# DiVa: An Iterative Framework to Harvest More Diverse and Valid Labels from User Comments for Music

Hongru Liang
Sichuan University
Chengdu, China
lianghongru@scu.edu.cn

Jingyao Liu
Sichuan University
Chengdu, China
liujingyao@stu.scu.edu.cn

Yuanxin Xiang
Sichuan University
Chengdu, China
whutxyx@gmail.com

Jiachen Du
Tencent Music Entertainment
Shenzhen, China
jacobdu@tencent.com

Lanjun Zhou
Tencent Music Entertainment
Shenzhen, China
jedzhou@tencent.com

Shushen Pan
Tencent Music Entertainment
Shenzhen, China
samsonpan@tencent.com

Wenqiang Lei*
Sichuan University
Chengdu, China
wenqianglei@scu.edu.cn

## ABSTRACT

Towards sufficient music searching, it is vital to form a complete set of labels for each song. However, current solutions fail to resolve it as they cannot produce diverse enough mappings to make up for the information missed by the gold labels. Based on the observation that such missing information may already be presented in user comments, we propose to study the automated music labeling in an essential but under-explored setting, where the model is required to harvest more diverse and valid labels from the users' comments given limited gold labels. To this end, we design an iterative framework (DiVa) to harvest more Diverse and Valid labels from user comments for music. The framework makes a classifier able to form complete sets of labels for songs via pseudo-labels inferred from pre-trained classifiers and a novel joint score function. The experiment on a densely annotated testing set reveals the superiority of the DiVa over state-of-the-art solutions in producing more diverse labels missed by the gold labels. We hope our work can inspire future research on automated music labeling.

## CCS CONCEPTS

• **Information systems** → **Music retrieval**; • **Computing methodologies** → **Information extraction**.

## KEYWORDS

Automated Music Labeling, User Comments, Music Searching

*Corresponding author.

## 1 INTRODUCTION

Towards sufficient music searching, especially at the industrial scale, an ideal solution is to hire experts to annotate the complete set of all possible labels for a song. As shown in Figure 1, an expert is hired to select the gold labels for a song (i.e., *We Are The World*) from the gold label vocabulary. However, such gold labels only take part of all possible labels of the song. The information missed by the gold labels exists in two scenarios. First, there lack of mappings to labels out of the gold label vocabulary (OOV labels), e.g., "**unity**". This is because we can never provide the annotator with a complete vocabulary covering all music characteristics and user searching preferences [22]. Second, even with a complete vocabulary, the mappings to some labels may also be missed by the expert (missed supervisions), e.g., "**hope**". This is because the expert can never be serious and careful enough to pinpoint all the right labels from thousands or more candidate labels in the vocabulary [28]. Therefore, to form the complete set of labels for a song, it is needed to obtain more diverse mappings to OOV labels and missed supervisions.

However, current multi-label classifiers [13, 14], including recent extreme multi-label learning methods [23, 26], fail to resolve the above challenge. Because they target to approach a vocabulary-size vector, where the elements indicating gold labels are set to 1 and other elements are set to 0. A few studies try to obtain more mappings using generative models [22], self-training methods [27, 28, 32], etc. However, such new mappings are not diverse enough to make up for the missing information. Because these newly mapped labels are either semantically similar to the gold labels or are copied from instances with similar features. To our best knowledge, it remains an open question how to automatically obtain the complete set of labels for instance given a small number of gold labels.

**(a) *We Are The World*** **(b) Complete set of labels** **(c) User Comments**
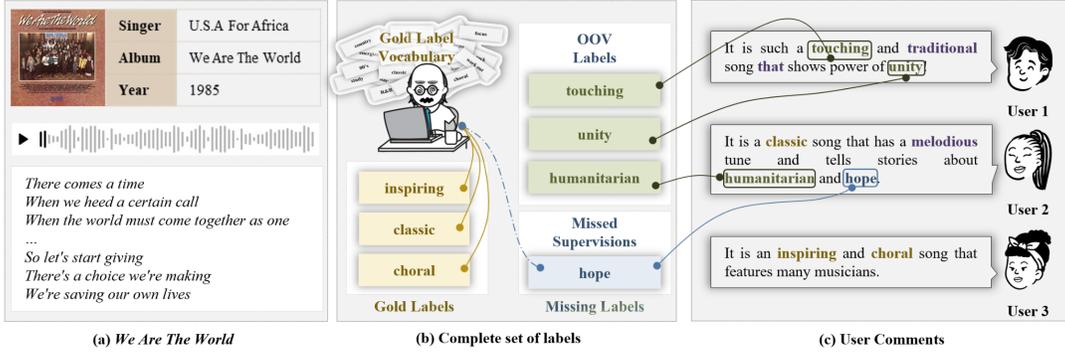
Figure 1: Illustration of (a) a song titled *We Are The World*, (b) the complete set of labels for the song, and (c) the user comments on the song. The gold labels annotated by the expert only occupy part of all possible labels of the song. There lack of mappings to labels out of the gold label vocabulary (OOV labels, e.g., "touching") and some labels in the gold label vocabulary (missed supervisions, e.g., "hope"). Fortunately, all possible labels for the song are already presented in the user comments.

Fortunately, there is a good chance to solve this question with the help of users. A helpful observation is that users themselves are "experts" in music searching. Specifically, the vocabulary the users used to comment on a song naturally matches the vocabulary they used to search for the song [5]. Besides, the user comments are continuously updated along with the changing of user's feelings and the latest Internet buzzwords, so they can reflect the dynamic changes in user searching preferences [16]. As such, the complete set of OOV labels, missed supervisions, and gold labels for a song may already be presented in the user comments of the song, cf., Figure 1 (c). This observation inspires us to harvest more labels from user comments. We believe such "user-generated" labels can largely enrich the expert-generated gold labels and support more sufficient music searching. Specifically, towards sufficient music searching, the newly generated labels should have two basic advantages. First, they should be adequately **valid** to serve as a label (e.g., "**touching**" vs. "**that**") and to benefit music searching (e.g., "**touching**" vs. "**melodious**"). Second, they should be **diverse** enough to semantically distinguish from the gold ones (e.g., "**classic**" vs. "**traditional**"). This encourages us to study automated music labeling in a vital but under-explored setting — given the limited gold labels of a song, the model is required to harvest more diverse and valid labels from the user comments for the song.

To this end, we propose DiVa, which is an iterative framework for harvesting more <u>Di</u>verse and <u>Va</u>lid labels for music. It consists of two key components, i.e., a binary classifier and a joint score function. The binary classifier takes the concatenation of user comments and a candidate label as input and outputs a confidence score indicating how likely the candidate label is a suitable label for the song. In the beginning, the binary classifier, which has been trained on the gold labels, has trouble predicting more diverse mappings than the gold ones. Inspired by Xie et al. [27], we use the pre-trained classifier to infer the pseudo-mappings to missed supervisions. In addition, we design a novel joint score function, which comprehensively measures the diversity and validness of a candidate label to infer the pseudo-mappings to OVV labels. In this way, the binary classifier can be fine-tuned on more diverse and valid labels than the gold ones. We repeat this process several iterations until the classifier can infer a complete set of labels for a song, namely, the classifier can barely infer any new pseudo-labels. We further collect

a corpus[1] from a Chinese online music platform[2]. It consists of 16,372 songs, each of which consists of user comments and gold labels. Towards data-driven evaluation, we also develop a testing set of label 1,000 songs, which have been densely annotated by experts based on the user comments. Experiment results show that, compared with state-of-the-art methods, the DiVa framework is more powerful in harvesting more diverse and valid labels for music than the state-of-the-art methods. We hope our work can provide insights into automated music labeling and advance the research on finding the complete set of labels in real-world applications. In summary, we highlight our contributions as follows.

