

# The Effect of Metadata on Scientific Literature Tagging: A Cross-Field Cross-Model Study

Yu Zhang, Bowen Jin, Qi Zhu, Yu Meng, Jiawei Han  
University of Illinois at Urbana-Champaign  
{yuz9, bowenj4, qiz3, yumeng5, hanj}@illinois.edu

## ABSTRACT

Due to the exponential growth of scientific publications on the Web, there is a pressing need to tag each paper with fine-grained topics so that researchers can track their interested fields of study rather than drowning in the whole literature. Scientific literature tagging is beyond a pure multi-label text classification task because papers on the Web are prevalently accompanied by metadata information such as venues, authors, and references, which may serve as additional signals to infer relevant tags. Although there have been studies making use of metadata in academic paper classification, their focus is often restricted to one or two scientific fields (e.g., computer science and biomedicine) and to one specific model. In this work, we systematically study the effect of metadata on scientific literature tagging across 19 fields. We select three representative multi-label classifiers (i.e., a bag-of-words model, a sequence-based model, and a pre-trained language model) and explore their performance change in scientific literature tagging when metadata are fed to the classifiers as additional features. We observe some ubiquitous patterns of metadata’s effects across all fields (e.g., venues are consistently beneficial to paper tagging in almost all cases), as well as some unique patterns in fields other than computer science and biomedicine, which are not explored in previous studies.

## CCS CONCEPTS

• Information systems → Digital libraries and archives; Data mining; World Wide Web.

## KEYWORDS

scientific literature tagging; metadata; text classification

### ACM Reference Format:

Yu Zhang, Bowen Jin, Qi Zhu, Yu Meng, Jiawei Han. 2023. The Effect of Metadata on Scientific Literature Tagging: A Cross-Field Cross-Model Study. In *Proceedings of the ACM Web Conference 2023 (WWW ’23)*, May 1–5, 2023, Austin, TX, USA. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3543507.3583354>

## 1 INTRODUCTION

A variety of academic service platforms, such as Google Scholar, AMiner [40], Microsoft Academic [38], Semantic Scholar [1], and

<sup>†</sup>Code and Datasets are available at <https://github.com/yuzhimanhua/MAPLE> and <https://doi.org/10.5281/zenodo.7611544>, respectively.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
WWW ’23, May 1–5, 2023, Austin, TX, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-1-4503-9416-1/23/04...\$15.00  
<https://doi.org/10.1145/3543507.3583354>

The screenshot shows a paper titled "Graph structure in the Web" from the 2000 The Web Conference. The page is annotated with red dashed boxes and labels on the left side. The labels include: "Title (text)" pointing to the title; "Venue (metadata)" pointing to the conference name; "Authors (metadata)" pointing to the list of authors; "Abstract (text)" pointing to the abstract paragraph; "References (metadata)" pointing to the list of references; and "Tags" pointing to a set of topic tags at the bottom. The tags include "Webgraph", "Graph database", "Link farm", "Crawling", "World Wide Web", "Computer science", "Graph", and "Sociological Phenomena".

Figure 1: A scientific paper with metadata from Microsoft Academic [38]. The goal of scientific literature tagging is to predict its related topics.

PubMed [26], are available on the Web with great attention received. One major goal of these platforms is to help researchers query and track academic information and resources. Meanwhile, the volume of scientific publications is growing exponentially, doubling every 12 years [12] and reaching 240,000,000 by 2019 [42]. In such an information explosion era, it becomes more important than ever to accurately tag each scientific paper with its relevant topics so that researchers can track their interested fields of study instead of getting overwhelmed by the whole literature. Figure 1 shows an example of the scientific literature tagging task, which aims to predict the tags such as “World Wide Web”, “Webgraph”, and “Link Farm” given the paper “Graph structure in the Web”.

Previous studies [47, 59] have pointed out that scientific literature tagging is beyond a multi-label text classification task because academic papers are accompanied by metadata, which make them more complex than plain text sequences. Such metadata information, including venues, authors, and references, can be strong indicators of each paper’s related topics. For example, in Figure 1, the venue “WWW” implies the paper’s relevance to “Computer Science” and “World Wide Web”, while the authors and references may further indicate fine-grained tags such as “Webgraph” and “Link Farm”.

Although existing studies on scientific literature tagging [45, 47, 57, 59, 60] have proposed to incorporate metadata features into the tagger, they still have two major limitations. First, their focus is restricted to one or two scientific fields only (e.g., computer science and biomedicine). Empirically examining the effect of metadata in other fields (e.g., art, economics, mathematics, physics, etc.) has remained elusive, mainly owing to the lack of benchmark datasets for fine-grained paper tagging in these fields. Second, the proposed

metadata-aware taggers mainly train an RNN [5] or Transformer [41] architecture from scratch. It is still unclear whether the proposed usage of metadata can be generalized to other approaches, such as bag-of-words classifiers (which are still broadly studied in large-scale multi-label text classification when efficiency is concerned [23]) and pre-trained language models.

**Contributions.** To address the aforementioned two limitations of previous studies, in this work, we conduct a systematic *cross-field cross-model* study on the effect of metadata on scientific literature tagging. First, we construct a large-scale scientific literature tagging benchmark, MAPLE (Metadate-Aware Paper collEction), from the Microsoft Academic Graph [38]. MAPLE covers 19 scientific fields and consists of more than 11.9 million papers. The number of candidate tags in each field ranges between  $\sim 700$  and  $\sim 64,000$ . Then, we consider three major types of multi-label classification approaches: bag-of-words classifiers [2, 16, 19, 31–34, 39, 49, 50, 54], sequence-based classifiers [21, 44, 45, 47, 53, 59], and pre-trained language models [4, 17, 48, 52, 56, 60]. We select one representative model from each of the three categories, namely Parabel [33], Transformer [44], and OAG-BERT [24], that can be modified in a straightforward way to jointly take text and metadata as input for classification. Based on MAPLE, we explore the effect of venues, authors, and references on paper tagging in the 19 fields when using the three selected classifiers. We have the following major observations:

- The effect of metadata varies significantly across different fields and classifiers. In general, venues are consistently beneficial in almost all cases, while the benefit of authors and references is highly dependent on the field and the classifier. For example, author information is evidently beneficial to paper tagging in the Philosophy field while harmful in Geology; references are helpful in Mathematics when we use Parabel but become disadvantageous in the same field when Transformer is adopted.
- The effect of metadata tends to be similar in two fields that belong to the same high-level scientific area. For example, Biology and Medicine are both life sciences (according to [35, 51]), and the effects of venues, authors, and references are largely aligned in the two fields. This finding implies that the experience of using metadata in one field can be extrapolated to a similar field.
- We also study the effect of metadata when predicting tags at different granularity levels. In a number of fields, venues improve the performance for not only coarse-grained tags but also very fine-grained ones, which may be unusual in the Computer Science field. The major reason is that some venues in these fields can indicate very fine tags (e.g., “*Journal of Roman Archaeology*” in History, “*Mediaeval Studies*” in Philosophy).

To summarize, this work makes the following contributions: (1) We construct a large-scale benchmark, MAPLE, for scientific literature tagging across 19 fields, whose field coverage is much broader than the datasets used in previous paper tagging studies [45, 47, 59]. (2) We comprehensively evaluate the performance of different types of multi-label classifiers in scientific literature tagging after incorporating metadata features. (3) Our empirical findings demonstrate some ubiquitous patterns of metadata’s effects across all fields, as well as some unique patterns in fields other than computer science and biomedicine, which are not explored in previous studies. Our systematic studies are meant for providing insights to practitioners to build scientific literature taggers that can benefit all fields.

**Table 1: Statistics of the 20 datasets in MAPLE across 19 fields. There are 2 datasets in the Computer Science field, one of which is collected from top conferences and the other from top journals.**

Field	Paper Source	#Papers	#Labels	#Venues	#Authors	#References
Art	Journal	58,373	1,990	98	54,802	115,343
Philosophy	Journal	59,296	3,758	98	36,619	198,010
Geography	Journal	73,883	3,285	98	157,423	884,632
Business	Journal	84,858	2,392	97	100,525	685,034
Sociology	Journal	90,208	1,935	98	85,793	842,561
History	Journal	113,147	2,689	99	84,529	284,739
Political Science	Journal	115,291	4,990	98	93,393	480,136
Environmental Science	Journal	123,945	694	100	265,728	1,217,268
Economics	Journal	178,670	5,205	97	135,247	1,042,253
Engineering	Journal	270,006	10,683	100	430,046	1,867,276
Psychology	Journal	372,954	7,641	100	460,123	2,313,701
Computer Science	Conference	263,393	13,613	75	331,582	1,084,440
Computer Science	Journal	410,603	15,540	96	634,506	2,751,996
Geology	Journal	431,834	7,883	100	471,216	1,753,762
Mathematics	Journal	490,551	14,271	98	404,066	2,150,584
Materials Science	Journal	1,337,731	6,802	99	1,904,549	5,457,773
Physics	Journal	1,369,983	16,664	91	1,392,070	3,641,761
Biology	Journal	1,588,778	64,267	100	2,730,547	7,086,131
Chemistry	Journal	1,849,956	35,538	100	2,721,253	8,637,438
Medicine	Journal	2,646,105	36,619	100	4,345,385	7,405,779

## 2 PRELIMINARIES

**Text and Metadata.** We represent the text information of a paper  $p$  as a single text sequence  $\mathcal{T}_p = w_1 w_2 \dots w_{|\mathcal{T}_p|}$  by concatenating its title and abstract. The metadata of a paper  $p$  is represented as a set  $\mathcal{M}_p = \{m_1, m_2, \dots, m_{|\mathcal{M}_p|}\}$  consisting of its venue, author(s), and reference(s).

