by using the segments or title from the instance, but the true labels' text may not always be found within the document. Nevertheless, we can extract the positive label straight from the predefined label set with LMTX instead of relying on segments or titles. This strategy mitigates the challenge of aligning the training dataset with the desired task.

Furthermore, the results for LF-wikiseeAlso-320k are also comparable to those of state-of-the-art methods. It is important to note that the LF-wikiseeAlso-320k task slightly deviates from the usual tagging task. Instead, it involves identifying related Wikipedia titles for the current Wikipedia page. Despite these difficulties, our LMTX model adeptly handles this task, underscoring the robustness and versatility of our approach. These outcomes emphatically affirm that our proposed method is not only efficient but also effective in zero-shot scenarios for tackling diverse tagging and categorization tasks.

Comparison with LLM-based baseline: We also compared our results with an approach ICXML[22] based on LLM model. For a fair comparison, we used the identical LLM model for ICXML and focused solely on the retrieval stage, as our model operates only in this phase. As shown in Table 3, our model outperforms the LLM-based approach even without using the LLM during inference, making our method more efficient and cost-effective. Our approach enables the use of small open-source models for zero-shot XMC.

Dataset	LLM models	P@1	P@3	P@5	R@1	R@3	R@5
	Wizardlm	26.18	17.45	13.39	5.52	10.95	13.83
AmazonCat-13K	Vicuna	25.01	16.51	12.70	5.24	10.23	12.95
	Llama2	25.21	16.44	12.76	5.34	10.27	13.22
	Wizardlm	19.09	14.19	11.12	11.49	21.76	26.59
LF-WikiSeeAlso-320K	Vicuna	17.76	14.03	11.07	10.91	21.64	26.58
	Llama2	17.59	13.52	10.64	10.46	20.54	25.27
	Wizardlm	40.27	23.12	16.89	13.57	22.25	26.30
LF-Wikipedia-500K	Vicuna	39.37	22.80	16.78	13.47	21.88	26.04
	Llama2	41.67	23.53	17.20	14.37	22.74	26.86

Table 5: Comparison of using different LLM models as the teachers

Training time: In Table 4, we present the training time for our model when trained with a subset of the training set. The table shows that LMTX's time efficiency is competitive or even superior compared to other models, especially in the context of larger datasets. These results underscore the effectiveness of LMTX, even with the incorporation of the LLM model. Our method's quick inference and retrieval are further supported by the inference time in Table 4.

2.10 Ablations

In this section, we present a series of analysis studies conducted to evaluate various components of our methodology. These experiments are designed to assess: (i) the impact of different Large Language Models (LLMs) chosen as the teacher model, (ii) the impact of training sample size on model performance across various datasets, and (iii) evaluation of robustness of our proposed approach against model initialization bias, (iv) the impact of hard negatives.

LMTX efficacy with diverse open-source LLMs: Several open-source LLM models are available for consideration as potential teachers. We've also detailed the performance results using different recently released LLM models. The outcomes of utilizing different teachers to evaluate the shortlist are presented in Table 5. Notably, all of these LLM models have the same parameters size of 13 billion. In the realm of zero-shot extreme text classification, Llama2 outperform WizardLM-13B-V1.0 on certain datasets. Importantly, the results emphasize that our proposed method is not restricted to a particular LLM model. This adaptability allows us to choose a more suitable teacher for achieving enhanced performance.

Impact of training sample size on model performance and training time: To make training more efficient and cost-effective with the inclusion of the LLM, especially for large datasets, we reduce the number of documents used to train the bi-encoder. We randomly select training data from the entire dataset. The performance and training time with varying training sample size is depicted in Figure 4. Better performance is observed with more

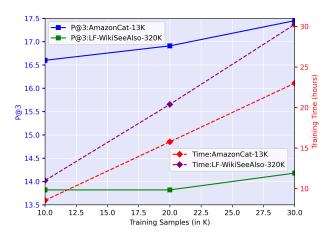


Figure 4: Evaluating training sample size effects of LMTX on efficacy and training time

documents, but it's worth noting that the training time also significantly increases as the number of training data points grows.

