

VLMine: Long-Tail Data Mining with Vision Language Models

Mao Ye¹, Gregory P. Meyer¹, Zaiwei Zhang¹, Dennis Park¹, Siva Karthik Mustikovela¹, Yuning Chai², Eric M Wolff¹

¹Cruise LLC.

²Meta Inc.

Abstract

Ensuring robust performance on long-tail examples is an important problem for many real-world applications of machine learning, such as autonomous driving. This work focuses on the problem of identifying rare examples within a corpus of unlabeled data. We propose a simple and scalable data mining approach that leverages the knowledge contained within a large vision language model (VLM). Our approach utilizes a VLM to summarize the content of an image into a set of keywords, and we identify rare examples based on keyword frequency. We find that the VLM offers a distinct signal for identifying long-tail examples when compared to conventional methods based on model uncertainty. Therefore, we propose a simple and general approach for integrating signals from multiple mining algorithms. We evaluate the proposed method on two diverse tasks: 2D image classification, in which inter-class variation is the primary source of data diversity, and on 3D object detection, where intra-class variation is the main concern. Furthermore, through the detection task, we demonstrate that the knowledge extracted from 2D images is transferable to the 3D domain. Our experiments consistently show large improvements (between 10% and 50%) over the baseline techniques on several representative benchmarks: ImageNet-LT, Places-LT, and the Waymo Open Dataset.

1 Introduction

In many real-world robotic applications, such as autonomous driving, the model must be able to handle a wide variety of situations. Ensuring the model’s performance in long-tail scenarios remains a challenging problem. Previous efforts on long-tail learning include improving model optimization (Cui et al. 2019a, 2021), model modification (Menon et al. 2020; Wang et al. 2020), and fine-tuning from pre-trained models (Long et al. 2022; Shi et al. 2023). While these works significantly improve the model’s performance on long-tail data, they mainly focus on how to improve the model given a fixed dataset. Alternatively, this paper focuses on mining for long-tail examples that can be added a training set to improve model performance. In many real-world applications, there exists a relatively small amount of labeled data and a vast amount of unlabeled data. Take autonomous driving for example, the vehicles collect enormous quantities of raw sensor data while driving. Most of this data is uninteresting; thus, annotating it would be costly with-

out providing much benefit. However, a tiny portion of the data consists of interest long-tail situations unavailable in our existing training set. As a result, automatic mining for this long-tail data is of the utmost importance.

Identifying long-tail examples from a large pool of data can be challenging. Existing approaches primarily utilize model uncertainty as a key signal (Gal, Islam, and Ghahramani 2017; Jiang et al. 2022; Choi et al. 2021). The underlying assumption is that the model tends to be less confident in its predictions for long-tail examples. However, an example that results in a high model uncertainty might not be a long-tail example but a hard example. For instance, in the study by (Jiang et al. 2022) on LiDAR-based outdoor 3D detection, hard examples generally consist of long-range or occluded vehicles. These hard examples are common, yet the model produces highly uncertain predictions for them, which significantly reduces the effectiveness of using uncertainty as a criterion for selecting long-tail examples.

Recently, foundation models such as large language models (LLMs) (Chowdhery et al. 2023; Touvron et al. 2023; OpenAI 2023; Chiang et al. 2023) and large vision language models (VLMs) (Alayrac et al. 2022; Dai et al. 2023; Awadalla et al. 2023; Liu et al. 2024; Zhu et al. 2024) have emerged. These models are trained with internet-scale data and have demonstrated impressive few-shot and zero-shot performance on recognition tasks (Radford et al. 2021; Liu et al. 2024). Leveraging the knowledge of VLM also benefits other tasks such as image clustering (Kwon et al. 2024) and out-of-domain detection (Jiang et al. 2024).