**Problem Definition.** The scientific literature tagging task can be cast as a large-scale multi-label classification problem, where all candidate tags (e.g., “Organic Chemistry”, “Coupling Reaction”, “Suzuki Reaction”) constitute a large label space  $\mathcal{L}$  (e.g., with  $10^3$ – $10^5$  labels). Given a scientific paper with its text and metadata information, the task is to find a set of labels from  $\mathcal{L}$  that are relevant to the paper. Formally, the problem is defined as follows.

*Definition 2.1.* (Problem Definition) Given a training corpus  $\mathcal{D}$  and the label space  $\mathcal{L}$ , where each paper  $p \in \mathcal{D}$  is associated with its text  $\mathcal{T}_p$ , metadata information  $\mathcal{M}_p$ , and relevant tags  $\mathcal{L}_p \subseteq \mathcal{L}$ , our objective is to learn a multi-label text classifier  $f_{\text{class}}$  that can map an unlabeled paper  $p' \notin \mathcal{D}$  to its relevant tags  $\mathcal{L}_{p'} \subseteq \mathcal{L}$ .

## 3 DATASET CONSTRUCTION

Previous studies on scientific literature tagging [9, 31, 45, 47, 52, 59] mainly use computer science and biomedicine papers to evaluate their proposed model, meanwhile paying less attention to other scientific fields. To bridge this gap, we construct MAPLE, a multi-field benchmark for evaluating scientific literature tagging. MAPLE is built upon data from the Microsoft Academic Graph (MAG) [38], which has been widely adopted in scientific text mining [6, 51, 59]. MAG covers 19 academic fields, which are listed in Table 1. For each field, we conduct the following steps to construct a dataset.

**Venue Selection.** MAG maintains a list of the top-100 journals in each field according to the  $h$ -index [13]. When constructing MAPLE, we focus on papers published in these top journals. Note that some preprint services (e.g., “arXiv”, “bioRxiv”, “SSRN”, “NBER Working Paper”) [43] are also viewed as top journals in MAG, but we exclude

them from our consideration. As a result, in Table 1, some of the constructed datasets contain less than 100 venues. Among the 19 fields, computer science (CS) has a unique publication culture: CS papers often appear first or exclusively in conferences rather than journals. Thus, for the Computer Science field, besides a collection of journal papers, we construct another dataset with papers from 75 top conferences according to CSRankings<sup>1</sup>. That being said, we will construct 20 datasets in total for the 19 fields.

**Label Space Construction.** For each field, we need a set of candidate labels for paper tagging. MAG has a directed acyclic graph (DAG)-structured label taxonomy  $\mathcal{L}_{\text{MAG}}$  [37]. The taxonomy has 6 levels and more than  $10^5$  labels, where each of the 19 fields is a Layer-0 label (i.e., the most coarse-grained label). Given a field  $F$ , we extract its descendant labels in the taxonomy  $\mathcal{L}_{\text{MAG}}$  as the label space  $\mathcal{L}_F$  of this field, but we exclude  $F$  itself from  $\mathcal{L}_F$  because the root label is trivial to predict in classification. Since a label may have more than one parent in the DAG-structured taxonomy, it may appear in the label space of two or more fields. For example, the label “Anatomy” is a child of both “Biology” and “Medicine”, so it is a candidate label in both the Biology and the Medicine datasets in MAPLE.

**Paper Selection.** Each paper  $p$  in MAG is tagged with its relevant labels  $\mathcal{L}_p \subseteq \mathcal{L}_{\text{MAG}}$  [37]<sup>2</sup>. To be included in MAPLE (given a field  $F$ ), a paper  $p$  needs to satisfy the following two criteria: (1)  $p$  is published in a selected venue of  $F$ ; (2)  $p$  is labeled with  $F$  and at least one of  $F$ 's descendants (i.e.,  $F \in \mathcal{L}_p$  and  $|\mathcal{L}_p \cap \mathcal{L}_F| \geq 1$ ). When studying the scientific literature tagging task in a field  $F$ , we focus on labels related to that field only. Therefore, the ground truth labels of  $p$  are defined as  $\mathcal{L}_{p|F} = \mathcal{L}_p \cap \mathcal{L}_F$ . Note that a paper may appear in more than one field, and its ground truth labels are different when we consider different fields. For example, if a paper is tagged with “Medicine”, “Polypharmacy”, “Computer Science”, and “Graph Embedding”, given that “Polypharmacy” is a candidate label in the Medicine field and “Graph Embedding” is a candidate label in the Computer Science field, the paper will appear in both the Medicine and the Computer Science datasets in MAPLE. However, when we perform tagging in the Medicine field, the ground-truth label of the paper is “Polypharmacy”; when we perform tagging in the Computer Science field, the ground-truth label of the paper is “Graph Embedding”.

**Text and Metadata Extraction.** For each selected paper  $p$ , we extract its title, abstract, venue, author(s), and reference(s) from MAG. The title and abstract are concatenated as text information  $\mathcal{T}_p$ . The venue, author(s), and reference(s) constitute metadata features  $\mathcal{M}_p$  of the paper.

**Training-Validation-Testing Split.** MAPLE contains academic papers published between Jan. 1, 1981 and Dec. 31, 2020. We use papers from 1981 to 2015 for training and validation, and papers from 2016 to 2020 for testing.

Statistics of the constructed 20 datasets in MAPLE can be found in Table 1. More details are shown in Appendix A.1. From now on, for convenience of discussion, we use the terms “field” and “dataset” interchangeably if there is no ambiguity. (In other words, we treat

“Computer Science (Conference)” and “Computer Science (Journal)” as different fields.)

## 4 MODELS

Large-scale multi-label text classification (LMTC) has been extensively studied over the past decade. Various approaches are proposed and can be applied to scientific literature tagging. Based on how text is used as features in the classifier, existing LMTC approaches can be divided into three major categories: (1) *Bag-of-words classifiers* [2, 16, 19, 31–34, 39, 49, 50, 54] treat each document as a multiset of tokens while disregarding word position and order. Trees, embeddings, and linear layers are commonly used in these models. (2) *Sequence-based classifiers* [21, 44, 45, 47, 53, 59] take each document as a sequence of tokens and train a CNN, RNN, or Transformer architecture from scratch to build a multi-label classifier. (3) *Pre-trained language model classifiers* [4, 17, 48, 52, 56, 60] aim at transferring the knowledge learned from web-scale corpora (e.g., Wikipedia and PubMed) to the LMTC task, which can complement the text information from the training corpus. Since the goal of this paper is to study the effect of metadata, we select one LMTC model from each of the three categories that can be easily augmented with metadata information. Now we introduce the three selected models – Parabel [33] (bag-of-words), Transformer [44] (sequence-based), and OAG-BERT [24] (pre-trained language model) – and how they can take text and metadata features as input for classification.

### 4.1 Bag-of-Words Classifier: Parabel [33]

**4.1.1 Using Text Only.** In general, bag-of-words classifiers represent each document  $p$  as a  $|\mathcal{V}_{\mathcal{D}}|$ -dimensional vector  $\mathbf{x}_p$ , where  $\mathcal{V}_{\mathcal{D}}$  is the vocabulary of the training corpus  $\mathcal{D}$ . Given a word  $w \in \mathcal{V}_{\mathcal{D}}$ , its corresponding entry in  $\mathbf{x}_p$  is defined using the tf-idf score:

$$x_{p,w} = \text{tf}(w, p) \cdot \text{idf}(w, \mathcal{D}). \quad (1)$$

Here,  $\text{tf}(w, p)$  is the term frequency of  $w$  in document  $p$ ;  $\text{idf}(w, \mathcal{D}) = \log \frac{|\mathcal{D}|}{|\{p' \in \mathcal{D} : w \in \mathcal{T}_{p'}\}|}$  is the inverse document frequency of  $w$ .

The relevant labels of each document are represented by an  $|\mathcal{L}|$ -dimensional vector  $\mathbf{y}_p$ , where  $y_{p,l} = 1$  if  $l$  is a tag relevant to  $p$  (i.e.,  $l \in \mathcal{L}_p$ ), and  $y_{p,l} = 0$  otherwise.

To perform multi-label text classification, Parabel learns an ensemble of multiple label trees, each of which is obtained by recursively partitioning the labels into two balanced groups until each node contains less than a certain number of labels. The partition process is implemented by spherical 2-means clustering based on the label representation  $\mathbf{v}_l$ , which is a unit vector in the direction of the mean of the training points containing label  $l$ . After label space partitioning, Parabel learns a hierarchical discriminative classifier  $\Pr(\mathbf{y}_p | \mathbf{x}_p)$ . Specifically, at each non-leaf node, a distribution is learned to determine which child nodes should be traversed; at each leaf node, a distribution is learned to predict the set of relevant tags. For more technical details, one can refer to [33], but we omit them here as they are not directly related to the usage of metadata.