Evaluating the robustness of our method against initialization model bias: The selection of the initial model has an impact on both the quality of the initial labels shortlist in the first iteration and the training process of the bi-encoder model. To eliminate the advantages of the initialization, we adopted the identical initialization method employed in this paper for the baseline RTS. The outcomes of our experiments are presented in the table 6. These results indicate that our method performs better even when the baseline utilizes the same initialization.

Dataset	Models with different initialization	P@1	P@3	P@5	R@1	R@3	R@5
AmazonCat-13K	RTS	18.89	13.59	11.07	3.69	8.03	10.97
	RTS (same initialization)	17.87	12.78	10.35	3.57	7.81	10.61
	LMTX	26.18	17.45	13.39	5.52	10.95	13.83
LF-WikiSeeAlso-320K	RTS	18.64	15.14	12.07	10.86	22.68	28.29
	RTS (same initialization)	14.82	11.19	8.89	8.41	16.81	21.02
	LMTX	19.09	14.19	11.12	11.49	21.76	26.59

Table 6: Comparison of employing same initialization

Impact of different negative sampling strategies: The training of the bi-encoder uses in-batch negatives. We also investigated hard negatives, which are provided by the LLM model and marked with a tag "no". For each document, the negatives consist of the hard negatives with "no" labels and the pseudo-positive labels of other documents in the same batch. We present the results in Table 7. The results indicate that using hard negatives can hinder the training of the bi-encoder because these hard negatives might actually be false negatives.

Dataset	Negative Sampling Strategies	P@1	P@3	P@5	R@1	R@3	R@5
AmazonCat-13K	In-batch negatives	26.18	17.45	13.39	5.52	10.95	13.83
AmazonCat-13K	Hard negatives + in-batch negatives	25.21	16.44	12.76	5.34	10.27	13.22
LF-WikiSeeAlso-320K In-batch negatives Hard negatives + in-batch ne	In-batch negatives	19.09	14.19	11.12	11.49	21.76	26.59
	Hard negatives + in-batch negatives	12.88	9.14	7.33	7.15	13.53	17.06

Table 7: Comparison of different negative strategies

3 Related Work

Supervised extreme multi-label text classification: In the realm of supervised XMC, various algorithms have been proposed, which can be mainly categorized into three distinct groups. The first class of methods follow a one-vs-rest approach [43; 3; 44; 45] based on TF-IDF representation, achieving significant improvements over some of the earliest methods that relied on decision trees and random forests [46; 47]. The second category consists of tree-based methods [48; 49; 50; 4; 51; 52; 53; 54; 55]. These methods involve training distinct classifiers for different levels of the tree. Clustering algorithms are employed to group labels for these levels, and the text of labels is often not required for these algorithms. During inference, labels are filtered at various levels, leading to accelerated inference with logarithmic time complexity. The state-of-the-art performance [54; 53; 4] is based on a transformer encoder and multi-layered tree classifiers. On the other hand, for short-text inputs, it has been shown that state-of-the-art performance can be achieved by using a convolution architecture along the embedding dimension [56]. The third category encompasses embedding-based methods [57; 58; 59; 27; 60; 61; 62]. These methods involve training dense embedding encoders for documents and learning dense label embeddings through multi-stage modules and negative sampling. The labels text is often available and used for learning the labels embeddings. The predictions utilize ANNs [26] or trees to expedite the inference process in sub-linear time. For short-text input instances, when labels are endowed with textual features, a recent data augmentation has been proposed in a recent work [63]. While speeding up training and prediction stages has been a design goal for most algorithms in supervised XMC, recent works have also begun to focus on memory efficient computations for training on commodity GPUs [64]. All of these supervised approaches rely on well-annotated datasets and require comprehensive coverage of most of the labels in the training dataset. However, these traditional supervised approaches are unable to handle unseen labels effectively.