- We call attention to the information missed by the gold labels and study the automated music labeling in a challenging setting, where the model is required to produce more diverse and valid labels based on a small number of gold labels.
- We propose an iterative framework (DiVa) to harvest more diverse and valid labels from user comments. It makes a classifier infer more diverse and valid labels for music gradually with the help of a novel joint score function.
- The experiment reveals the superiority of DiVa in producing more diverse and valid mappings to OOV labels and missed supervisions as well as more complete sets of labels for songs.

## 2 PROBLEM FORMULATION

Let $\{d_i, X_i, \mathcal{Y}_i^G, \mathcal{Y}_i\}_{i=1}^N$ denote a dataset of $N$ songs, where $d_i$, $X_i$, $\mathcal{Y}_i^G$, and $\mathcal{Y}_i$ are the user comments, the set of words tokenized from user comments, gold labels, and the complete set of labels for the $i^{th}$ song, respectively. We assume $\mathcal{Y}_i$ is a proper superset of $\mathcal{Y}_i^G$ and a proper subset of $X_i$, cf., Figure 2. Besides, $\mathcal{Y}_i$ is unknown both at the training and testing phases. Let $\mathbb{Y}^G = \bigcup_{i=1}^N \mathcal{Y}_i^G$ denote the gold label vocabulary. The OOV labels in Figure 1 can be express as $\mathcal{Y}_i^O = \mathcal{Y}_i - \mathbb{Y}^G$ and the missed supervisions can be expressed as $\mathcal{Y}_i^S = \mathcal{Y}_i \cap \mathbb{Y}^G - \mathcal{Y}_i^G$. The task of harvesting more labels from the user comments is to find the $\mathcal{Y}_i^+ \subset X_i$ satisfying the following

---

[1]The corpus and our codes are available at https://github.com/jingyaolliu/DiVa.
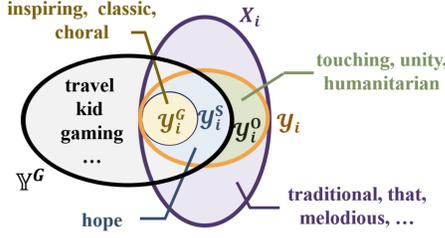[2]For convenience and saving spaces, all presented examples in this paper have been translated into English.

**Figure 2: Venn diagram for the labels of the $i^{th}$ song and the words in the user comments of the song**

conditions:

$$\arg\max_{\mathcal{Y}_i^+} \frac{|\mathcal{Y}_i \cap (\mathcal{Y}_i^G \cup \mathcal{Y}_i^+)|}{|\mathcal{Y}_i \cup (\mathcal{Y}_i^G \cup \mathcal{Y}_i^+)|}. \tag{1}$$

In other words, $\mathcal{Y}_i^+$ is required to cover as much as information missed by the gold labels (i.e., $\mathcal{Y}_i^O \cup \mathcal{Y}_i^S$) using as few as labels. Note that, semantically similar labels (e.g., "**classic**" and "**traditional**") are treated as the same label in this task.

## 3  THE DIVA FRAMEWORK

As shown in Figure 3, the DiVa framework starts from a binary classifier trained on paired user comments ($d_i$) and gold labels ($\mathcal{Y}_i^G$). The classifier is then used to infer the possibility distribution on candidate labels and select those with high confidence scores as pseudo-labels. Further, we design a joint score function to infer the possible distribution of the remaining candidates. All pseudo-labels are used to fine-tune the previous classifier. This process will repeat several iterations until the classifier can barely infer any new pseudo-labels or reach the maximum iterations. In the testing phase, we can harvest labels from user comments directly from the optimized binary classifier.

### 3.1  Binary Classifier

The multi-label classifier is limited to a close set of labels as it is required to fit a fixed-size multi-hot vector. To overcome this limitation, we lay our foundation on a binary classifier, which is required to fit a scalar-value label. Specifically, we pair the user comments ($d_i$) and a candidate label (denoted as $y_c$) of the $i^{th}$ song as a new instance. If the candidate label is a gold label to the $i^{th}$ song, we label the new instance as 1; otherwise, the new instance is labeled as 0. At this moment, the binary classifier is equivalent to the conventional multi-label classifier.

To make it different, we make two changes to the binary classifier. First, we change the candidate labels for the $i^{th}$ song from labels in the gold label vocabulary to the words from user comments plus its gold labels ($\mathcal{Y}_i^G$). As such, we enrich the label vocabulary without increasing the number of candidate labels for the $i^{th}$ song. Second, to better correlate document-level user comments with word-level labels, we use a contextual language model (i.e., XLNet [29]) to represent the user comments and use the improved static word embeddings (i.e., X2Static [8]) to represent the labels. We then concatenate the representations of user comments ($\vec{d_i}$) and a candidate label ($\vec{y_c}$) and use the concatenations as the input of the binary classifier, which is defined as

$$\mathcal{F}(d_i, y_c) = sigmoid(\vec{d_i} \oplus \vec{y_c}), \tag{2}$$

where $\oplus$ is the concatenate operation. The output of the $\mathcal{F}(d_i, y_c)$ is a confidence score indicating how likely $y_c$ is a gold label for the $i^{th}$ song. The objective of the classifier is to minimize the sum of the binary cross entropy between these confidence scores and the ground truth values (0 or 1).

### 3.2  Binary Classifier Inference

We then utilize the pre-trained classifier to infer the possibility distribution on the candidate labels, which involves the remaining labels of the gold vocabulary ($\mathbb{Y}^G$) and the remaining words tokenized from the user comments ($X_i$). Specifically, the candidate labels for the $i^{th}$ song are defined as $\mathbb{Y}^G \cup X_i - \mathcal{Y}_i^G$. The candidates with high confidence scores are picked out as pseudo-labels for the following fine-tuning. Note that, such pseudo-labels (e.g., "**hope**") can barely contribute to OOV labels missed by the annotation. This is because the classifier always assigns pseudo-labels to an instance by borrowing gold labels from other instances with similar features. This also explains why the candidates not in the vocabulary often get very low scores. An exception is the OVV candidate label "**traditional**", which gets a relatively high score. This is because "**traditional**" is semantically similar to the gold label "**classic**". As a result, though the inferred pseudo-labels contribute to missed supervision, they cannot virtually form a complete set of labels for a song.

### 3.3  Joint Score Function

We believe the above limitation will be largely mitigated if we "teach" the classifier what has not been acquired by it before, namely, fine-tuning the classifier on more diverse labels that are unknown at the previous training phase. Specifically, given the remaining candidate labels (e.g., $X_i - \mathcal{Y}_i^G - \{$"**hope**"$\}$), we design a joint score function to measure the diversity (statistical importance and semantic novelty) and validness (practical value and discrimination ability) of a candidate label. This function computes a joint value of the following scores.

***Statistical importance score*** *(SI)*  Following the tradition of information searching, we use TF-IDF to measure the statistical importance score of a candidate label. The higher the *SI* score, the more important the candidate label is to a song from the view of data distribution.