**4.1.2 Using Text + Metadata.** It is straightforward to generalize bag-of-words classifiers to take metadata features. Let  $\mathcal{U}_{\mathcal{D}}$  denote the set of metadata instances appearing in  $\mathcal{D}$ . Given a metadata instance  $m \in \mathcal{U}_{\mathcal{D}}$  and a paper  $p$ , we define  $\text{tf}(m, p)$  and  $\text{idf}(m, \mathcal{D})$  as follows:

$$\text{tf}(m, p) = \mathbf{1}(m \in \mathcal{M}_p), \quad \text{idf}(m, \mathcal{D}) = \log \frac{|\mathcal{D}|}{|\{p' \in \mathcal{D} : m \in \mathcal{M}_{p'}\}|}, \quad (2)$$

<sup>1</sup><https://csrankings.org/>

<sup>2</sup>The paper tags  $\mathcal{L}_p$  come from the predictions of a system [37] proposed by Microsoft Academic. The tags are accurate as checked by humans [37] and have been used to support notable findings [18, 51]. Meanwhile, we also conduct experiments on three datasets with MeSH labels [7], which are curated by biomedical experts. Discussions on the additional three datasets can be found in Appendix A.5.

where  $\mathbf{1}(\cdot)$  is the indicator function. One can find that the definitions in Eq. (2) are well aligned with the definitions of  $\text{tf}(w, p)$  and  $\text{idf}(w, \mathcal{D})$ , except that a metadata instance does not appear multiple times in a document.

Now we can have a “bag-of-metadata” representation  $\tilde{x}_p$  for each paper.  $\tilde{x}_p$  is a  $|\mathcal{U}_{\mathcal{D}}|$ -dimensional vector, where

$$\tilde{x}_{p,m} = \text{tf}(m, p) \cdot \text{idf}(m, \mathcal{D}). \quad (3)$$

Finally, we represent each paper  $p$  as a  $(|\mathcal{U}_{\mathcal{D}}| + |\mathcal{V}_{\mathcal{D}}|)$ -dimensional vector, that is, the concatenation of bag-of-words and “bag-of-metadata” representations  $x_p \parallel \tilde{x}_p$ . This vector is fed into Parabel for training and prediction.

## 4.2 Sequence-based Classifier: Transformer [44]

**4.2.1 Using Text Only.** To apply Transformer [41] to LMTC, we follow the architecture proposed in [44], which adds multiple “[CLS]” tokens in front of document text  $\mathcal{T}_p$  as the input sequence. Formally, given  $\mathcal{T}_p = w_1 w_2 \dots w_{|\mathcal{T}_p|}$ , the input sequence of paper  $p$  is

$$\tilde{\mathcal{I}}_p = [\text{CLS}_1] [\text{CLS}_2] \dots [\text{CLS}_C] w_1 w_2 \dots w_{|\mathcal{T}_p|} \quad (4)$$

The motivation here is that when the label space is large (e.g., with  $10^4$  tags), the output representation of one “[CLS]” token (e.g., a vector with several hundred dimensions) may not carry enough information to predict relevant labels. Therefore, multiple “[CLS]” tokens are needed to probe the semantics of text from different perspectives.

After Transformer encodes the input sequence  $\tilde{\mathcal{I}}_p$ , each token  $w \in \tilde{\mathcal{I}}_p$  will have an output representation  $h_w$ . We concatenate the representations of all “[CLS]” tokens together as the paper embedding:

$$h_p = h_{[\text{CLS}_1]} \parallel h_{[\text{CLS}_2]} \parallel \dots \parallel h_{[\text{CLS}_C]}. \quad (5)$$

$h_p$  is further fed into a fully connected layer to perform multi-label classification:

$$\hat{y}_p = \text{Sigmoid}(\mathbf{W}^T h_p + b). \quad (6)$$

Here,  $\hat{y}_p$  is an  $|\mathcal{L}|$ -dimensional vector, in which  $\hat{y}_{p,l}$  is the predicted probability that paper  $p$  is relevant to label  $l$ . The classifier is trained to minimize the following binary cross-entropy (BCE):

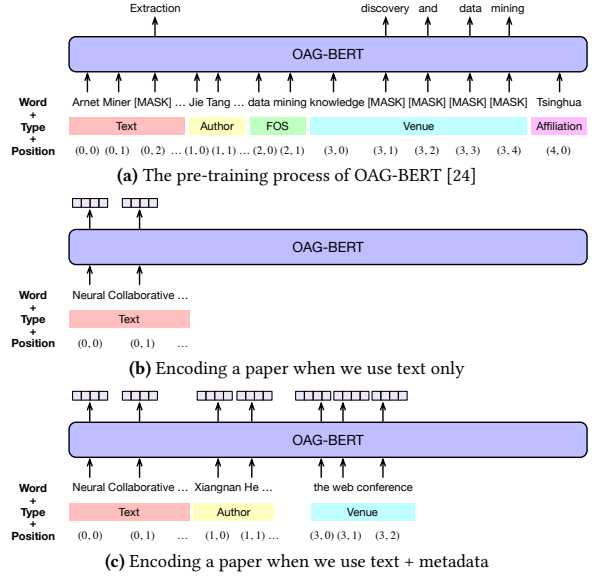
$$-\sum_{l \in \mathcal{L}} (y_{p,l} \log \hat{y}_{p,l} + (1 - y_{p,l}) \log(1 - \hat{y}_{p,l})). \quad (7)$$

**4.2.2 Using Text + Metadata.** The fully connected attention mechanism in Transformer paves an easy way to incorporate metadata. To be specific, following [59], we can directly insert metadata tokens into the input sequence. Given paper text  $\mathcal{T}_p = w_1 \dots w_{|\mathcal{T}_p|}$  and metadata  $\mathcal{M}_p = \{m_1, \dots, m_{|\mathcal{M}_p|}\}$ , the input sequence is

$$\tilde{\mathcal{I}}_p = [\text{CLS}_1] \dots [\text{CLS}_C] m_1 \dots m_{|\mathcal{M}_p|} [\text{SEP}] w_1 \dots w_{|\mathcal{T}_p|} \quad (8)$$

Unlike in CNN [21] or RNN [45, 53] classifiers, in Transformer, the order of metadata instances does not matter much because each pair of (metadata, word) or (metadata, metadata) can interact with each other via the fully connected attention mechanism during encoding. In our experiments, when listing authors in  $\tilde{\mathcal{I}}_p$ , we follow the authorship order (i.e., 1<sup>st</sup> author, 2<sup>nd</sup> author, ...); when listing references in  $\tilde{\mathcal{I}}_p$ , we adopt a random order because the citation order and inline contexts are not stored in MAG.

Following [59], we treat each metadata instance  $m_i \in \tilde{\mathcal{I}}_p$  as one token during Transformer encoding. For example, the venue “WWW” is represented as one token “[VENUE\_1135342153]” instead of its textual name “the web conference” containing multiple tokens.



**Figure 2: The pre-training process and our usage of OAG-BERT [24].**

After  $\tilde{\mathcal{I}}_p$  is fed to Transformer, the remaining designs of the metadata-aware classifier exactly follow Eqs. (5)-(7).

## 4.3 Pre-trained Language Model Classifier: OAG-BERT [24]

Previous classifiers built upon pre-trained language models (PLMs) mainly utilize BERT [10], BioBERT [20], SciBERT [3], XLNet [46], or RoBERTa [25] to derive document representations. However, these PLMs mostly focus on text information during pre-training and are not specifically designed to deal with metadata. To bridge this gap, we adopt OAG-BERT [24], an entity-augmented academic language model, which can jointly encode scientific text and venue/author information.

The pre-training process of OAG-BERT is briefly illustrated in Figure 2(a). It places text and metadata information in a single sequence for masked language modeling. Besides, it proposes three strategies to deal with metadata entities: (1) heterogeneous entity type embedding makes the model aware of different metadata types; (2) span-aware entity masking selects a continuous span within long entities (e.g., the venue “knowledge discovery and data mining”); (3) 2-dimensional positional embedding jointly models inter and intra-entity token orders. The model is first trained on 5 million full-text papers and then on 120 million paper titles/abstracts with metadata from the Open Academic Graph [55].

**4.3.1 Using Text Only.** When considering paper text  $\mathcal{T}_p$  only for classification, we first use OAG-BERT to encode the text sequence (as illustrated in Figure 2(b)).

$$H_p = \text{OAG-BERT}(\mathcal{T}_p). \quad (9)$$

Here,  $H_p = [h_{p,1}, \dots, h_{p,N}]$  contains the output vectors of all tokens. We then adopt mean pooling to obtain the paper representation  $x_p$ .

$$x_p = \frac{1}{N} \sum_{i=1}^N h_{p,i}. \quad (10)$$

Directly fine-tuning the PLM using a BCE loss (i.e., Eq. (7)) is very time-consuming, especially when the label space is large. Considering efficiency and model simplicity, we fix the paper representation  $\mathbf{x}_p$  to train a Parabel classifier  $f_{\text{Parabel}}(\mathbf{x}_p) = \Pr(\mathbf{y}_p|\mathbf{x}_p)$ .