Zero-shot extreme multi-label text classification

The zero-shot XMC means that the the model is capable of handling unseen labels which are not in the training dataset. ZestXML [10] was the first paper that attempted to address the issue of unseen labels. They generated TF-IDF features for both documents and labels, then trained a linear model to project documents into the labels' TF-IDF space, enabling retrieval of unseen labels using their TF-IDF features. However, this method still relies on well-annotated training datasets to learn the linear model and is not suitable for the cold start scenario where no annotated data is available. Another extreme setting in zero-shot XMC is Extreme Zero-shot Extreme Multi-label Text Classification (EZ-XMC) [11]. EZ-XMC is specifically designed for the zero-shot scenario, particularly tailored for the cold start scenario without the need for a well-annotated training dataset. The key distinction between zero-shot XMC and EZ-XMC lies in whether annotated labels are employed in the training process. Unlike zero-shot XMC, EZ-XMC does not utilize any annotated labels. We adopt the EZ-XMC setting in this paper. The existing EZ-XMC models mainly concentrate on

creating effective representation for documents and labels by employing a bi-encoder. When searching for a document, the labels can be efficiently retrieved through the MIPS [24] constructed over label embeddings. The bi-encoder is trained with pseudo positive labels constructed from the document structure or metadata. For example, MACLR [11] proposes a multi-stage self-supervised approach using pseudo pairs of (title, document). On the other hand, RTS [12] introduces a randomized text segmentation method to construct pseudo positive labels with segments within one document. Additionally, [65] introduced a metadata-induced contrastive learning method for training the bi-encoder.

Dense sentence embedding

In the domains of open domain question answering and information retrieval, several unsupervised sentence embedding methods have recently emerged. The key aspect of these unsupervised methods lies in constructing pseudo positive and negative passages. For instance, ICT (Inverse Cloze Task) [38] constructs positive passages by extracting random sentences and their corresponding contexts from the documents. MSS[66] shows that the ICT encoder can be improved by predicting the masked salient spans with a reader. Spider [67] adopts sentences that contain recurring spans as positive passage. Both HLP [68] and WLP[69] utilize hyperlinks within Wikipedia pages to construct positive passages. Another direction is to assign relevance score for the retrieved top k passages, ART [20] tries to guide the training of bi-encoder via the question reconstruction score. Additionally, there are works that focus on sentence similarity, including (i) SimCSE [40] introduces a contrastive learning framework that employs dropout noise as augmented positives, and (ii) Sentence-BERT [39] introduces a fine-tuned siamese transformer sentence embedding framework that can be utilized for training models on various downstream tasks as well.

Large language models applications

LLM models such as GPT-3 [70], ChatGPT, and GPT-4 [34] have demonstrated their zero-shot effectiveness in various NLP downstream tasks. Certain fields, such as retrieval and recommendation, have begun exploring the applications of LLM. The directions for utilizing LLM can be broadly classified into two categories. The first category involves employing LLM for augmentation and generating training pairs for smaller dense retrieval models [13; 71; 14]. On the other hand, the second category focuses on harnessing the reasoning abilities in relevance matching for pairs, mimicking human-like reasoning for ranking items or passages [15; 16; 17; 18; 19; 21]. Some researchers have delved into utilizing large language models in XMC. In the work [23], the LLM is employed to construct a thesaurus for labels in a few-shot setting. ICXML [22], on the other hand, directly applied the LLM for inference in EZ-XMC setting. This approach predominantly focuses on recommendation datasets and relies on the costly GPT-3.5 and GPT-4 for inference. In contrast, our methods concentrate on tagging tasks and emphasize swift inference through a lightweight bi-encoder.

4 Conclusion

This paper presents an innovative approach aimed at tackling the EZ-XMC tagging and categorization problem. We make use of LLM model as a teacher to facilitate the training of the bi-encoder. In contrast to existing methods, our novel approach effectively addresses the challenge of discrepancies between the training pairs created and the intended task. Furthermore, our algorithm directly employs pseudo positive labels identified by LLM to train the model, eliminating the necessity to generate low-quality pairs through document segments. The outcome of our performance evaluation shows that our method attains state-of-the-art predictive performance across multiple datasets. This achievement highlights the remarkable efficiency of our proposed algorithm in resolving EZ-XMC tagging and categorization issues, while maintaining a faster inference speed during the prediction process. Moreover, ablation experiments underscore its capacity to achieve even better performance with alternative teacher models. For future work, exploring more efficient ways to integrate the LLM model is interesting, given the large model size and inference latency in many online applications.