We argue that VLMs offer a more comprehensive understanding of long-tail examples due to their rich semantic extraction capabilities. However, the challenge lies in leveraging VLMs to enhance the performance of task-specific models on these examples. Directly deploying VLMs is often impractical due to latency concerns and the complexity of adapting them to task with different modalities. This paper demonstrates that the knowledge embedded in VLMs can serve as an effective signal for mining long-tail examples, thereby improving the quality of training data and enhancing the performance of task-specific models. The intuition is straightforward: since VLMs have been exposed to a highly diverse set of examples, they are capable of providing a more comprehensive understanding of semantic information. By comparing the semantic information described by the VLM,

we can identify examples with less frequent semantic descriptions, which serves as a stronger signal for detecting long-tail instances. To measure the rareness of an example based on these descriptions, we employ a simple keyword-based approach that summarizes the description into a set of keywords. The rareness of an example can then be approximated by the frequency of these keywords. We name our approach VLMine, and to the best of our knowledge, this is the first work to utilize a VLM for data mining.

VLMine does not use any information from the task-specific model. For that reason, the long-tail signal provided by our approach tends to be complementary to the signals obtained from traditional model-based mining algorithms (Beluch et al. 2018; Choi et al. 2021; Jiang et al. 2022). To leverage both types of mining techniques, we propose an algorithm, referred to as Pareto mining, to integrate long-tail signals from any number of sources.

The contributions of this work are summarized below.

- We propose VLMine, a simple model-agnostic long-tail data mining algorithm that leverages the world knowledge from a VLM.
- Additionally, we propose an algorithm for integrating the long-tail signal from both model-based and model-agnostic mining techniques.
- We evaluate the data mining performance with end-to-end metrics on 2D image classification and 3D object detection benchmarks, which demonstrates the ability of our proposed method to identify both inter- and intra-class long-tail variations on multiple domains. We show consistent improvement over baselines methods on ImageNet-LT, Places-LT, and the Waymo Open Dataset.

2 Related Work

Long-tail Perception Most existing work on long-tail visual recognition considers objects with less frequent class labels as long-tail examples and develops approaches to improve prediction accuracy in such label-imbalance situations. Applications include 2D tasks such as image classification (Cui et al. 2019b; Kang et al. 2019; Liu et al. 2019), segmentation (Hsieh et al. 2021; Wang et al. 2021; Wu et al. 2020), and object detection (Li et al. 2020; Tan et al. 2020, 2021a; Hyun Cho and Krähenbühl 2022). Typical approaches for long-tail perception include data reweighting and resampling (Cui et al. 2019a; Cao et al. 2019; Khan et al. 2019; Huang et al. 2019; Zhang et al. 2021; Hyun Cho and Krähenbühl 2022; Chawla et al. 2002; Han, Wang, and Mao 2005), gradient balancing (Tang, Huang, and Zhang 2020; Tan et al. 2021b), model or logits modifications (Wang et al. 2020; Alshammari et al. 2022; Menon et al. 2020; Ren et al. 2020), knowledge transfer or distillation (Wang, Ramanan, and Hebert 2017; Chu et al. 2020; Kim, Jeong, and Shin 2020; Liu et al. 2021; Xiang, Ding, and Han 2020; Li, Wang, and Wu 2021), representation learning (Zhong et al. 2021; Cui et al. 2021; Du et al. 2024), and fine-tuning pre-trained models (Ma et al. 2021; Tian et al. 2022; Dong et al. 2022; Long et al. 2022; Shi et al. 2023).

Compared with 2D vision, long-tail 3D object detection is less explored. (Peri et al. 2023) developed feature shar-

ing and camera-LiDAR fusion mechanisms to improve 3D detection on objects with infrequent class labels. (Ma et al. 2023) proposed a late-fusion approach of camera and LiDAR features for 3D detection on uncommon classes. Both (Peri et al. 2023; Ma et al. 2023) focus on inter-class long-tail examples. (Sick, Walter, and Abhau 2023) used shape-aware anchor distributions and heatmaps combined with camera-LiDAR fusion to improve 3D detection on objects with less frequent shapes. Those approaches aim at improving the detection on long-tail examples by modifying the model architecture rather than mining long-tail data. (Jiang et al. 2022) was the first to study intra-class long-tail 3D detection and proposed a density-based (using a flow model) and uncertainty-based approach for mining long-tail examples. Compared with (Jiang et al. 2022), we consider a new data mining approach by leveraging the knowledge from a pre-trained VLM.