***Semantic novelty score*** *(SN)*  If only an expert says a given candidate label is semantically different from the existing ones, it may be subjective. However, if many other experts also agree with that, the candidate label is most likely to be a semantic novel label. The analogy to this, we conduct $m$ groups of $K$-means clustering on the labels to compute the semantic novelty score of a candidate label against the existing ones. In this way, we get $m$ "experts", each of which has its own taxonomy and opinions about the labels. The confidence of $y_c$ being a semantic novel label is defined as

$$SN = \frac{1}{2} * \sum_{i=1}^{m} \frac{1 - \min\left[\mathcal{D}(y_c, C_i^1), \mathcal{D}(y_c, C_i^2), ..., \mathcal{D}(y_c, C_i^K)\right]}{m}, \tag{3}$$

where $C_1^K$ means the $K^{th}$ cluster of the $i^{th}$ "expert" and $\mathcal{D}(y_c, C_i^K)$ calculates the cosine similarity between $y_c$ and the center of $C_i^K$.
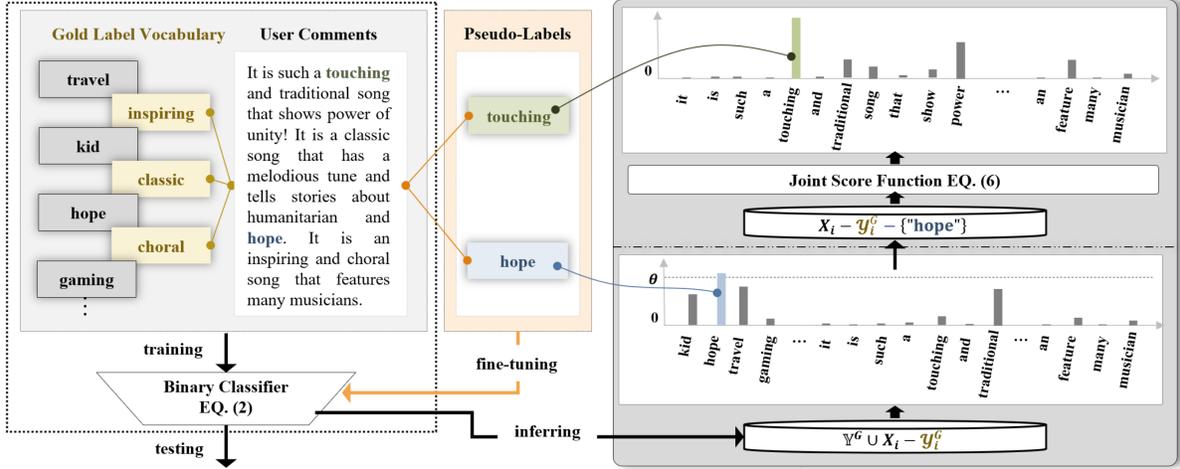
**Figure 3: The proposed iterative framework (`DiVa`). It consists of two key components, i.e., a binary classifier and a joint score function. The binary classifier is first trained on user comments and gold labels and then fine-tuned on pseudo-labels inferred from the pre-trained classifier and a novel joint score. This process continues until the classifier can barely produce any new pseudo-labels. During testing, the labels for a song are directly inferred from user comments via the optimized classifier.**

***Practical value score** (PV)*  Similar to $SN$, if all experts agree that a candidate label is not suitable to be used as a label, then we think this candidate label is not valid enough from the point of practical values. Here, we treat the songs as the "experts" and aim to penalize candidates with low practical values (e.g., below a pre-defined threshold $\tau$). Specifically, the practical value score depends on the confidence scores of the candidate label predicted on all songs and is defined as

$$PV = \begin{cases} 1, & \sum_{i=1}^{N} \frac{\mathcal{F}(X_i, y_c)}{N} \geq \tau, \\ 0, & \sum_{i=1}^{N} \frac{\mathcal{F}(X_i, y_c)}{N} < \tau. \end{cases} \quad (4)$$

***Discrimination ability score** (DA)*  We believe a valid label, even presented in many songs, should be able to distinguish different songs. As such, we want to penalize candidates with low discrimination ability. Specifically, we use the coefficient of variation to describe the discrimination ability of a candidate label and define the $DA$ score as

$$DA = \begin{cases} 1, & \frac{\sigma(y_c)}{\mu(y_c)} \geq \tau, \\ 0, & \frac{\sigma(y_c)}{\mu(y_c)} < \tau, \end{cases} \quad (5)$$

where $\tau$ is the same threshold used in $PV$, $\sigma(y_c)$ and $\mu(y_c)$ computes the standard deviation and mean of the number of $y_c$ appearing in songs, respectively.

Intuitively, we define the joint score function of the candidate label $y_c$ based on the above scores as

$$\mathcal{G}(y_c) = SI * SN * PV * DA. \quad (6)$$

The joint score function describes a unified objective to harvest labels with high semantic novelty and semantic novelty and avoid those with low practical value and discrimination ability. The candidate labels that get high joint scores are then selected as pseudo-labels, which involve OOV labels (e.g., "**touching**") and a few new supervisions missed before. These pseudo-labels, together with those selected by the pre-trained classifier, are used to fine-tune the classifier.

### 3.4 Iteration

Towards fast and stable optimization, we adopt the negative sampling strategy at the training/fine-tuning phase of the binary classifier. To construct a negative sample for a given song, we randomly select a label from the candidate labels except for the gold and pseudo ones and pair it with the user comments on the song. Note that, even if a candidate label, which may be a suitable label for a song, is accidentally used as a negative sample in the current iteration of training, the bad effect can be ignored as it will be selected as a pseudo-label via the following inference and be used to fine-tune the classifier. Besides, to balance the importance of rare and frequent pseudo-labels, we further leverage the subsampling strategy [17].

Besides reaching the maximum, the iterations stop when the classifier can barely infer any new pseudo-labels. A straightforward solution is to define a threshold (e.g., 50) — if the number of new labels inferred by the binary classifier is below the threshold, we stop the iterations. However, it is hard to find a perfect threshold for different datasets, for example, 50 may be too large for datasets with a small-sized label vocabulary and too small for datasets with a large-sized label vocabulary. As a remedy, we estimate the stopping condition by evaluating the classifier's ability to infer new labels on long-tail labels, e.g., the performance on the PSP and PSnDCG metrics [10]. This is reasonable as the tail labels hold a big part of the labels missed by the gold ones.

## 4 EXPERIMENT

### 4.1 Experiment Settings

***Dataset***  To perform a systematic data-driven study, we gather a corpus of 16,372 songs, each of which consists of user comments and gold labels summarized by experts from the information retrieval department of a Chinese online music platform. The statistics are shown in Table 1. We then randomly select a testing set of 1000 songs with more than 10 user comments. Three senior experts from the same department are hired to annotate the testing set. At each

**Table 1: The number of labels w.r.t., the original corpus, training set and two testing set**

| Dataset | label | labels / song | | | user comments / song | | | words / user comment | | |
|---|---|---|---|---|---|---|---|---|---|---|
| (Number of songs) | vocabulary size | Avg. | Min. | Max. | Avg. | Min. | Max. | Avg. | Min. | Max. |
| Corpus (16,372) | 820 | 5.82 | 1 | 31 | 38.83 | 6 | 60 | 44.05 | 9 | 335 |
| Training (15,372) | 816 | 5.82 | 1 | 31 | 38.84 | 6 | 60 | 44.04 | 9 | 335 |
| test-1 (1,000) | 425 | 5.78 | 1 | 27 | 38.76 | 11 | 60 | 44.15 | 9 | 303 |
| test-2 (1,000) | 3312 | 19.41 | 4 | 44 | 38.76 | 11 | 60 | 44.15 | 9 | 303 |

time, a worker will be shown a song with its metadata (title, artist, etc), user comments, and candidate labels (stopwords are removed in advance), he is required to select proper labels for the song from the candidates after carefully go through the metadata and read the user comments. If needed, he can also use his own knowledge and searching necessary information online. We give a training session to the workers to help them fully understand our annotation requirements. The annotation tasks are released to the workers via an in-house website. On average, it takes about 20 minutes for a worker to annotate a song. Each song is randomly assigned to two workers. For the $i^{th}$ song, we calculate the Kappa coefficient score between the different workers' annotations. If the score is higher than 75%, all selected candidates together with the gold labels are used as the complete set of labels $\mathcal{Y}_i$. Otherwise, we let the two workers discuss their disagreements and come up with the agreed results. After the annotation, we get two testing sets — one has gold labels for every song (denoted as test-1) and the other one has the complete set of labels for every song (denoted as test-2). On average, there are about 5.8 gold labels for every song and about 19.4 labels for a song in test-2.