References

- [1] K. Bhatia, K. Dahiya, H. Jain, P. Kar, A. Mittal, Y. Prabhu, and M. Varma. The extreme classification repository: Multi-label datasets and code, 2016.
- [2] Wei-Cheng Chang, Daniel Jiang, Hsiang-Fu Yu, Choon Hui Teo, Jiong Zhang, Kai Zhong, Kedarnath Kolluri, Qie Hu, Nikhil Shandilya, Vyacheslav Ievgrafov, et al. Extreme multi-label learning for semantic matching in product search. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 2643–2651, 2021.
- [3] Rohit Babbar and Bernhard Schölkopf. Dismec: Distributed sparse machines for extreme multi-label classification. In *Proceedings of the tenth ACM international conference on web search and data mining*, pages 721–729, 2017.

- [4] Ting Jiang, Deqing Wang, Leilei Sun, Huayi Yang, Zhengyang Zhao, and Fuzhen Zhuang. Lightxml: Transformer with dynamic negative sampling for high-performance extreme multi-label text classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 7987–7994, 2021.
- [5] Mohammadreza Qaraei, Erik Schultheis, Priyanshu Gupta, and Rohit Babbar. Convex surrogates for unbiased loss functions in extreme classification with missing labels. In *Proceedings of the Web Conference 2021*, pages 3711–3720, 2021.
- [6] Erik Schultheis and Rohit Babbar. Unbiased loss functions for multilabel classification with missing labels. *arXiv* preprint arXiv:2109.11282, 2021.
- [7] Erik Schultheis, Marek Wydmuch, Rohit Babbar, and Krzysztof Dembczynski. On missing labels, long-tails and propensities in extreme multi-label classification. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1547–1557, 2022.
- [8] Marek Wydmuch, Kalina Jasinska-Kobus, Rohit Babbar, and Krzysztof Dembczynski. Propensity-scored probabilistic label trees. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2252–2256, 2021.
- [9] Himanshu Jain, Yashoteja Prabhu, and Manik Varma. Extreme multi-label loss functions for recommendation, tagging, ranking & other missing label applications. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 935–944, 2016.
- [10] Nilesh Gupta, Sakina Bohra, Yashoteja Prabhu, Saurabh Purohit, and Manik Varma. Generalized zero-shot extreme multi-label learning. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 527–535, 2021.
- [11] Yuanhao Xiong, Wei-Cheng Chang, Cho-Jui Hsieh, Hsiang-Fu Yu, and Inderjit Dhillon. Extreme Zero-Shot learning for extreme text classification. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5455–5468, Seattle, United States, July 2022. Association for Computational Linguistics.
- [12] Tianyi Zhang, Zhaozhuo Xu, Tharun Medini, and Anshumali Shrivastava. Structural contrastive representation learning for zero-shot multi-label text classification. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4937–4947, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- [13] Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, and Rodrigo Nogueira. Inpars: Data augmentation for information retrieval using large language models. *arXiv preprint arXiv:2202.05144*, 2022.
- [14] Jon Saad-Falcon, Omar Khattab, Keshav Santhanam, Radu Florian, Martin Franz, Salim Roukos, Avirup Sil, Md Arafat Sultan, and Christopher Potts. Udapdr: Unsupervised domain adaptation via llm prompting and distillation of rerankers. *arXiv preprint arXiv:2303.00807*, 2023.
- [15] Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and Jimmy Lin. Zero-shot listwise document reranking with a large language model. *arXiv preprint arXiv:2305.02156*, 2023.
- [16] Zhen Qin, Rolf Jagerman, Kai Hui, Honglei Zhuang, Junru Wu, Jiaming Shen, Tianqi Liu, Jialu Liu, Donald Metzler, Xuanhui Wang, et al. Large language models are effective text rankers with pairwise ranking prompting. arXiv preprint arXiv:2306.17563, 2023.
- [17] Weiwei Sun, Lingyong Yan, Xinyu Ma, Pengjie Ren, Dawei Yin, and Zhaochun Ren. Is chatgpt good at search? investigating large language models as re-ranking agent. *arXiv preprint arXiv:2304.09542*, 2023.
- [18] Sunhao Dai, Ninglu Shao, Haiyuan Zhao, Weijie Yu, Zihua Si, Chen Xu, Zhongxiang Sun, Xiao Zhang, and Jun Xu. Uncovering chatgpt's capabilities in recommender systems. *arXiv preprint arXiv:2305.02182*, 2023.
- [19] Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian McAuley, and Wayne Xin Zhao. Large language models are zero-shot rankers for recommender systems. *arXiv preprint arXiv:2305.08845*, 2023.
- [20] Devendra Singh Sachan, Mike Lewis, Dani Yogatama, Luke Zettlemoyer, Joelle Pineau, and Manzil Zaheer. Questions are all you need to train a dense passage retriever. *Transactions of the Association for Computational Linguistics*, 11:600–616, 2023.