During the prediction stage, given a paper  $p'$ , the trained classifier predicts the probability that  $p'$  is relevant to each label  $l$ :

$$\hat{y}_{p',l} = f_{\text{Parabel}}(\mathbf{x}_{p'}). \quad (11)$$

Then, we can rank the labels according to  $\hat{y}_{p',l}$  to predict a list of the most relevant labels. However, in practice, we find this strategy does not yield competitive prediction accuracy. This finding is consistent with the design in related studies [4, 31] that discrete lexical features should be considered together with continuous representations during classification. To add lexical features, motivated by [60], we propose a simple heuristic to re-rank the labels: Given a paper  $p'$ , all labels whose name appears in the paper text  $\mathcal{T}_{p'}$  will be ranked higher than those not appearing in  $\mathcal{T}_{p'}$ . This heuristic is equivalent to ranking the labels according to a modified score  $\hat{z}_{p',l} = \hat{y}_{p',l} + 1(t_l \in \mathcal{T}_{p'})$ , where  $t_l$  stands for the textual name of label  $l$ .

**4.3.2 Using Text + Metadata.** When metadata information is available, we use OAG-BERT to jointly encode paper text and metadata (as illustrated in Figure 2(c)).

$$\tilde{\mathbf{H}}_p = \text{OAG-BERT}(\mathcal{T}_p, \mathcal{M}_p). \quad (12)$$

The remaining steps exactly follow the text-only model. It is worth noting that since references are not involved during the pre-training of OAG-BERT, we can only use venues and authors as metadata here. Also, unlike the Transformer classifier that views each venue/author as one token, OAG-BERT tokenizes the textual name of venues/authors based on its vocabulary during encoding.

## 5 EXPERIMENTS

### 5.1 Setup

**Datasets and Compared Methods.** We have introduced the 20 datasets in Section 3. For each of the three classifiers, we test its performance when using text only, text+venue, text+author, and text+reference in order to check the effect of each metadata type separately. (Recall that OAG-BERT cannot take references as metadata features, so it does not have the text+reference variant.) For detailed hyperparameter settings, one can refer to Appendix A.2.

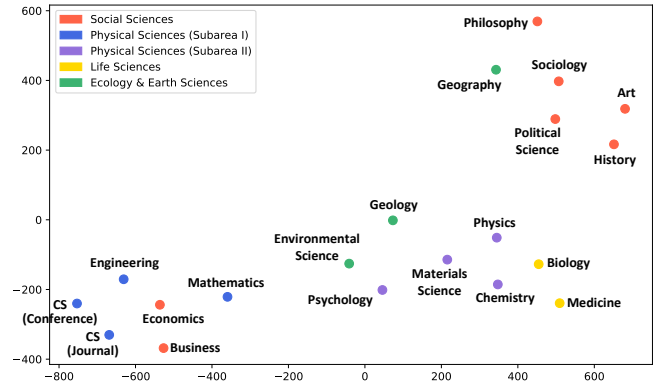
**Evaluation Metrics.** Following previous studies on multi-label text classification [21, 44, 59], we adopt  $P@k$  and  $\text{NDCG}@k$  as evaluation metrics, where  $k = 1, 3, 5$ . Given a paper  $p$ , let  $\text{rank}(i)$  be the index of the  $i$ -th highest predicted label according to each classifier, then

$$P@k = \frac{1}{k} \sum_{i=1}^k y_{p,\text{rank}(i)}. \quad (13)$$

$$\text{DCG}@k = \sum_{i=1}^k \frac{y_{p,\text{rank}(i)}}{\log(i+1)}, \quad \text{NDCG}@k = \frac{\text{DCG}@k}{\sum_{i=1}^{\min(k, \|\mathbf{y}_p\|_0)} \frac{1}{\log(i+1)}}.$$

### 5.2 Overall Analysis

Table 2 shows the  $P@k$  and  $\text{NDCG}@k$  scores of the three classifiers on the 20 datasets. We run each experiment 5 times with the average score reported. We conduct two-tailed t-tests to check the statistical significance of metadata’s effect. To be specific, given a field and a classifier, if a score is significantly improved (with  $p$ -value  $< 0.05$ ) after using a certain type of metadata in comparison with using text



**Figure 3:** We represent each field with a 24-dimensional vector based on the effect of venue, author, and reference information on the three classifiers. Then, we apply t-SNE [27] to visualize these fields in a 2-dimensional space. The color scheme highlights several high-level scientific areas, following the major clusters of science detected by [35, 51] and suggesting similar effects of metadata within each area.

only, we mark the score as **blue** in Table 2; if a score significantly deteriorates (with  $p$ -value  $< 0.05$ ), we mark the score as **red**.

From Table 2, we can observe that: (1) The effect of metadata varies remarkably across different fields. For instance, author information is significantly helpful in Computer Science and Engineering when Parabel is utilized as the classifier but becomes harmful in Biology and Medicine when the same classifier is adopted. References are useful in Business and Economics when Parabel is employed but become disadvantageous in the same fields if Transformer is the tagger. (2) Venues are consistently beneficial to scientific literature tagging in almost all cases. To be specific, all 20 fields significantly benefit from venues<sup>3</sup> when Parabel is used, 18 fields when Transformer is used, and 18 fields when OAG-BERT is used. In the remaining 4 cases where venue information is not beneficial, neither is it harmful. (3) Authors are helpful in a majority of (15 and 17, respectively) fields when Parabel and OAG-BERT are the taggers but rarely help when Transformer is adopted. References are useful in even fewer cases. The overall performance of the three types of metadata is reflected by the macro average of  $P@k$  and  $\text{NDCG}@k$  scores over the 20 datasets, which are shown in Table 3. The reason why authors and references do not work in the Transformer classifier is that their embeddings need to be trained from scratch without good initialization. In contrast, leveraging venues only incurs a small number of additional parameters. Therefore, if one aims to utilize author and reference information, some embedding pre-training techniques [59] may help.

### 5.3 Effect in Different Fields

In this section, we examine the effect of metadata in different fields. Given a field  $F$ , to describe the effect of metadata in  $F$ , we construct a 24-dimensional vector in the following way: There are three classifiers (i.e., Parabel, Transformer, and OAG-BERT) and three types of metadata (i.e., Venue, Author, and Reference) studied

<sup>3</sup>We say one field *significantly benefits from* one type of metadata (when using a certain classifier) if at least one of the five metrics is significantly improved (with  $p$ -value  $< 0.05$ ) after incorporating that type of metadata.



**Table 3: Macro average of P@k and NDCG@k over the 20 datasets. Blue, Red, and “-”: the same meaning as in Table 2.**

	Input	Parabel [33]					Transformer [44]					OAG-BERT [24]				
		P@1	P@3	P@5	N@3	N@5	P@1	P@3	P@5	N@3	N@5	P@1	P@3	P@5	N@3	N@5
Macro Average	Text	0.7513	0.5811	0.4678	0.7076	0.6977	0.7510	0.5673	0.4507	0.6896	0.6712	0.6983	0.5354	0.4275	0.6549	0.6429
	+Venue	0.7554	0.5858	0.4717	0.7130	0.7033	0.7599	0.5753	0.4572	0.6995	0.6812	0.7004	0.5391	0.4318	0.6588	0.6480
	+Author	0.7512	0.5817	0.4687	0.7079	0.6983	0.7442	0.5594	0.4433	0.6809	0.6617	0.6984	0.5374	0.4305	0.6569	0.6461
	+Reference	0.7487	0.5809	0.4689	0.7065	0.6977	0.7277	0.5469	0.4328	0.6652	0.6459	-	-	-	-	-

in our experiments, so there are  $3 \times 3 - 1 = 8$  (classifier, metadata) combinations, since (OAG-BERT, Reference) is not applicable. For each of the 8 (classifier, metadata) combinations, we calculate the relative performance change of P@1, P@3, and P@5 by comparing the classifier using text only and the classifier using text together with the metadata. For example, given the (Parabel, Venue) combination, we calculate the following 3 values:

$$\frac{P@k(\text{Parabel, Text + Venue}) - P@k(\text{Parabel, Text})}{P@k(\text{Parabel, Text})}, \quad k = 1, 3, 5. \quad (14)$$

In total, we will have  $3 \times 8 = 24$  values, which can form a 24-dimensional vector  $u_F$ . We compute  $u_F$  for all 20 fields and then apply t-SNE [27] to visualize these vectors in a 2-dimensional space.