***Evaluation Metrics***   Following the tradition, we evaluate DiVa on the commonly used multi-label classification metrics, which involve precision (P), recall (R), F1, and nDCG, and extreme multi-label learning metrics, which involve normalized PSP and PSnDCG [10]. The PSP and PSnDCG scores are not applicable for test-2 as we assume that each song in test-2 has a complete set of labels. Inspired by Tandon et al. [24] and Simig et al. [22], to eliminate the inaccurate mismatch caused by synonym labels, we also employ the soft matching metrics. Given the predicted labels (denoted as $\mathcal{Y}_i^* = \{y_j^*\}_{j=1}^{|\mathcal{Y}_i|}$) and the gold labels (denoted as $\mathcal{Y}_i^G = \{y_k\}_{k=1}^{|\mathcal{Y}_i^G|}$) for the $i^{th}$ song, the soft precision is defined as

$$SP = \sum_{j=1}^{|\mathcal{Y}_i^*|} \frac{\max\left[\mathcal{D}(y_j^*, y_1), ..., \mathcal{D}(y_j^*, y_{|\mathcal{Y}_i^G|})\right]}{|\mathcal{Y}_i^*|}, \qquad (7)$$

where $\mathcal{D}(y_j^*, y_1)$, with the same definition in EQ. (3), returns the cosine similarity between $\vec{y}_j^*$ and $\vec{y}_1$. The soft recall is defined as

$$SR = \sum_{k=1}^{|\mathcal{Y}_i^G|} \frac{\max\left[\mathcal{D}(y_k, y_1^*), ..., \mathcal{D}(y_k, y_{|\mathcal{Y}_i^*|}^*)\right]}{|\mathcal{Y}_i^G|}. \qquad (8)$$

Accordingly, the soft F1 score is defined as $SF1 = \frac{2*SP*SR}{SP+SR}$. Note that, the soft matching metrics are not suitable for methods that have represented labels via their own models. In line with EQ. (1), we also want to measure whether the predicted label can sufficiently make up for the information missed by the gold labels. As such, we design a coverage score on test-2, which is defined as Jaccard similarity between predicted labels and the complete set of labels,

i.e., $Coverage = \frac{1}{N} * \sum_{i=1}^{N} \frac{\mathcal{Y}_i^* \cap \mathcal{Y}_i}{\mathcal{Y}_i^* \cup \mathcal{Y}_i}$. The coverage score for the gold labels in test-2 is 0.284. For all metrics, higher values indicate better performances. Particularly, we place emphasis on the F1/SF1 and coverage scores, which are directly related to our ultimate goal, i.e., the ability to obtain a complete set of labels for each song.

***Baselines***   Towards a thorough comparison, we compare DiVa with the following baseline methods.

- TF-IDF, which identifies labels based on the statistic features of the candidate label. It also contributes to the statistical importance score ($SI$) in our novel joint score function defined in EQ. (6).
- Multi-label Classification (MLC), which takes user comments as input and outputs a label vocabulary-sized vector.
- LightXML [11], which is a tree-based extreme multi-label learning method that uses dynamic negative sampling to boost the model performance on massive labels.
- Our binary classifier, which is define by EQ. (2).
- nnPU [12], which avoids biased noises in binary classification via a non-negative risk estimator. We deploy this strategy based on our binary classifier.
- GROOV [22], which labels songs through a generative model.
- ChatGPT [19], which can be instructed to annotate text data via prompts [9].
- Noisy Student Training (NST) [27], which trains a classifier iteratively using pseudo-labels predicted on unlabeled instances. We deploy the NST framework on both MLC and binary classifiers, denoted as NST w/ MLC and NST w/ (EQ. 2), respectively.
- Two variants of DiVa: DiVa-static, which, before starting the iterations, combines the pseudo-labels inferred from the initialized binary classifier and the joint score function as labels; DiVa-light, which only uses the pseudo-labels inferred from the current iteration to fine-tune the binary classifier.

***Implementation Details***   The DiVa framework and the neural baselines are trained/fine-tuned on 8 Nvidia A6000 GPUs using Adam optimizer with a learning rate=0.001. We use the pre-trained XLNet-base model provided by Yang et al. [29]) to represent the user comments. As for candidate labels, following Gupta and Jaggi [8], we get the X2Static embeddings based on the XLNet-base model and CBOW model. All hyperparameters are tuned on the training set. Each baseline is optimized to gain its best performance.

## 4.2   Results and Observations

As mentioned in Sec. 2, our ultimate goal is to predict the complete set of labels for a song. Table 2 demonstrates that the proposed DiVa framework is the most desired solution to the goal. Because it obtains superior results (best F1 score, best coverage score, and second best SF1 score) on test-2, suggesting that DiVa, compared

**Table 2: Performances of all methods on test-1/test-2 w.r.t., precision (P), recall (R), F1, nDCG, normalized PSP [10], normalized PSnDCG [10], soft precision (SP), soft recall (SR), soft F1 (SF1), and the coverage score. We mark the values indicating the best performance in bold and the values indicating the second-best performance with underlines.**

| Method | MLC | | | | XML | | Soft matching | | | Coverage |
|---|---|---|---|---|---|---|---|---|---|---|
| | P | R | F1 | nDCG | PSP | PSnDCG | SP | SR | SF1 | |
| TF-IDF | 0.015/0.082 | 0.090/0.346 | 0.021/0.119 | 0.046/0.247 | 0.061/- | 0.057/- | 0.064/0.245 | 0.164/0.415 | 0.081/0.282 | 0.139 |
| MLC | 0.044/0.063 | 0.013/0.007 | 0.018/0.012 | 0.044/0.062 | 0.152/- | 0.137/- | 0.050/0.071 | 0.019/0.012 | 0.025/0.020 | 0.283 |
| LightXML [11] | 0.405/0.819 | 0.488/0.271 | 0.411/0.396 | 0.539/<u>0.834</u> | 0.449/- | 0.418/- | <u>0.441</u>/<u>0.877</u> | 0.521/0.316 | 0.449/0.454 | 0.395 |
| EQ. (2) | <u>0.443</u>/<u>0.836</u> | 0.691/0.360 | **0.510**/0.492 | **0.686/0.858** | 0.653/- | **0.594**/- | **0.475/0.893** | 0.713/0.404 | **0.544**/0.546 | 0.420 |
| nnPU [12] w/ EQ. (2) | 0.070/0.189 | **0.980**/<u>0.737</u> | 0.128/0.298 | <u>0.654</u>/0.690 | **0.931**/- | <u>0.513</u>/- | 0.126/0.318 | **0.982**/<u>0.766</u> | 0.220/0.446 | 0.178 |
| GROOV [22] | **0.512/0.844** | 0.412/0.200 | 0.426/0.310 | 0.295/0.244 | 0.327/- | 0.203/- | - | - | - | 0.206 |
| ChatGPT [19] | 0.128/0.270 | 0.141/0.098 | 0.110/0.129 | 0.130/0.177 | 0.105/- | 0.105/- | - | - | - | 0.239 |
| NST [27] w/ MLC | 0.096/0.155 | 0.022/0.013 | 0.034/0.024 | 0.098/0.155 | 0.189/- | 0.185/- | 0.105/0.184 | 0.029/0.021 | 0.044/0.037 | 0.286 |
| NST [27] w/ EQ. (2) | 0.337/0.775 | 0.746/0.483 | <u>0.444</u>/0.584 | 0.603/0.802 | 0.645/- | 0.462/- | 0.378/0.841 | 0.765/0.521 | <u>0.487</u>/0.633 | 0.484 |
| DiVa-static | 0.175/0.566 | <u>0.897</u>/**0.811** | 0.281/<u>0.657</u> | 0.462/0.703 | <u>0.835</u>/- | 0.404/- | 0.222/0.707 | <u>0.906</u>/**0.831** | 0.344/**0.757** | 0.517 |
| DiVa-light | 0.231/0.694 | 0.722/0.629 | 0.331/0.648 | 0.564/0.753 | 0.644/- | 0.451/- | 0.274/0.784 | 0.745/0.663 | 0.383/0.709 | <u>0.550</u> |
| DiVa | 0.239/0.727 | 0.716/0.631 | 0.341/**0.665** | 0.569/0.782 | 0.643/- | 0.444/- | 0.282/0.819 | 0.739/0.665 | 0.392/<u>0.724</u> | **0.570** |