- [21] Devendra Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. Improving passage retrieval with zero-shot question generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3781–3797, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- [22] Yaxin Zhu and Hamed Zamani. Icxml: An in-context learning framework for zero-shot extreme multi-label classification. *arXiv preprint arXiv:2311.09649*, 2023.
- [23] Nan Xu, Fei Wang, Mingtao Dong, and Muhao Chen. Dense retrieval as indirect supervision for large-space decision making. *arXiv preprint arXiv:2310.18619*, 2023.
- [24] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547, 2019.
- [25] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv* preprint arXiv:1910.01108, 2019.
- [26] Yu A Malkov and Dmitry A Yashunin. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and machine intelligence*, 42(4):824–836, 2018.
- [27] Kunal Dahiya, Nilesh Gupta, Deepak Saini, Akshay Soni, Yajun Wang, Kushal Dave, Jian Jiao, Gururaj K, Prasenjit Dey, Amit Singh, et al. Ngame: Negative mining-aware mini-batching for extreme classification. In Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining, pages 258–266, 2023.
- [28] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [29] R. Manmatha, Chaoxia Wu, Alex Smola, and Philipp Krähenbühl. Sampling matters in deep embedding learning. 2017 IEEE International Conference on Computer Vision (ICCV), pages 2859–2867, 2017.
- [30] Sashank J. Reddi, Satyen Kale, Felix Yu, Daniel Holtmann-Rice, Jiecao Chen, and Sanjiv Kumar. Stochastic negative mining for learning with large output spaces. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 1940–1949. PMLR, 16–18 Apr 2019.
- [31] Wei-Cheng Chang, Daniel L. Jiang, Hsiang-Fu Yu, Choon-Hui Teo, Jiong Zhang, Kai Zhong, Kedarnath Kolluri, Qie Hu, Nikhil Shandilya, Vyacheslav Ievgrafov, Japinder Singh, and Inderjit S. Dhillon. Extreme multi-label learning for semantic matching in product search. In Feida Zhu, Beng Chin Ooi, and Chunyan Miao, editors, *KDD '21: The 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, Singapore, August 14-18, 2021*, pages 2643–2651. ACM, 2021.
- [32] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018.
- [33] Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions, 2023.
- [34] OpenAI. Gpt-4 technical report. ArXiv, abs/2303.08774, 2023.
- [35] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [36] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023.
- [37] Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics.