Active Learning This work is closely related to pool-based active learning, where the goal is to actively select samples from a large set of unlabeled data to be labeled. Current active learning approaches can be classified into two categories: uncertainty-based and density-based. Uncertainty-based approaches use a model’s level of uncertainty as a signal to identify hard examples. The methods to measure uncertainty vary, with some representative approaches including ensemble variance (Beluch et al. 2018; Choi et al. 2021), entropy of the predicted distribution (Holub, Perona, and Burl 2008; Segal et al. 2022), and Bayesian modeling (Gal, Islam, and Ghahramani 2017; Harakeh, Smart, and Waslander 2020) of the logits. The other category selects examples such that the data distribution formed by those examples closely approximates the population data distribution (Sinha, Ebrahimi, and Darrell 2019; Sener and Savarese 2017; Gudovskiy et al. 2020). These approaches are not designed specifically for long-tail data mining. Despite some close connections between long-tail data mining and active learning, the goals are different. Active learning aims to reduce labeling costs by properly selecting a smaller amount of data, while data mining aims to discover long-tail data. To the best of our knowledge, we are the first to utilize a VLM to select data for labeling.

3 Method

Our goal is to find long-tail examples from a large pool of unlabeled data. Depending on the problem and data, the long-tail examples can be either inter-class, i.e. less common class categories, or intra-class, i.e. less common variations within a class. The main idea of our approach is to leverage the world knowledge from a VLM and LLM to extract keywords that describe the examples, and the frequency of the keywords is used to select the data to be labeled. Furthermore, our proposed Pareto mining can be used to incorporate other long-tail signal. An overview of our proposed method is illustrated in Figure 1.

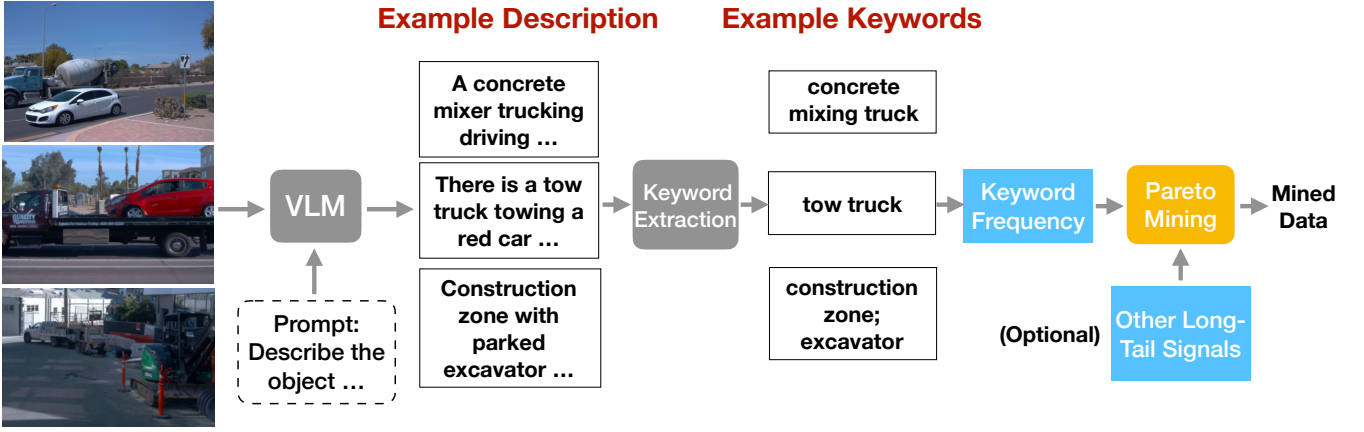


Figure 1: An overview of our proposed method. First, we prompt a VLM to describe the image, then the descriptions are summarized into a set of representative keywords using a rule-based heuristic or LLM. The frequency of the keywords are used to score the novelty of the images. Afterwards, the score can be combined with other long-tail signals, and our proposed Pareto mining is used to select the long-tail data to be labeled.

3.1 VLMine: VLM Knowledge for Long-tail Example Mining

VLMine is a straightforward mining procedure consisting of two steps: extraction of representative keywords and example selection via keyword frequency.