to other methods, can cover more diverse information using fewer labels. We believe this is achieved by the iterative training/fine-tuning of the gold labels and pseudo-labels inferred from the pre-trained binary classifier and the novel joint score function. We also make the following observations.

- When evaluating on test-1, EQ. (2) obtains better F1 and SF1 scores yet much worse PSP score than of DiVa-static. These results indicate that if we equip EQ. (2) with the joint score function, it is more powerful to handle an incomplete set of labels. This assumption is further demonstrated by the evaluation on test-2, where DiVa-static has higher F1, SF1, and coverage scores than EQ. (2).

- The DiVa-light method gets worse worse F1, SF1, coverage scores than DiVa. We conjecture this is because only using pseudo-labels to fine-tune the classifier makes the classifier over-fit to pseudo-labels during the iterations. As a result, though the binary classifier learned from DiVa-light has the ability to produce new information missed by the gold labels, it loses the ability to produce gold labels.

- From the comparison between NST [27] and DiVa, we notice that NST always works better than DiVa on test-1 but works worse than DiVa on test-2. This is because the goal of self-training methods is to approach the mappings to gold labels and thus these methods, including NST, are good at handling missed supervision. Whereas, the DiVa framework forces the classifier to learn mappings to OOV labels as well as missed supervisions and thus is better at forming complete sets of labels for songs.

- Compared with DiVa and other methods, nnPu [12] w/ EQ. (2) gets the best PSP scores on test-1. Thus, it has a big chance to get a good performance on test-2. However, its coverage score is only slightly better than the worst one (TF-IDF) and much worse than the best one (DiVa). This is because the labels produced by nnPu contain too much noise. Another piece of evidence is that nnPu gets very high R and SR scores but rather low P and SP scores on test-1 and test-2.

- We further observe that a better understanding of labels can greatly improve performance. Specifically, the methods using indexes to represent labels (MLC and NST w/ MLC) always work
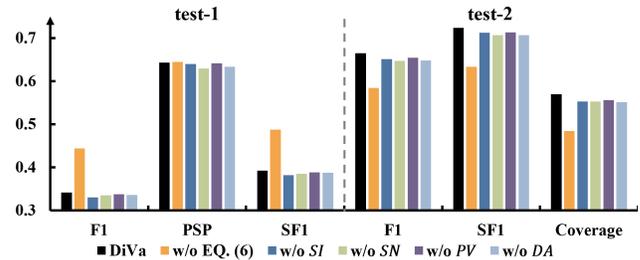


**Figure 4: Performances w.r.t., different scores of EQ. (6) on test-1 (left) and test-2 (right)**

worse than those using tree (lightXML [11]) or embeddings to represent labels (GROOV [22] and DiVa).

- To better instruct ChatGPT [19], we have tried all off-the-shelf prompts that are designed for text annotation and utilize the prompt with the best performances [6, 9, 18] to predict labels for songs from user comments. However, ChatGPT doesn't produce promising results as it has done in other tasks. We find a possible explanation in Gupta and Jaggi [8] that contextual language models are not always welcome by word-level interpretability.

- Compared to conventional MLC, our binary classifier defined by EQ. (2) performs well not only on test-1 but also test-2. This is because instead of using the labels in the gold label vocabulary, we use words from use comments as candidate labels, the number of which is less than the gold vocabulary size. So it can learn more reliable mappings in a lower-dimensional space. Besides, we use X2Static [8] to encode labels into vectors. So the labels have a better correlation with each other and user comments. This is also why EQ. (2), as a supervised method, can produce some labels beyond the gold label vocabulary.

### 4.3 Ablation Study on different scores of EQ. (6)

As shown in Figure 4, we perform an ablation study to investigate the effect of the novel joint score function (EQ. (6)) and the effect of $SI$, $SN$, $PV$, $DA$ scores in measuring a candidate label. Note that, if we remove the novel joint score function from DiVa during the iterations, we will get NST [27] w/o EQ. (6). Figure 4 presents the performances with different joint score functions on test-1 and
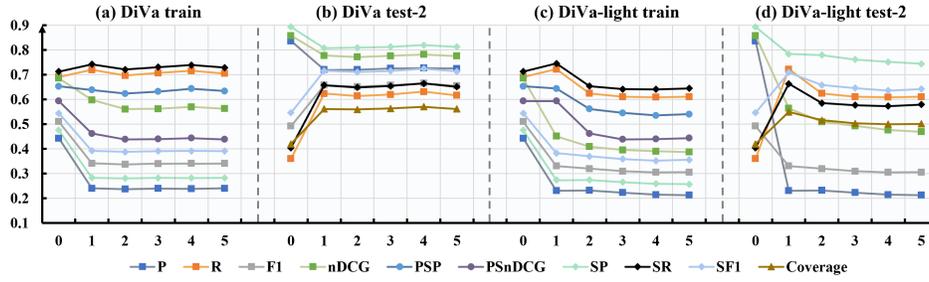
**Figure 5: Performances of (a) `DiVa` on the training set, (b) `DiVa` on `test-2`, (c) `DiVa-light` on the training set, and (d) `DiVa-light` on `test2` over iterations.**

`test-2`, respectively. From the comparisons between `DiVa` against `DiVa` w/o *SI*, w/o *SN*, w/o *PV*, and w/o *DA*, we can see that any one of the scores has little impact on `DiVa`. However, when these scores are assembled in the joint score function, there presents a strong influence on `DiVa`. Specifically, we see a big rise (about 0.1) of F1 and SF1 scores when comparing `DiVa` and `DiVa` w/o EQ. (6) on `test-1`. Conversely, we see a big drop (about 0.08) of all metrics when comparing two methods on `test-2`. This is because EQ. (6) is designed to find more diverse and valid pseudo-labels beyond gold labels to fine-tune the binary classifier. Thus, these pseudo-labels are "noises" for the ground truth values in `test-1` but the right answers for the ground truth values in `test-2`, resulting in the rise and drop in Figure 4.