- [38] Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. Latent retrieval for weakly supervised open domain question answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6086–6096, 2019.
- [39] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, 2019.
- [40] Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 6894–6910, 2021.
- [41] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. Mpnet: Masked and permuted pre-training for language understanding. *Advances in Neural Information Processing Systems*, 33:16857–16867, 2020.
- [42] Nils Reimers and Iryna Gurevych. The curse of dense low-dimensional information retrieval for large index sizes. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 605–611, 2021.
- [43] Ian En-Hsu Yen, Xiangru Huang, Pradeep Ravikumar, Kai Zhong, and Inderjit Dhillon. Pd-sparse: A primal and dual sparse approach to extreme multiclass and multilabel classification. In *International conference on machine learning*, pages 3069–3077. PMLR, 2016.
- [44] Rohit Babbar and Bernhard Schölkopf. Data scarcity, robustness and extreme multi-label classification. *Machine Learning*, 108(8-9):1329–1351, 2019.
- [45] Erik Schultheis and Rohit Babbar. Speeding-up one-versus-all training for extreme classification via mean-separating initialization. *Machine Learning*, 111(11):3953–3976, 2022.
- [46] Yashoteja Prabhu and Manik Varma. Fastxml: A fast, accurate and stable tree-classifier for extreme multi-label learning. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 263–272, 2014.
- [47] Rahul Agrawal, Archit Gupta, Yashoteja Prabhu, and Manik Varma. Multi-label learning with millions of labels: Recommending advertiser bid phrases for web pages. In *Proceedings of the 22nd international conference on World Wide Web*, pages 13–24, 2013.
- [48] Ronghui You, Zihan Zhang, Ziye Wang, Suyang Dai, Hiroshi Mamitsuka, and Shanfeng Zhu. Attentionxml: Label tree-based attention-aware deep model for high-performance extreme multi-label text classification. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [49] Hsiang-Fu Yu, Kai Zhong, Jiong Zhang, Wei-Cheng Chang, and Inderjit S Dhillon. Pecos: Prediction for enormous and correlated output spaces. *the Journal of machine Learning research*, 23(1):4233–4264, 2022.
- [50] Wei-Cheng Chang, Hsiang-Fu Yu, Kai Zhong, Yiming Yang, and Inderjit S Dhillon. Taming pretrained transformers for extreme multi-label text classification. In *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 3163–3171, 2020.
- [51] Xuanqing Liu, Wei-Cheng Chang, Hsiang-Fu Yu, Cho-Jui Hsieh, and Inderjit Dhillon. Label disentanglement in partition-based extreme multilabel classification. *Advances in Neural Information Processing Systems*, 34:15359–15369, 2021.
- [52] Nilesh Gupta, Patrick Chen, Hsiang-Fu Yu, Cho-Jui Hsieh, and Inderjit Dhillon. Elias: End-to-end learning to index and search in large output spaces. Advances in Neural Information Processing Systems, 35:19798–19809, 2022.
- [53] Jiong Zhang, Wei-Cheng Chang, Hsiang-Fu Yu, and Inderjit Dhillon. Fast multi-resolution transformer fine-tuning for extreme multi-label text classification. Advances in Neural Information Processing Systems, 34:7267–7280, 2021.
- [54] Siddhant Kharbanda, Atmadeep Banerjee, Erik Schultheis, and Rohit Babbar. Cascadexml: Rethinking transformers for end-to-end multi-resolution training in extreme multi-label classification. *Advances in Neural Information Processing Systems*, 35:2074–2087, 2022.