Extraction of representative keywords. We summarize the examples with a set of keywords. Given an image, we first prompt the VLM to describe the image, and then we summarize the description into a set of keywords. This can be done using a rule-based approach or by querying a LLM. After extracting the keywords, we also apply some standard post-processing, such as lemmatization and stop word removal. Note that each example might have a different number of keywords.

Example selection via keyword frequency. We compute the frequency of the keywords within all of the examples. The novelty score of each example is defined as a reverse monotonic transformation of the pooled keyword frequency. Depending on the data and the task, the pooling operator can be average pooling (i.e. averaging the frequency of all the keywords) or min-pooling (i.e. choosing the frequency of the least common keyword). Concretely, suppose an example has n keywords with frequency k_1, k_2, \dots, k_n , depends on the choice of pooling operator, its novelty score s is defined as

$$s := \phi(\min_i k_i) \quad \text{or} \quad s := \phi(\sum_i k_i / n), \quad (1)$$

where ϕ is any reverse monotonic transformation, e.g. $\phi(x) = -x$. Examples with higher novelty scores are assumed to be more rare. The choice of pooling function is task dependent. For example, in object detection, each example can contain multiple objects with most being commonly occurring objects. Thus, a majority of the keywords will be associated to common objects. In this case, the min-pooling operator is a better choice. However, for image classification, we expect most of the keywords to be associated

with the class and thus average pooling becomes a more reasonable choice.

We could directly ask the VLM to identify examples that are less common. However, such an approach selects examples based on the VLM’s knowledge, which introduces bias. The examples that the VLM believes to be uncommon might not actually be rare in the task database.

3.2 Pareto Mining: Combining Multiple Data Mining Signals

The novelty score determined by VLMine is model-agnostic. It purely uses the knowledge from a VLM and possibly a LLM without utilizing any information from the task-specific model we want to improve, e.g. an image classifier or object detector. We find that the signal provided by VLMine is often complementary to the signals obtained from traditional model-based approaches such as uncertainty-based mining. On the other hand, uncertainty-based mining algorithms, which all rely on the task-model’s knowledge, are often correlated with each other. Intuitively, if multiple signals are highly correlated, a long-tail example that is missed by one technique is likely to be missed by another. However, if the signals are complementary, even if an example is missed by one algorithm, we have a higher chance to detect it with another. Therefore, it is beneficial to develop a mining strategy that can integrated multiple signals.

An obvious approach is to combine the signals from multiple algorithms is to use a linear combinations of their scores. However, the final result will heavily rely on the weights of the linear combination, which may be hard to determine since the magnitude of each signal might be very different and there is not an explicit objective to optimize. Refer to Appendix A.3 for an analysis.

We propose a simple and general approach to combine the signals from multiple algorithms. The idea is motivated by the concept of a Pareto frontier in multi-objective opti-

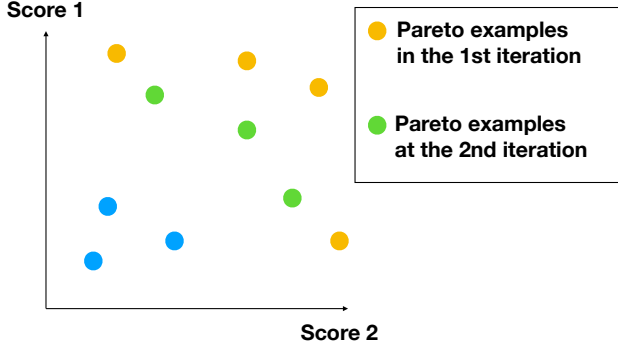


Figure 2: Illustration of the Pareto examples identified by Pareto mining.

mization. For each example, suppose we receive n scores $\mathbf{s} = [s_1, s_2, \dots, s_n]$ from different algorithms that represent the novelty of the example. Without loss of generality, we assume that a signal with a higher value means a higher chance that the example is part of the long-tail. We say that an example with scores \mathbf{s} is dominated by an example with scores \mathbf{s}' if $\forall i \in 1, \dots, n, s_i \leq s'_i$ and there exists $j \in 1, \dots, n$ such that $s_j < s'_j$.