### 4.4 Performance over iterations

We set the iteration number of `DiVa` and `DiVa-light` based on their results, mainly the PSP and PSnDCG scores, on the training set. Note that, the classifier at the $0^{th}$ iteration is only trained on gold labels. As we can see from Figure 5 (a) and Figure 5 (c), the best iteration number for `DiVa` is 4 and that for `DiVa-light` is 1. For example, the results of `DiVa` shown in Table 2 are obtained from the binary classifier fine-tuned at the $4^{th}$ iteration. For auxiliary analysis, we also report the results of `DiVa` and `DiVa-light` on `test-2` produced by the binary classifier optimized in every iteration in Figure 5 (b) and Figure 5 (d). In keeping with the training set, `DiVa` also gets its best performance at the $4^{th}$ iteration. We observe a trend that the precision scores are becoming smaller gradually with the increase of iterations on `test-1` and `test-1`. This is because, at each iteration, the classifier is fine-tuned on more pseudo-labels, many of which are regarded as "noises" by the gold ones. These "noises" has a negative effect on the precision scores. However, the decrease presented in Figure 5 (c) is much sharper than that presented in Figure 5 (a). This is because, as mentioned in Sec. 4.2, the classifier, only fine-tuned on pseudo-labels, is overfitting to the pseudo-labels and loses the ability to predict gold labels. Accordingly, it cannot get good performances on `test-2`, as shown in Figure 5 (d).

### 4.5 Case Study

We present two examples of two songs (*The Last Of The Mohicans* and *Mermaid*) with gold labels, complete set of labels and labels produced by different methods in Table 3. Among all methods, our `DiVa` framework can cover most labels of the complete set with the least bad labels. As explained in Sec. 4.2, the supervised methods (lightXML, our binary classifier EQ. (2), and NST w/ EQ. (2))

have the ability to infer missed supervision because of their extra attentions on labels. Surprisingly, although nnPU w/ EQ. (2) and GROOV have produced many labels, most of them are bad ones. We also have an interesting observation that, ChatGPT, even with carefully designed prompts, tends to model the automated music labeling as a summarization task and tends to produce phrases rather than words.

## 5 RELATED WORK

***Extreme Multi-label Learning***   The widely used XML methods can be divided into three types [15, 26]: the embedding-based methods try to map both the features and the labels into a joint low-dimensional space [7, 23], the tree-based methods partition labels or instances into a tree [4, 11, 30], and one-vs-all methods learns a binary classifier for each label separately [2, 31]. Although such XML methods can learn more reliable mappings from massive labels, they require full label coverage and full supervision to train the classifiers. Note that the binary classifier used in `DiVa` differs from the one-vs-all XML methods in two aspects. First, instead of learning a classifier for each label, we only design one binary classifier for all the concatenation of user comments and candidate labels. This largely reduces the number of parameters that need to be learned. Second, instead of using the full label set as candidates, we only use the gold labels and the words from user comments as the candidate labels of the song in the training phase. This largely reduces the number of candidates that need to be estimated.

***Open Vocabulary Classification***   Simig et al. [22] introduces the open vocabulary XML classification for the first time. It also proposes GROOV, a generative model that can tag textual context with a set of labels from an open vocabulary. Although the GROOV model can enlarge the number of labels, the newly generated labels share a lot of similarities with the gold ones in the semantic space. Another work on open vocabulary classification is Xiong et al. [28], which is also the only work, as far as we know, that concerns both incomplete label coverage and incomplete supervision. However, this approach is more suitable for instances with a small number of labels and can hardly apply to automated music labeling. Because it trains the model on too many noise pseudo labels that could lead the user searching to wrong songs with big chances. It also has a strict assumption that all instances must have label-like (e.g., title) phrases. This is not always satisfied in real-world scenarios.

***Self-training***   The self-training methods work by training a classifier iteratively by assigning pseudo-labels to unlabeled instances [1,

**Table 3: Examples of two songs with gold labels, complete set of labels and labels produced by different methods.**

| Title: *The Last of The Mohicans* | Artist: Alexandro Querevalú |
|---|---|
| Gold labels | soul-touching, aesthetic, trouble, passionate, sad |
| Complete set of labels | melodious, proud, helpless, life, era, soul-touching, male, surging, America, shot, future, world, moving, the sound of nature, turbulent, passionate, sad, nation, lonely, touching, story, aesthetic, trouble, slave, war |
| LightXML [11] | aesthetic, passionate, single, touching, lonely, moving, world |
| EQ. (2) | passionate, sad, single, touching, story, lonely, moving, world, soul-touching, aesthetic |
| nnPU [12] w/ EQ. 2 | passionate, war, proud, sad, single, life, story, cry, moving, world, future, video, era, touching, witness, peace, eat, pure, hope, the sound of nature, scared, profit, street, trouble, remember, poor, the Orient |
| GROOV [22] | healing, peaceful, memory, classic, love song, love, affectionate, single, story, galaxy, youth, sunshine, self-improvement, enthusiasm, aesthetics, light music, adoration, bgm, soul-touching, archaic rhyme |
| ChatGPT [19] | original ecological music, folk music, humanistic concern, historical significance, reflection on human civilization, cultural diversity, tragic, struggle and frustration |
| NST [27] w/ EQ. (2) | world, story, single, sad, passionate, moving, touching, lonely, aesthetic |
| DiVa | moving, sad, single, touching, story, lonely, passionate, aesthetic, turbulent, world, male, slave, Spanish, play wind instruments, the sound of nature, nation, era, soul-stirring, shock, good faith, sad, war, democratic, soul-touching, life, surging, melodious |
| Title: *Mermaid* | Artist: Skott |
| Gold labels | meet |
| Complete set of labels | deep love, spacious, peaceful, chilly, appeal, thunder, clip, forest, bewitching, ethereal, meet, beautiful, happy, affectionate, pirate, story, stunning, love |
| LightXML [11] | affectionate, deep love, beautiful |
| EQ. (2) | love, deep love, affectionate, story, beautiful |
| nnPU [12] w/ EQ. 2 | affectionate, deep love, story, video, beautiful, stunning, share, peaceful, forest, ethereal, hope, happy, vocal cords, style, worth, anticipation, tone colour, meet, clean, treasure, wish, collect, human heart |
| GROOV [22] | healing, galaxy, memory, peaceful, love song, affectionate, love, sad, classic, youth, sunshine, meet, single, bgm, rhythm, passionate, soul-touching |
| ChatGPT [19] | famous singer Thunder sister, mermaid, sea monster, Skott, Bjork, Huawei, Edited video of Ice and Snow 2 |
| NST [27] w/ EQ. (2) | story, love, deep love, affectionate, ethereal, forest, beautiful, happy, peaceful |
| DiVa | beautiful, love, story, affectionate, deep love, pirate, ethereal, forest, bewitching, spacious, chilly, on the road, peaceful, appeal, happy, stunning, meet |

25, 32]. For example, Noisy Student Training [27] has three steps: train a teacher classifier on labelled instances, predict the pseudo-labels on unlabeled instances, and train a student classifier on both labelled and pseudo-labelled data. However, the self-training classifier tends to assign similar labels to instances with similar features. This is not helpful to capture the complete picture of user searching preferences on a single song.

***Positive-Unlabeled (PU) Learning*** The incomplete supervision issue can be also found in the PU learning setting, where only some of the positive instances are labeled [3]. One of the most popular PU learning methods is nnPU [12], which treats the unlabeled instances as negative instances and weights the unlabeled data via a non-negative risk estimator to avoid biased noises. Generally, such methods are designed for binary classification settings. In this paper, we can change the multi-label classification task to a binary classification task by pairing the user comments and a candidate label as an instance. This makes it possible to apply the PU learning methods in multi-label classification settings. However, this doesn't mean that such methods are suitable for the proposed task. Because these methods aim to get better performances on the seen supervision instead of digging the unseen supervision. Besides, they require the prior class distribution to optimize the classifier.

***Active learning*** The proposed iterative framework (DiVa) is also slightly relevant to active learning, which selects instances for annotating and training in an interactive process [20, 21]. However, instead of reducing the cost during the annotation stage, our work aims to make the best of the labelled data after the annotation stage. Besides, instead of hiring experts for annotation, we obtain pseudo-labels from a joint score function.