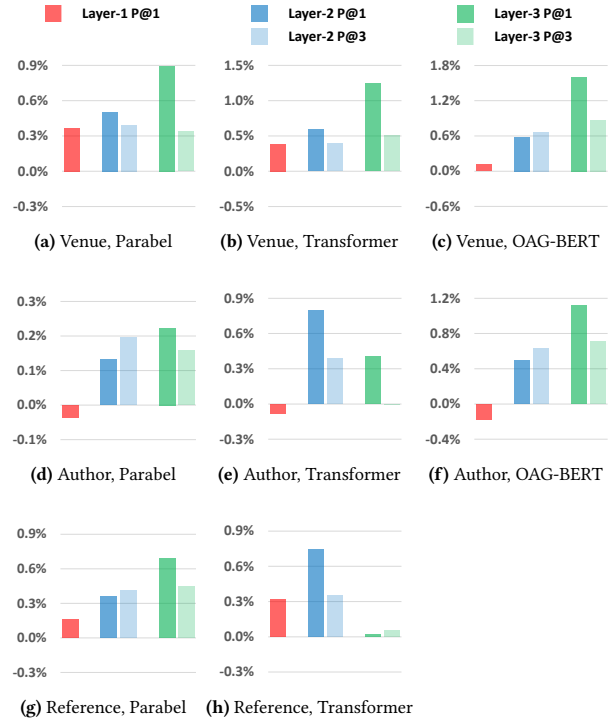
The visualization result is shown in Figure 3, where we color each field according to high-level scientific areas/subareas detected by [35, 51]. To be specific, Rosvall and Bergstrom [35] cluster scientific fields into four major areas – social sciences, physical sciences, life sciences, and ecology & earth sciences – based on a large-scale paper citation network. The cluster of physical sciences is further split into two subareas, one of which is related to mathematics and computer science, and the other to physics and chemistry; Yin et al. [51] observe a similar partitioning of the 19 fields when studying public uses of scientific papers in each field.

In Figure 3, we find that fields with the same color are often embedded closer, indicating that metadata have similar effects in fields belonging to the same area/subarea. This finding is very useful when we need to extrapolate the experience of using metadata in one field to a similar field. For example, one may have known that references are beneficial and authors are harmful when using Parabel in the Medicine field. Based on our finding, the same pattern can be deduced when using Parabel in the Biology field because Medicine and Biology are both life sciences. According to Table 2, this deduction is correct.

Meanwhile, we observe two exceptions: (1) The area of social sciences is split into two clusters, one of which contains Art, History, Philosophy, Sociology, and Political Science, and the other contains Economics and Business. The effects of metadata in the Business and Economics fields are more similar to those in Mathematics and Computer Science. This is possibly because a large proportion of Economics and Business papers rely on quantitative analysis of massive data to draw conclusions, which is different from the paradigm in most Art, History, and Philosophy papers. (2) Geography is embedded closer to social sciences than it is to Geology and Environmental Science. This is possibly because Geography is an interdisciplinary field. The physical geography subfield is more related to earth sciences while the human geography subfield is closer to social sciences. Overall, we still observe commonalities in the effects of metadata within high-level areas/subareas.

#### 5.4 Effect at Different Label Granularities

Now we examine the effect of metadata when predicting labels at different granularity levels. The MAG taxonomy [37] has 6 layers (from the most coarse-grained Layer 0 to the most fine-grained



**Figure 4: The effect of metadata on Layer- $j$  P@ $k$  scores (averaged over the fields that significantly benefit from leveraging that type of metadata using the classifier).**

**Table 4: The fields that benefit the most from leveraging venue information in terms of Layer-1 P@1, Layer-2 P@1, and Layer-3 P@1. “T”: text only. “+V”: text+venue. “Δ”: absolute performance improvement.**

		L1 P@1	L2 P@1		L3 P@1		
		<b>Parabel [33]</b>					
Physics	T	0.7178	T	0.7846	T	0.6118	
	+V	0.7284	Philosophy +V	0.8170	Philosophy +V	0.7209	
	Δ	1.06%	Δ	3.24%	Δ	10.91%	
		<b>Transformer [44]</b>					
Physics	T	0.7393	T	0.6728	Philosophy T	0.4105	
	+V	0.7495	Philosophy +V	0.7504	Philosophy +V	0.6071	
	Δ	1.02%	Δ	7.76%	Δ	19.66%	
		<b>OAG-BERT [24]</b>					
Geography	T	0.6626	Art T	0.6747	History T	0.3455	
	+V	0.6679	+V	0.6879	+V	0.4130	
	Δ	0.53%	Δ	1.32%	Δ	6.75%	

Layer 5). Layer-0 labels are the 19 fields excluded from our label space. Across all 20 datasets, 85.2%, 86.5%, 70.4%, 35.0%, and 13.2% of the papers have ground-truth Layer-1, Layer-2, Layer-3, Layer-4, and Layer-5 labels, respectively. Due to the low proportions of papers related to Layer-4 and Layer-5 tags, we only study the effect of metadata on predicting Layer-1, Layer-2, and Layer-3 tags.

Given one classifier and one type of metadata, we consider all fields that significantly benefit from the metadata using the classifier. In these fields, we calculate the Layer- $j$  P@ $k$  scores  $((j, k) = (1, 1), (2, 1), (2, 3), (3, 1), (3, 3))$ . The definition of Layer- $j$  P@ $k$  is very similar to that of P@ $k$  except that Layer- $j$  P@ $k$  considers the  $k$  most confident labels at Layer  $j$  instead of in the whole label space. When calculating Layer- $j$  P@ $k$ , we focus on papers with at least one ground-truth Layer- $j$  label. Then, we compute the absolute performance change of Layer- $j$  P@ $k$  scores after incorporating the considered type of metadata. The results are shown in Figure 4.

From Figure 4, we find that: (1) On average, venues can help scientific literature tagging for not only coarse-grained tags but also fine-grained ones. This may be counterintuitive from computer scientists' perspective because CS venues (e.g., "WWW") can hardly indicate very fine tags (e.g., "Link Farm"). Indeed, in the Computer Science (Conference) dataset, the contribution of venues on Layer-3 P@1 is subtle (e.g., 0.19% when using Parabel, -0.11% when using Transformer). However, in fields other than Computer Science, some venues do carry very fine-grained signals. For example, there are two venues "Journal of Roman Archaeology" and "Medieval Studies" in History and Philosophy, respectively. These two venues may strongly imply a paper's relevance to "Classical Archaeology" and "Medievalism", which are Layer-2 and Layer-3 tags, respectively. In Table 4, we list the fields that benefit the most from leveraging venue information in terms of Layer-1 P@1, Layer-2 P@1, and Layer-3 P@1, where we do observe that Philosophy and History are the biggest beneficiaries in terms of Layer-3 P@1. (2) Different from venues, authors are beneficial to fine-grained tagging but harmful to coarse-grained prediction. This observation consistently holds across all three classifiers.

## 6 RELATED WORK

**Extreme Multi-Label Text Classification.** In keeping with our discussion in Section 4, we divide related studies on extreme multi-label classification into three major categories. (1) *Bag-of-words classifiers* [2, 16, 19, 31–34, 39, 49, 50, 54] take sparse tf-idf features as input. To improve model efficiency, 1-vs-all approaches such as DiSMEC [2] and PPDSParse [49] explore parallelism and model size reduction via model weight truncation. In another direction, tree-based approaches apply various partitioning techniques on the large label space. For example, Parabel [33] partitions the labels to a balanced tree structure using 2-means clustering; Bonsai [19] improves Parabel by allowing multi-way and unbalanced partitions; XR-Linear [54] improves Parabel by incorporating various hard negative sampling schemes; AnnexML [39] partitions the labels via graph-based nearest neighbor indices. (2) *Sequence-based classifiers* [21, 44, 45, 47, 53, 59] employ deep neural architectures such as CNNs (e.g., XML-CNN [21]), RNNs (e.g., MeshProbeNet [45] and AttentionXML [53]), and Transformers (e.g., BertXML [44]) to learn semantic representations of input text sequences for classification. There are also studies aggregating shallow word embeddings and/or applying MLP layers to obtain document embeddings, such as Slice [15], DeepXML [8], DECAF [29], GalaXC [36], and ECLARE [30]. (3) *Pre-trained language model classifiers* [4, 17, 48, 52, 56, 60] propose to transfer the knowledge learned by PLMs from web-scale corpora to the classification task. For example, X-Transformer [4] complements the output of BERT [10], XLNet [46], or RoBERTa [25] with sparse tf-idf features; LightXML [17] adopts PLMs as the text encoder and

performs label shortlist and re-ranking with the same PLM; XR-Transformer [56] further proposes fast multi-resolution PLM fine-tuning. Despite the success of these models in extreme multi-label classification, they mainly focus on classifying plain text sequences and are less aware of document metadata. In contrast, our work proposes several straightforward ways to enhance text classifiers with metadata signals.

**Scientific Literature Tagging.** Classifying academic papers is a common evaluation task in text mining [6, 28, 58] and graph mining [11, 14] studies. However, most studies consider coarse-grained paper classification only (e.g., with 5-20 categories in the label space), the result of which is not subdivided enough to satisfy users' fine-grained interests. To tag papers on PubMed with fine-grained medical subject headings, the task of MeSH indexing has been extensively studied [9, 22, 31, 45, 52]. However, these models only use text or text+venue as input, leaving the effect of authors and references unexplored. Recently, Zhang et al. [59, 60] and Ye et al. [47] make use of metadata to tag papers with fields of study in MAG. Still, these studies are restricted to computer science and biomedicine fields only. In comparison, our work conducts a systematic study across 19 fields and three major types of classifiers.

## 7 CONCLUSIONS AND FUTURE WORK

In this work, we examine the effect of metadata on scientific literature tagging in 19 fields using three classifiers. Our results provide the following insights to practitioners aiming at building accurate scientific literature taggers: First, while previous metadata-aware approaches often directly use all types of available metadata, we show that not all of them are always beneficial. It is important to select useful metadata features based on the classifier's type, the field, and the granularity level of predicted tags. Second, although the state-of-the-art models for text sequence modeling have been dominated by Transformer-based models, we demonstrate that simple bag-of-words classifiers work comparably well or even better in many cases for large-scale fine-grained paper tagging, and may be more effective in leveraging different types of metadata. Third, despite the varying effects of different metadata types in different fields, the gain or loss induced by each metadata type is rather consistent across similar fields.