- [55] Sujay Khandagale, Han Xiao, and Rohit Babbar. Bonsai: diverse and shallow trees for extreme multi-label classification. *Machine Learning*, 109:2099–2119, 2020.
- [56] Siddhant Kharbanda, Atmadeep Banerjee, Devaansh Gupta, Akash Palrecha, and Rohit Babbar. Inceptionxml: A lightweight framework with synchronized negative sampling for short text extreme classification. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 760–769, 2023.
- [57] Deepak Saini, Arnav Kumar Jain, Kushal Dave, Jian Jiao, Amit Singh, Ruofei Zhang, and Manik Varma. Galaxc: Graph neural networks with labelwise attention for extreme classification. In *Proceedings of the Web Conference* 2021, pages 3733–3744, 2021.
- [58] Kunal Dahiya, Ananye Agarwal, Deepak Saini, K Gururaj, Jian Jiao, Amit Singh, Sumeet Agarwal, Purushottam Kar, and Manik Varma. Siamesexml: Siamese networks meet extreme classifiers with 100m labels. In *International Conference on Machine Learning*, pages 2330–2340. PMLR, 2021.
- [59] Anshul Mittal, Noveen Sachdeva, Sheshansh Agrawal, Sumeet Agarwal, Purushottam Kar, and Manik Varma. Eclare: Extreme classification with label graph correlations. In *Proceedings of the Web Conference 2021*, pages 3721–3732, 2021.
- [60] Anshul Mittal, Kunal Dahiya, Sheshansh Agrawal, Deepak Saini, Sumeet Agarwal, Purushottam Kar, and Manik Varma. Decaf: Deep extreme classification with label features. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, pages 49–57, 2021.
- [61] Nilesh Gupta, Fnu Devvrit, Ankit Singh Rawat, Srinadh Bhojanapalli, Prateek Jain, and Inderjit S Dhillon. Efficacy of dual-encoders for extreme multi-label classification. In *The Twelfth International Conference on Learning Representations*, 2023.
- [62] Siddhant Kharbanda, Devaansh Gupta, Pankaj Malhotra, Cho-Jui Hsieh, Rohit Babbar, et al. Unidec: Unified dual encoder and classifier training for extreme multi-label classification. *arXiv preprint arXiv:2405.03714*, 2024.
- [63] Siddhant Kharbanda, Devaansh Gupta, Erik Schultheis, Atmadeep Banerjee, Cho-Jui Hsieh, and Rohit Babbar. Learning label-label correlations in extreme multi-label classification via label features. *arXiv preprint arXiv:2405.04545*, 2024.
- [64] Erik Schultheis and Rohit Babbar. Towards memory-efficient training for extremely large output spaces—learning with 670k labels on a single commodity gpu. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 689–704. Springer, 2023.
- [65] Yu Zhang, Zhihong Shen, Chieh-Han Wu, Boya Xie, Junheng Hao, Ye-Yi Wang, Kuansan Wang, and Jiawei Han. Metadata-induced contrastive learning for zero-shot multi-label text classification. In *Proceedings of the ACM Web Conference* 2022, pages 3162–3173, 2022.
- [66] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. Retrieval augmented language model pre-training. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference* on Machine Learning, volume 119 of Proceedings of Machine Learning Research, pages 3929–3938. PMLR, 13–18 Jul 2020.
- [67] Ori Ram, Gal Shachaf, Omer Levy, Jonathan Berant, and Amir Globerson. Learning to retrieve passages without supervision. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2687–2700, 2022.
- [68] Jiawei Zhou, Xiaoguang Li, Lifeng Shang, Lan Luo, Ke Zhan, Enrui Hu, Xinyu Zhang, Hao Jiang, Zhao Cao, Fan Yu, et al. Hyperlink-induced pre-training for passage retrieval in open-domain question answering. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7135–7146, 2022.
- [69] Wei-Cheng Chang, X Yu Felix, Yin-Wen Chang, Yiming Yang, and Sanjiv Kumar. Pre-training tasks for embedding-based large-scale retrieval. In *International Conference on Learning Representations*, 2019.
- [70] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

[71] Vitor Jeronymo, Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, Roberto Lotufo, Jakub Zavrel, and Rodrigo Nogueira. Inpars-v2: Large language models as efficient dataset generators for information retrieval. *arXiv* preprint arXiv:2301.01820, 2023.

5 Appendix

5.1 Prompts for LLM

- EURLex-4k and Wiki10-31K: "document = {doc}. Is the tag {label_text} relevant to the document? answer yes or no"
- AmazonCat-13K: "document = {doc}. The document is amazon product description, Is the tag {label_text} relevant to the document? answer yes or no"
- **LF-WikiSeeAlso-320K**: "document = {doc}. The document is the wikipedia page. Does another wikipedia page name "{label_text}" has the relation to the document? answer yes or no"
- **LF-Wikipedia-500K**:"document = {doc}, the document is the wikipedia page. Is the tag "{label_text}" relevant to the document? answer yes or no".