## 6 CONCLUSION

In this paper, we call attention to the information missed by the gold labels and highlight the importance of the complete set of labels for a song towards sufficient music searching. After a detailed investigation of the missing labels, user comments, and current studies, we'd like to take a step forward and study automated music labeling in an essential yet under-exploring setting — given limited gold labels, the model is required to harvest more diverse and valid labels from the user comments to form a complete set of labels for the song. Further, we develop an iterative framework (DiVa) to learn a desired classifier for the proposed automated music labeling task. We also design a novel joint score function that can offer more diverse and valid pseudo-labels to fine-tune the classifier. Experiments on two testing sets reveal the superiority of DiVa in covering as much information with as few labels. Besides, the comparison with other iterative methods demonstrates the effectiveness of the novel joint score function in harvesting more diverse and valid labels. We hope our work can inspire future research on automated music labeling and shed light on multi-label learning with missing information in real-world scenarios.

# REFERENCES

[1] Massih-Reza Amini, Vasilii Feofanov, Loic Pauletto, Emilie Devijver, and Yury Maximov. 2022. Self-training: A survey. *arXiv preprint arXiv:2202.12040* (2022).

[2] Rohit Babbar, Ioannis Partalas, Eric Gaussier, and Massih-Reza Amini. 2013. On flat versus hierarchical classification in large-scale taxonomies. In *Proceedings of the 26th International Conference on Neural Information Processing Systems-Volume 2*. 1824–1832.

[3] Jessa Bekker and Jesse Davis. 2020. Learning from positive and unlabeled data: A survey. *Machine Learning* 109 (2020), 719–760.

[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 *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 3163–3171.

[5] Chongming Gao, Wenqiang Lei, Xiangnan He, Maarten de Rijke, and Tat-Seng Chua. 2021. Advances and challenges in conversational recommender systems: A survey. *AI Open* 2 (2021), 100–126.

[6] Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. Chatgpt outperforms crowd-workers for text-annotation tasks. *arXiv preprint arXiv:2303.15056* (2023).

[7] Chuan Guo, Ali Mousavi, Xiang Wu, Daniel Holtmann-Rice, Satyen Kale, Sashank Reddi, and Sanjiv Kumar. 2019. Breaking the glass ceiling for embedding-based classifiers for large output spaces. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. 4943–4953.

[8] Prakhar Gupta and Martin Jaggi. 2021. Obtaining Better Static Word Embeddings Using Contextual Embedding Models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 5241–5253.

[9] Fan Huang, Haewoon Kwak, and Jisun An. 2023. Is ChatGPT Better than Human Annotators? Potential and Limitations of ChatGPT in Explaining Implicit Hate Speech. In *Companion Proceedings of the ACM Web Conference 2023*. Association for Computing Machinery, New York, NY, USA, 294–297. https://doi.org/10.1145/3543873.3587368

[10] Himanshu Jain, Yashoteja Prabhu, and Manik Varma. 2016. 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*. 935–944.

[11] 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 *Proceedings of the AAAI Conference on Artificial Intelligence*. 7987–7994.

[12] Ryuichi Kiryo, Gang Niu, Marthinus C du Plessis, and Masashi Sugiyama. 2017. Positive-unlabeled learning with non-negative risk estimator. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*. 1674–1684.

[13] Hongru Liang, Wenqiang Lei, Paul Yaozhu Chan, Zhenglu Yang, Maosong Sun, and Tat-Seng Chua. 2020. Pirhdy: Learning pitch-, rhythm-, and dynamics-aware embeddings for symbolic music. In *Proceedings of the 28th ACM international conference on multimedia*. 574–582.

[14] Hongru Liang, Haozheng Wang, Jun Wang, Shaodi You, Zhe Sun, Jin-Mao Wei, and Zhenglu Yang. 2018. JTAV: Jointly learning social media content representation by fusing textual, acoustic, and visual features. In *Proceedings of the 27th International Conference on Computational Linguistics*. 1269–1280.

[15] Weiwei Liu, Haobo Wang, Xiaobo Shen, and Ivor W Tsang. 2021. The emerging trends of multi-label learning. *IEEE transactions on pattern analysis and machine intelligence* 44, 11 (2021), 7955–7974.

[16] Melese Mihret and Muluneh Atinaf. 2019. Sentiment Analysis Model for Opinion-ated Awngi Text: Case of Music Reviews. In *Proceedings of the 2019 Workshop on Widening NLP*. Association for Computational Linguistics, Florence, Italy, 49–51. https://aclanthology.org/W19-3617

[17] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Proceedings of the 26th International Conference on Neural Information Processing Systems-Volume 2*. 3111–3119.

[18] Anders Giovanni Møller, Jacob Aarup Dalsgaard, Arianna Pera, and Luca Maria Aiello. 2023. Is a prompt and a few samples all you need? Using GPT-4 for data augmentation in low-resource classification tasks. *arXiv preprint arXiv:2304.13861* (2023).

[19] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* 35 (2022), 27730–27744.

[20] Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Brij B Gupta, Xiaojiang Chen, and Xin Wang. 2021. A survey of deep active learning. *ACM computing surveys (CSUR)* 54, 9 (2021), 1–40.

[21] Burr Settles. 1995. Active Learning Literature Survey. *Science* 10, 3 (1995), 237–304.

[22] Daniel Simig, Fabio Petroni, Pouya Yanki, Kashyap Popat, Christina Du, Sebastian Riedel, and Majid Yazdani. 2022. Open Vocabulary Extreme Classification Using

[23] Yukihiro Tagami. 2017. AnneXML: Approximate nearest neighbor search for extreme multi-label classification. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*. 455–464.

[24] Niket Tandon, Keisuke Sakaguchi, Bhavana Dalvi, Dheeraj Rajagopal, Peter Clark, Michal Guerquin, Kyle Richardson, and Eduard Hovy. 2020. A Dataset for Tracking Entities in Open Domain Procedural Text. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 6408–6417.

[25] Jesper E Van Engelen and Holger H Hoos. 2020. A survey on semi-supervised learning. *Machine learning* 109, 2 (2020), 373–440.

[26] Tong Wei, Zhen Mao, Jiang-Xin Shi, Yu-Feng Li, and Min-Ling Zhang. 2022. A Survey on Extreme Multi-label Learning. *arXiv preprint arXiv:2210.03968* (2022).

[27] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. 2020. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10687–10698.

[28] Yuanhao Xiong, Wei-Cheng Chang, Cho-Jui Hsieh, Hsiang-Fu Yu, and Inderjit Dhillon. 2022. 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*. 5455–5468.

[29] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems* 32 (2019).

[30] Ronghui You, Zihan Zhang, Ziye 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. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. 5820–5830.

[31] Min-Ling Zhang, Yu-Kun Li, Xu-Ying Liu, and Xin Geng. 2018. Binary relevance for multi-label learning: an overview. *Frontiers of Computer Science* 12 (2018), 191–202.

[32] Yang Zou, Zhiding Yu, Xiaofeng Liu, BVK Kumar, and Jinsong Wang. 2019. Confidence regularized self-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5982–5991.

Generative Models. In *Findings of the Association for Computational Linguistics: ACL 2022*. 1561–1583.

## A IMPLEMENTATION DETAILS

### A.1 Baselines

***TF-IDF*** We deploy the classic TF-IDF, which is defined as

$$\text{TF-IDF}(y_c, d_i) = \frac{u_{y_c}}{\sum_{d_i} u_{y_c}} \times \log \frac{N}{|\{d_j : y_c \in d_j\}_{j=1}^N|}, \qquad (9)$$

$u_{y_c}$ and $\sum_{d_i} u_{y_c}$ denote the number of occurrences of $y_c$ in the current user song comments $d_i$ and the total number of occurrences of all words tokenized from the comments, respectively. The left side of $*$ indicates the relative frequency of the candidate label $y_c$ within the current song comments, and the right side of $*$ is related to the number of songs with $y_c$ in their user comments.