This study explores each type of metadata separately to avoid confounders. However, different types of metadata (e.g., a venue and an author) may interact with each other and provide additional hints for classification. In the future, it is of our interest to study the composite effect of multiple types of metadata.

## ACKNOWLEDGMENTS

We thank anonymous reviewers for their valuable and insightful feedback. Research was supported in part by the IBM-Illinois Discovery Accelerator Institute, US DARPA KAIROS Program No. FA8750-19-2-1004 and INCAS Program No. HR001121C0165, National Science Foundation IIS-19-56151, IIS-17-41317, and IIS 17-04532, and the Molecule Maker Lab Institute: An AI Research Institutes program supported by NSF under Award No. 2019897, and the Institute for Geospatial Understanding through an Integrative Discovery Environment (I-GUIDE) by NSF under Award No. 2118329. Any opinions, findings, and conclusions or recommendations expressed herein are those of the authors and do not necessarily represent the views, either expressed or implied, of DARPA or the U.S. Government.



## A APPENDIX

### A.1 Datasets

Supplementary to Table 1, Table 5 summarizes more dataset statistics in MAPLE, including the average numbers of authors, references, and labels per paper as well as the numbers of training, validation, and testing papers.

**Table 5: More statistics of the 20 datasets in MAPLE.**

Field	#Authors/ Paper	#References/ Paper	#Labels/ Paper	#Train Papers	#Valid Papers	#Test Papers
Art	1.314	2.813	2.425	39,901	9,975	8,497
Philosophy	1.117	8.121	3.715	40,696	10,174	8,426
Geography	3.625	31.793	2.231	44,710	11,178	17,995
Business	2.346	37.635	3.556	49,481	12,370	23,007
Sociology	1.592	24.474	2.318	60,815	15,204	14,189
History	1.218	4.405	1.938	80,830	20,208	12,109
Political Science	1.496	10.198	3.107	73,650	18,413	23,228
Environmental Science	4.301	34.156	2.138	70,479	17,620	35,846
Economics	1.966	26.608	4.961	119,309	29,827	29,534
Engineering	3.055	19.703	4.987	179,804	44,951	45,251
Psychology	3.805	41.201	5.032	245,288	61,322	66,344
CS (Conference)	3.389	17.417	6.089	136,322	34,081	92,990
CS (Journal)	3.296	21.968	5.976	214,616	53,654	142,333
Geology	3.552	36.510	5.957	284,022	71,006	76,806
Mathematics	2.151	19.068	6.433	338,454	84,613	67,484
Materials Science	4.829	26.838	4.428	783,289	195,822	358,620
Physics	9.841	23.509	7.271	905,613	226,403	237,967
Biology	5.611	36.137	8.306	1,125,605	281,401	181,772
Chemistry	4.293	28.382	6.288	1,252,531	313,133	284,292
Medicine	5.647	12.888	5.376	1,620,165	405,041	620,899

### A.2 Hyperparameters

In each of the 20 datasets, we remove metadata instances that appear in less than 5 papers. We adopt the same hyperparameter configuration when using text only, text+venue, text+author, and text+reference as input. The code source and hyperparameters of each classifier are introduced below.

**A.2.1 Parabel.**<sup>4</sup> We remove words appearing in less than 5 papers. All parameters are set by default. Specifically, the number of threads  $T = 1$ ; the number of trees  $t = 3$ ; the maximum number of labels in a leaf node  $m = 100$ ; the beam search width in prediction  $B = 10$ .

**A.2.2 Transformer.**<sup>5</sup> The number of Transformer layers is 1; the number of attention heads is 2; the number of [CLS] tokens  $C = 10$ ; the embedding dimension is 100; the maximum sequence length is 500; the batch size is 256. We adopt GloVe.6B.100d as initialized word embeddings. For (the number of training epochs, the number of warm-up epochs), we use (100, 20) for Art, Philosophy, and Geography, (80, 16) for Business, Sociology, History, and Political Science, (60, 12) for Environmental Science, Economics, Engineering, and Computer Science, (40, 8) for Psychology, Geology, and Mathematics, (20, 4) for Materials Science, Physics, and Biology, and (15, 3) for Chemistry and Medicine. Other hyperparameters are set by default.

**A.2.3 OAG-BERT.**<sup>6</sup> We use “oagbert-v2” as the PLM. After PLM encoding, we fix paper embeddings to train a Parabel classifier. All parameters of Parabel are set by default.

<sup>4</sup><http://manikvarma.org/code/Parabel/download.html>

<sup>5</sup><https://github.com/XunGuangxu/CorNet> (We use the BertXML classifier in this GitHub repository. Although the model is called BertXML, it trains a Transformer architecture from scratch without BERT initialization.)

<sup>6</sup><https://github.com/THUDM/OAG-BERT>

### A.3 More on the Effect of Metadata at Different Label Granularities

Supplementary to Table 4, Table 6 shows the fields that benefit the most from leveraging author or reference information in terms of Layer-1 P@1, Layer-2 P@1, and Layer-3 P@1. We find that authors and references are beneficial to the Art, Philosophy, and History fields in many cases. The possible reason is that each paper in these fields has a small number of authors and references (according to the statistics in Table 5). Therefore, the author or reference list may contain fewer confounding signals and be more topic-indicative.

**Table 6: The fields that benefit the most from leveraging author or reference information in terms of Layer-1 P@1, Layer-2 P@1, and Layer-3 P@1. “T”: text only. “+A”: text+author. “+R”: text+reference. “Δ”: absolute performance improvement.**

		L1 P@1		L2 P@1		L3 P@1		
<b>Parabel [33]</b>								
Mathematics	T	0.5960	Political Science	T	0.7928	Philosophy	T	0.6118
	+A	0.5971		+A	0.7969		+A	0.6160
	Δ	0.11%		Δ	0.41%		Δ	0.42%
Psychology	T	0.6992	Business	T	0.6704	Business	T	0.6124
	+R	0.7086		+R	0.6801		+R	0.6265
	Δ	0.94%		Δ	0.97%		Δ	1.41%
<b>Transformer [44]</b>								
History	T	0.5471	Philosophy	T	0.6728	Art	T	0.4859
	+A	0.5489		+A	0.6900		+A	0.4910
	Δ	0.18%		Δ	1.72%		Δ	0.51%
History	T	0.5471	Art	T	0.7859	History	T	0.4857
	+R	0.5537		+R	0.7953		+R	0.4890
	Δ	0.66%		Δ	0.94%		Δ	0.33%
<b>OAG-BERT [24]</b>								
Art	T	0.6241	Art	T	0.6747	History	T	0.3455
	+A	0.6254		+A	0.6875		+A	0.3819
	Δ	0.13%		Δ	1.28%		Δ	3.64%

### A.4 Effect of Metadata on Efficiency

Now we report the effect of metadata on model efficiency. Table 7 shows the average relative training time increase of the three classifiers across the 20 datasets after incorporating venues, authors, and references, respectively. All models are run on Intel Xeon E5-2680 v2 @ 2.80GHz and one NVIDIA GeForce GTX 1080 Ti GPU (if a GPU is needed). We can observe a significant increase in training time after adding references as features. When Parabel is the classifier, the increase is due to a much longer bag-of-words vector used to represent a paper; when Transformer is the classifier, the increase is caused by a large number of additional parameters (i.e., reference embeddings), which make the model converge more slowly.

**Table 7: Average relative increase in training time across 20 datasets after incorporating one type of metadata.**

	Parabel [33]	Transformer [44]	OAG-BERT [24]
+Venue	+0.15%	+0.13%	+0.11%
+Author	+0.72%	+1.86%	+0.11%
+Reference	+22.68%	+10.75%	-

Figure 5 shows the training time of the three classifiers on the 20 datasets. The reported training time is an average over 20 runs (i.e., 5 runs × {text only, text+venue, text+author, text+reference}) when Parabel and Transformer are tested, or 15 runs (i.e., 5 runs × {text only, text+venue, text+author}) when OAG-BERT is tested. Among the three classifiers, Parabel is consistently the most efficient across the 20 datasets; Transformer spends more time than OAG-BERT on

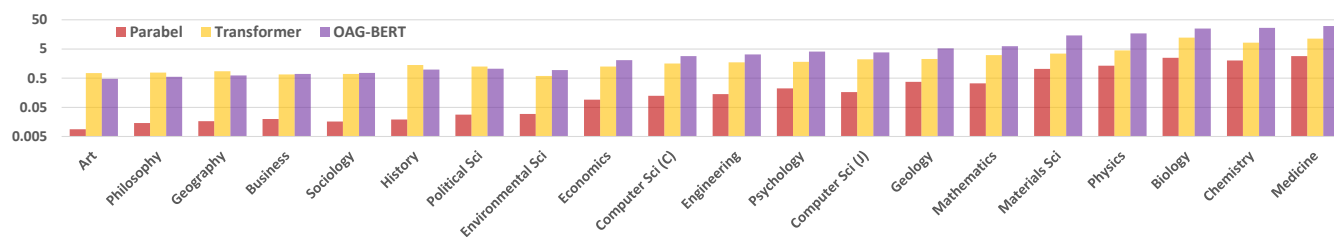


Figure 5: Training time (in hours) of the three classifiers on the 20 datasets.