***Multi-label Classification (MLC)*** We use the user comments of a song to predict a vector, whose dimension is equal to the size of gold label vocabulary. Each element of the vector indicates the probability of $y_c$ being an accurate label for the $i^{th}$ song. The multi-label classifier is formulated as

$$F(d_i, \mathbb{Y}^G) = sigmoid(\vec{d_i}), \qquad (10)$$

where $\vec{d_i}$ is encoded from XLNet [29]. The objective of the classifier is to minimize the sum of the binary cross entropy between the predicted and ground truth vectors.

***LightXML [11]*** We use the codes released by the original paper https://github.com/kongds/LightXML.

***nnPU [12] w/ EQ. (2)*** The original loss function of EQ. (2 is defined by the binary cross entropy, i.e.,

$$\mathcal{L} = \sum_{i=1}^N \sum_{y_c} BCE(y^*, \mathcal{F}(d_i, y_c))$$

$$= -\sum_{i=1}^N \sum_{y_c} y^* log(\mathcal{F}(d_i, y_c)) + (1 - y^*) \times \log(1 - \mathcal{F}(d_i, y_c)), \qquad (11)$$

where $y^* = 1$ if $y_c$ is the gold label of the $i^{th}$ song, $y^* = 0$ if $y_c$ is a sampled negative label of the $i^{th}$ song. We replace this loss function with the nnPU loss [12], which is defined as

$$\widetilde{R}_{pu}(\mathcal{F}) = \pi_p \hat{R}_p^+(\mathcal{F}) + \max\{0, \hat{R}_u^-(\mathcal{F}) - \pi_p \hat{R}_p^-(\mathcal{F})\},$$

$$\hat{R}_p^+(\mathcal{F}) = \frac{1}{|\chi_p|} \sum_{d_i, y_c \in \chi_p} BCE(1, \mathcal{F}(d_i, y_c)),$$

$$\hat{R}_u^-(\mathcal{F}) = \frac{1}{|\chi_u|} \sum_{d_i, y_c \in \chi_u} BCE(0, \mathcal{F}(d_i, y_c)), \qquad (12)$$

$$\hat{R}_p^-(\mathcal{F}) = \frac{1}{|\chi_p|} \sum_{d_i, y_c \in \chi_p} BCE(0, \mathcal{F}(d_i, y_c)),$$

where $\chi_p$ and $\chi_u$ represent the positive and unlabeled samples, respectively. Besides, the nnPu loss is deployed following the implementation of https://github.com/kiryor/nnPUlearning.

***GROOV [22]*** We use the codes released by the original paper https://github.com/facebookresearch/GROOV. Particularly, for a fair comparison, we have fine-tuned the generative model on our training set.

***ChatGPT [19]*** We have tried the prompts provided by Gilardi et al. [6], Huang et al. [9], Møller et al. [18]. To better instruct ChatGPT to perform the task, all prompts are translated into Chinese with minor changes. The results in Table 2 are obtained from the best-performing prompt, which is "请从歌曲评论{}中提取出多样的歌曲标签, 例如{}.format(user comments, gold labels)" ("Please select diverse labels from the user comments of song {}, for example , {} .format(user comments, gold labels)").

***Noisy Student Training (NST) [27]*** The codes released by the original paper https://github.com/google-research/noisystudent are based on image classification. In this paper, we replace the image classifier with MLC and EQ. (2).

### A.2 Metrics

Besides the metrics (P, SR, SF1, and Coverage scores) defined Sec. 4.1. We use the scikit-learn tool (https://scikit-learn.org/stable/modules/classes.html#module-sklearn.metrics) to compute the precision, recall, F1, and nDCG scores. Moreover, we use the implementation of https://github.com/kunaldahiya/pyxclib to compute the propensity scored precision (PSP) and propensity scored nDCG (PSnDCG) scores. Specifically, the propensity score of a label is estimated via label priors $\pi_{y^*} := \mathbb{P}[y^* = 1]$:

$$p_{y^*} := \frac{1}{1 + (\log N - 1) \times (b + 1)^a \times e^{-a \log(N \pi_{y^*} + b)}}, \qquad (13)$$

where $a = 0.55$ and $b = 1.5$, as recommended by Jain et al. [10].

## B CHINESE VERSION OF EXAMPLES IN SEC. 4.5

Tabel 4 presents the original Chinese version of the examples used Sec. 4.5.

**Table 4: Chinese version of the examples in Section 4.5.**

| Title: *The Last Of The Mohicans* | Artist: Alexandro Querevalú |
| --- | --- |
| **Gold labels** | 灵魂，唯美，烦恼，激情，伤感 |
| **Complete set of labels** | 悠扬，骄傲，无助，人生，时代，灵魂，男性，澎湃，美洲，呐喊，未来，人间，感人，天籁，动荡，激情，伤感，民族，孤独，感动，故事，唯美，烦恼，奴隶，战争 |
| **LightXML [11]** | 唯美，激情，一个人，感动，孤独，感人，人间 |
| **EQ. (2)** | 激情，伤感，一个人，感动，故事，孤独，感人，人间，灵魂，唯美 |
| **nnPU [12] w/ EQ. 2** | 激情，战争，骄傲，伤感，一个人，人生，故事，哭泣，感人，人间，未来，视频，时代，感动，见证，和平，食，纯，希望，天籁，害怕，利益，大街，苦难，记住，可怜，东方 |
| **GROOV [22]** | 治愈，宁静，回忆，经典，情歌，爱情，深情，一个人，故事，星河，青春，阳光，励志，热血，唯美，轻音乐，热爱，bgm，灵魂，古韵 |
| **ChatGPT [19]** | 原生态音乐，民族音乐，人文关怀，历史意义，反思人类文明，文化多样性，悲壮情感，抗争与挫折 |
| **NST [27] w/ EQ. (2)** | 人间，故事，一个人，伤感，激情，感人，感动，孤独，唯美 |
| **DiVa** | 感人，伤感，一个人，感动，故事，孤独，激情，唯美，动荡，人间，男性，奴隶，西班牙，吹奏，天籁，民族，时代，荡气回肠，震撼，诚信，伤感，战争，民主，灵魂，人生，澎湃，悠扬 |
| **Title:** *Mermaid* | **Artist: Skott** |
| **Gold labels** | 相见 |
| **Complete set of labels** | 情深，空旷，宁静，清冷，感染力，打雷，剪辑，森林，蛊惑，空灵，相见，美好，开心，深情，海盗，故事，惊艳，爱情 |
| **LightXML [11]** | 深情，情深，美好 |
| **EQ. (2)** | 爱情，情深，深情，故事，美好 |
| **nnPU [12] w/ EQ. 2** | 深情，情深，故事，视频，美好，惊艳，分享，宁静，森林，空灵，希望，开心，声线，风格，值得，期待，音色，相见，干净，宝藏，渴望，收藏，人心 |
| **GROOV [22]** | 治愈，星河，回忆，宁静，情歌，深情，爱情，伤感，经典，青春，阳光，相见，一个人，bgm，节奏，激情，灵魂 |
| **ChatGPT [19]** | 知名歌手打雷姐，美人鱼，海妖，Skott，Bjork，华为，冰雪2剪辑视频 |
| **NST [27] w/ EQ. (2)** | 故事，爱情，情深，深情，空灵，森林，美好，开心，宁静 |
| **DiVa** | 美好，爱情，故事，深情，情深，海盗，空灵，森林，蛊惑，空旷，清冷，路上，宁静，感染力，开心，惊艳，相见 |