Table 8: Statistics of the three additional datasets with MeSH labels.

Field	Paper Source	#Papers	#Labels	#Venues	#Authors	#References	#Authors/ Paper	#References/ Paper	#Labels/ Paper	#Train Papers	#Valid Papers	#Test Papers
Biology-MeSH	Journal	1,379,393	25,039	100	2,486,814	6,876,739	5.686	40.043	13.870	985,364	246,341	147,688
Chemistry-MeSH	Journal	762,129	21,585	87	1,498,358	5,928,908	4.741	34.344	10.984	511,814	127,954	122,361
Medicine-MeSH	Journal	1,536,660	25,188	100	2,791,165	7,190,021	5.254	20.931	11.819	1,020,969	255,242	260,449

Table 9: P@k and NDCG@k scores on the three additional datasets. Blue, Red, and “-”: the same meaning as in Table 2.

Field	Input	Parabel [33]					Transformer [44]					OAG-BERT [24]				
		P@1	P@3	P@5	N@3	N@5	P@1	P@3	P@5	N@3	N@5	P@1	P@3	P@5	N@3	N@5
Biology-MeSH	Text	0.8924	0.7964	0.7088	0.8194	0.7576	0.9112	0.8098	0.7193	0.8341	0.7698	0.7506	0.6185	0.5460	0.6490	0.5945
	+Venue	0.8934	0.7976	0.7101	0.8206	0.7587	0.9119	0.8105	0.7194	0.8348	0.7701	0.7520	0.6200	0.5473	0.6504	0.5958
	+Author	0.8932	0.7985	0.7112	0.8212	0.7596	0.9090	0.8028	0.7093	0.8281	0.7614	0.7504	0.6184	0.5464	0.6488	0.5946
	+Reference	0.8976	0.8058	0.7198	0.8279	0.7674	0.9079	0.7988	0.7034	0.8248	0.7565	-	-	-	-	-
Chemistry-MeSH	Text	0.8447	0.7340	0.6407	0.7603	0.6947	0.8445	0.7113	0.6082	0.7427	0.6684	0.6971	0.5804	0.5099	0.6071	0.5557
	+Venue	0.8453	0.7354	0.6421	0.7615	0.6960	0.8465	0.7151	0.6121	0.7460	0.6720	0.6995	0.5826	0.5132	0.6093	0.5587
	+Author	0.8450	0.7350	0.6419	0.7611	0.6957	0.8439	0.7105	0.6066	0.7419	0.6670	0.6981	0.5819	0.5126	0.6086	0.5580
	+Reference	0.8490	0.7409	0.6491	0.7667	0.7023	0.8248	0.6802	0.5743	0.7139	0.6366	-	-	-	-	-
Medicine-MeSH	Text	0.9673	0.8410	0.7454	0.8712	0.8075	0.9683	0.8567	0.7583	0.8842	0.8191	0.7343	0.6229	0.5629	0.6491	0.6089
	+Venue	0.9673	0.8422	0.7472	0.8722	0.8090	0.9688	0.8585	0.7606	0.8856	0.8210	0.7370	0.6235	0.5607	0.6503	0.6081
	+Author	0.9671	0.8342	0.7424	0.8656	0.8039	0.9687	0.8495	0.7496	0.8786	0.8117	0.7347	0.6185	0.5549	0.6459	0.6027
	+Reference	0.9666	0.8382	0.7463	0.8687	0.8073	0.9662	0.8470	0.7458	0.8762	0.8084	-	-	-	-	-

small datasets, but the situation is reversed on moderate-sized and large datasets.

### A.5 Additional Datasets with MeSH Labels

To further strengthen our findings, we construct three datasets with Medical Subject Headings (MeSH) terms [7] as their labels, which are curated by experts from the National Library of Medicine. To be specific, we take the three datasets, Biology, Chemistry, and Medicine, from MAPLE and obtain the ground-truth MeSH labels of each paper<sup>7</sup>. After removing papers not having MeSH labels, we get three new datasets, Biology-MeSH, Chemistry-MeSH, and Medicine-MeSH. Their statistics are shown in Table 8.

We run Parabel, Transformer, and OAG-BERT on the three new datasets. When running Transformer, for (the number of training epochs, the number of warm-up epochs), we use (20, 4) for all three datasets. When running OAG-BERT, we need to rerank those labels appearing in the paper text higher than those not. Note that a MeSH label may have multiple label names (i.e., one canonical name and 0, 1, or several entry terms, see the MeSH label “COVID-19”<sup>8</sup> as an example). Given a MeSH label, if any of its label names appears in the paper text, we view it as occurring in the paper. All other hyperparameters are the same as in Appendix A.2.

The P@k and NDCG@k scores of Parabel, Transformer, and OAG-BERT on the three new datasets are demonstrated in Table 9. In general, we still find that venues are beneficial to scientific

literature tagging in almost all cases, while the effect of authors and references depends on the classifier’s type and the field.

The three additional datasets are also available in our MAPLE benchmark: <https://doi.org/10.5281/zenodo.7611544>.

## REFERENCES

- [1] Waleed Ammar, Dirk Groeneveld, Chandra Bhagavatula, Iz Beltagy, Miles Crawford, Doug Downey, Jason Dunkelberger, Ahmed Elgohary, Sergey Feldman, Vu Ha, et al. 2018. Construction of the Literature Graph in Semantic Scholar. In *NAACL-HLT’18*. 84–91.
- [2] Rohit Babbar and Bernhard Schölkopf. 2017. Dismec: Distributed sparse machines for extreme multi-label classification. In *WSDM’17*. 721–729.
- [3] Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. SciBERT: A Pretrained Language Model for Scientific Text. In *EMNLP’19*. 3615–3620.
- [4] Wei-Cheng Chang, Hsiang-Fu Yu, Kai Zhong, Yiming Yang, and Inderjit S Dhillon. 2020. Taming pretrained transformers for extreme multi-label text classification. In *KDD’20*. 3163–3171.
- [5] Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv:1412.3555* (2014).
- [6] Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel S Weld. 2020. SPECTER: Document-level Representation Learning using Citation-informed Transformers. In *ACL’20*. 2270–2282.
- [7] Margaret H Coletti and Howard L Bleich. 2001. Medical subject headings used to search the biomedical literature. *JAMIA* 8, 4 (2001), 317–323.
- [8] Kunal Dahiya, Deepak Saini, Anshul Mittal, Ankush Shaw, Kushal Dave, Akshay Soni, Himanshu Jain, Sumeet Agarwal, and Manik Varma. 2021. Deepxml: A deep extreme multi-label learning framework applied to short text documents. In *WSDM’21*. 31–39.
- [9] Suyang Dai, Ronghui You, Zhiyong Lu, Xiaodi Huang, Hiroshi Mamitsuka, and Shanfeng Zhu. 2020. FullMeSH: improving large-scale MeSH indexing with full text. *Bioinformatics* 36, 5 (2020), 1533–1541.
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In

<sup>7</sup><https://learn.microsoft.com/en-us/academic-services/graph/reference-data-schema#paper-mesh>

<sup>8</sup><https://meshb.nlm.nih.gov/record/ui?ui=D000086382>

- NAACL-HLT'19. 4171–4186.
- [11] Yuxiao Dong, Nitesh V Chawla, and Ananthram Swami. 2017. metapath2vec: Scalable representation learning for heterogeneous networks. In *KDD'17*. 135–144.
  - [12] Yuxiao Dong, Hao Ma, Zhihong Shen, and Kuansan Wang. 2017. A century of science: Globalization of scientific collaborations, citations, and innovations. In *KDD'17*. 1437–1446.
  - [13] Jorge E Hirsch. 2005. An index to quantify an individual's scientific research output. *PNAS* 102, 46 (2005), 16569–16572.
  - [14] Ziniu Hu, Yuxiao Dong, Kuansan Wang, and Yizhou Sun. 2020. Heterogeneous graph transformer. In *WWW'20*. 2704–2710.
  - [15] Himanshu Jain, Venkatesh Balasubramanian, Bhanu Chanduri, and Manik Varma. 2019. Slice: Scalable linear extreme classifiers trained on 100 million labels for related searches. In *WSDM'19*. 528–536.
  - [16] Himanshu Jain, Yashoteja Prabhu, and Manik Varma. 2016. Extreme multi-label loss functions for recommendation, tagging, ranking & other missing label applications. In *KDD'16*. 935–944.
  - [17] Ting Jiang, Deqing Wang, Leilei Sun, Huayi Yang, Zhengyang Zhao, and Fuzhen Zhuang. 2021. Lightxml: Transformer with dynamic negative sampling for high-performance extreme multi-label text classification. In *AAAI'21*. 7987–7994.
  - [18] Ching Jin, Yifang Ma, and Brian Uzzi. 2021. Scientific prizes and the extraordinary growth of scientific topics. *Nature Communications* 12, 1 (2021), 5619.
  - [19] Sujay Khandagale, Han Xiao, and Rohit Babbar. 2020. Bonsai: diverse and shallow trees for extreme multi-label classification. *Machine Learning* 109, 11 (2020), 2099–2119.
  - [20] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* 36, 4 (2020), 1234–1240.
  - [21] Jingzhou Liu, Wei-Cheng Chang, Yuexin Wu, and Yiming Yang. 2017. Deep learning for extreme multi-label text classification. In *SIGIR'17*. 115–124.
  - [22] Ke Liu, Shengwen Peng, Junqiu Wu, Chengxiang Zhai, Hiroshi Mamitsuka, and Shanfeng Zhu. 2015. MeSHLabeler: improving the accuracy of large-scale MeSH indexing by integrating diverse evidence. *Bioinformatics* 31, 12 (2015), i339–i347.
  - [23] Weiwei Liu, Haobo Wang, Xiaobo Shen, and Ivor W Tsang. 2021. The emerging trends of multi-label learning. *IEEE TPAMI* 44, 11 (2021), 7955–7974.
  - [24] Xiao Liu, Da Yin, Jingnan Zheng, Xingjian Zhang, Peng Zhang, Hongxia Yang, Yuxiao Dong, and Jie Tang. 2022. OAG-BERT: Towards a Unified Backbone Language Model for Academic Knowledge Services. In *KDD'22*. 3418–3428.
  - [25] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
  - [26] Zhiyong Lu. 2011. PubMed and beyond: a survey of web tools for searching biomedical literature. *Database* 2011 (2011).
  - [27] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *JMLR* 9 (2008), 2579–2605.
  - [28] Dheeraj Mekala, Xinyang Zhang, and Jingbo Shang. 2020. META: Metadata-Empowered Weak Supervision for Text Classification. In *EMNLP'20*. 8351–8361.
  - [29] Anshul Mittal, Kunal Dahiya, Sheshansh Agrawal, Deepak Saini, Sumeet Agarwal, Purushottam Kar, and Manik Varma. 2021. Decaf: Deep extreme classification with label features. In *WSDM'21*. 49–57.
  - [30] Anshul Mittal, Naveen Sachdeva, Sheshansh Agrawal, Sumeet Agarwal, Purushottam Kar, and Manik Varma. 2021. ECLARE: Extreme Classification with Label Graph Correlations. In *WWW'21*. 3721–3732.
  - [31] Shengwen Peng, Ronghui You, Hongming Wang, Chengxiang Zhai, Hiroshi Mamitsuka, and Shanfeng Zhu. 2016. DeepMeSH: deep semantic representation for improving large-scale MeSH indexing. *Bioinformatics* 32, 12 (2016), i70–i79.
  - [32] Yashoteja Prabhu, Anil Kag, Shilpa Gopinath, Kunal Dahiya, Shrutendra Harsola, Rahul Agrawal, and Manik Varma. 2018. Extreme multi-label learning with label features for warm-start tagging, ranking & recommendation. In *WSDM'18*. 441–449.
  - [33] Yashoteja Prabhu, Anil Kag, Shrutendra Harsola, Rahul Agrawal, and Manik Varma. 2018. Parabel: Partitioned label trees for extreme classification with application to dynamic search advertising. In *WWW'18*. 993–1002.
  - [34] Yashoteja Prabhu and Manik Varma. 2014. Fastxml: A fast, accurate and stable tree-classifier for extreme multi-label learning. In *KDD'14*. 263–272.
  - [35] Martin Rosvall and Carl T Bergstrom. 2011. Multilevel compression of random walks on networks reveals hierarchical organization in large integrated systems. *PLoS one* 6, 4 (2011), e18209.
  - [36] Deepak Saini, Arnav Kumar Jain, Kushal Dave, Jian Jiao, Amit Singh, Ruofei Zhang, and Manik Varma. 2021. GalaXC: Graph Neural Networks with Labelwise Attention for Extreme Classification. In *WWW'21*. 3733–3744.
  - [37] Zhihong Shen, Hao Ma, and Kuansan Wang. 2018. A Web-scale system for scientific knowledge exploration. In *ACL'18 System Demonstrations*. 87–92.
  - [38] Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, Bo-June Hsu, and Kuansan Wang. 2015. An overview of microsoft academic service (mas) and applications. In *WWW'15 Companion*. 243–246.
  - [39] Yukihiro Tagami. 2017. Annexml: Approximate nearest neighbor search for extreme multi-label classification. In *KDD'17*. 455–464.
  - [40] Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, and Zhong Su. 2008. Arnetminer: extraction and mining of academic social networks. In *KDD'08*. 990–998.
  - [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS'17*. 5998–6008.
  - [42] Kuansan Wang, Zhihong Shen, Chiyan Huang, Chieh-Han Wu, Yuxiao Dong, and Anshul Kanakia. 2020. Microsoft academic graph: When experts are not enough. *Quantitative Science Studies* 1, 1 (2020), 396–413.
  - [43] Boya Xie, Zhihong Shen, and Kuansan Wang. 2021. Is preprint the future of science? A thirty year journey of online preprint services. *arXiv preprint arXiv:2102.09066* (2021).
  - [44] Guangxu Xun, Kishlay Jha, Jianhui Sun, and Aidong Zhang. 2020. Correlation networks for extreme multi-label text classification. In *KDD'20*. 1074–1082.
  - [45] Guangxu Xun, Kishlay Jha, Ye Yuan, Yaqing Wang, and Aidong Zhang. 2019. MeSHProbeNet: a self-attentive probe net for MeSH indexing. *Bioinformatics* 35, 19 (2019), 3794–3802.
  - [46] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *NeurIPS'19*.
  - [47] Chenchen Ye, Linhai Zhang, Yulan He, Deyu Zhou, and Jie Wu. 2021. Beyond Text: Incorporating Metadata and Label Structure for Multi-Label Document Classification using Heterogeneous Graphs. In *EMNLP'21*. 3162–3171.
  - [48] Hui Ye, Zhiyu Chen, Da-Han Wang, and Brian Davison. 2020. Pretrained generalized autoregressive model with adaptive probabilistic label clusters for extreme multi-label text classification. In *ICML'20*. 10809–10819.
  - [49] Ian EH Yen, Xiangru Huang, Wei Dai, Pradeep Ravikumar, Inderjit Dhillon, and Eric Xing. 2017. Pdp sparse: A parallel primal-dual sparse method for extreme classification. In *KDD'17*. 545–553.
  - [50] Ian En-Hsu Yen, Xiangru Huang, Pradeep Ravikumar, Kai Zhong, and Inderjit Dhillon. 2016. Pd-sparse: A primal and dual sparse approach to extreme multiclass and multilabel classification. In *ICML'16*. 3069–3077.
  - [51] Yian Yin, Yuxiao Dong, Kuansan Wang, Dashun Wang, and Benjamin F Jones. 2022. Public use and public funding of science. *Nature Human Behaviour* 6, 10 (2022), 1344–1350.
  - [52] Ronghui You, Yuxuan Liu, Hiroshi Mamitsuka, and Shanfeng Zhu. 2021. BERTMeSH: deep contextual representation learning for large-scale high-performance MeSH indexing with full text. *Bioinformatics* 37, 5 (2021), 684–692.
  - [53] Ronghui You, Zihan Zhang, Ziyi Wang, Suyang Dai, Hiroshi Mamitsuka, and Shanfeng Zhu. 2019. Attentionxml: Label tree-based attention-aware deep model for high-performance extreme multi-label text classification. *NeurIPS'19* (2019), 5820–5830.
  - [54] Hsiang-Fu Yu, Kai Zhong, Jiong Zhang, Wei-Cheng Chang, and Inderjit S Dhillon. 2022. PECOS: Prediction for enormous and correlated output spaces. *JMLR* 23, 98 (2022), 1–32.
  - [55] Fanjin Zhang, Xiao Liu, Jie Tang, Yuxiao Dong, Peiran Yao, Jie Zhang, Xiaotao Gu, Yan Wang, Bin Shao, Rui Li, et al. 2019. Oag: Toward linking large-scale heterogeneous entity graphs. In *KDD'19*. 2585–2595.
  - [56] Jiong Zhang, Wei-Cheng Chang, Hsiang-Fu Yu, and Inderjit Dhillon. 2021. Fast multi-resolution transformer fine-tuning for extreme multi-label text classification. In *NeurIPS'21*. 7267–7280.
  - [57] Yu Zhang, Xiuxi Chen, Yu Meng, and Jiawei Han. 2021. Hierarchical Metadata-Aware Document Categorization under Weak Supervision. In *WSDM'21*. 770–778.
  - [58] Yu Zhang, Shweta Garg, Yu Meng, Xiuxi Chen, and Jiawei Han. 2022. Motifclass: Weakly supervised text classification with higher-order metadata information. In *WSDM'22*. 1357–1367.
  - [59] Yu Zhang, Zhihong Shen, Yuxiao Dong, Kuansan Wang, and Jiawei Han. 2021. MATCH: Metadata-Aware Text Classification in A Large Hierarchy. In *WWW'21*. 3246–3257.
  - [60] Yu Zhang, Zhihong Shen, Chieh-Han Wu, Boya Xie, Junheng Hao, Ye-Yi Wang, Kuansan Wang, and Jiawei Han. 2022. Metadata-Induced Contrastive Learning for Zero-Shot Multi-Label Text Classification. In *WWW'22*. 3162–3173.