We say an example with scores \mathbf{s} is a Pareto example if it is not dominated by any other examples. Figure 2 illustrates the case when two novelty scores are used. We add samples from the Pareto examples into the mined dataset. Once all of the Pareto examples are added, we remove them from the frontier and determine the new Pareto examples until we obtain the desired amount of mined examples. In Figure 2, the examples colored yellow are the Pareto examples at the start of mining, and the examples colored green and blue are dominated. After the yellow Pareto examples are removed, the examples colored green become the new set of Pareto examples. See Algorithm 1 for the Pareto mining pseudo-code.

4 Experiments

Our experiments are organized as follows. In Section 4.1 and 4.2, we evaluate VLMine and Pareto mining using a pool-based active learning setup. Our goal is to demonstrate that the proposed method is effective in identifying both inter- and intra-class long-tail examples. Details on the prompt engineering and implementations are also discussed in these sections. Additional visualizations, studies, and remarks are presented in the appendix. In Appendix A.1, we provide visualizations to illustrate that VLMine successfully mines interesting long-tail examples. In Appendix A.2, we visualize examples identified by combining model-agnostic and model-based signals through Pareto mining. In Appendix A.3, we conduct an ablation study to understand how the design of the prompt influences the final performance, and to compare Pareto mining against other methods for combining novelty scores. In Appendix A.4 and A.5, we discuss limitations related to the VLM and their impact on VLMine.

Algorithm 1: Pseudo-code for Pareto Mining

Input: Set of examples \mathcal{I} , number of examples to mine N
Output: Set of mined examples \mathcal{M}
 $\mathcal{M} \leftarrow \emptyset$
while $|\mathcal{M}| < N$ **do**
 Identify the subset of Pareto examples $\mathcal{P} \subset \mathcal{I}$
 if $|\mathcal{P}| \leq N - |\mathcal{M}|$ **then**
 $\mathcal{M} \leftarrow \mathcal{M} \cup \mathcal{P}$
 else
 Create $\hat{\mathcal{P}}$ by sampling $N - |\mathcal{M}|$ examples from \mathcal{P}
 $\mathcal{M} \leftarrow \mathcal{M} \cup \hat{\mathcal{P}}$
 end if
end while

4.1 Long-tail 2D Image Classification

Setup. We consider two long-tail datasets for the image classification task: ImageNet-LT and Places-LT. The datasets were created by sampling a subset of the original dataset following a Pareto distribution (Liu et al. 2019). ImageNet-LT is an object-centric dataset with a training split consists of 115.8K images from 1000 categories and cardinality ranging from 5 to 1280. Places-LT is a scene-centric dataset with a training set that consists of 5 to 4980 images for each of its 365 classes for a total of 62.5K images. For both datasets, the validation and test splits are balanced and contain 20 and 100 examples per class, respectively.

We conduct pool-based active learning experiments following (Settles 2009; Jiang et al. 2022). For each class, we sub-sample 20% of the training examples, keeping the distribution of classes the same. This subset forms the “labeled” pool, and the rest of the training set and the entire validation set are considered the “unlabeled” pool from which we will mine data. The test set remains unchanged and is used for evaluation.

We use predictive entropy (Shannon 1948) and variation ratios (Freeman 1965) as the baselines, which have been shown to be strong acquisition functions based on the analysis in (Gal, Islam, and Ghahramani 2017). Predictive entropy,

$$-\sum_c p(y = c|x) \log p(y = c|x), \quad (2)$$

calculates the entropy of the predicted class distribution, and higher entropy values imply a higher novelty for the example. Likewise, the variation ratio,

$$1 - \max_c p(y = c|x), \quad (3)$$

measures the lack of confidence in a prediction, where a larger ratio is assumed to indicate a higher chance of an example being part of the long-tail.

Prompt and Implementation Details. For these experiments, we use LLaVA1.5-7B (Liu et al. 2024) as the VLM and a heuristic to generate keywords. For ImageNet-LT, we prompt the VLM with the following: What are the possible classes for this image? Give three possible answers. For Places-LT, the

following prompt is used: What are the possible scene categories for this image? Give three possible answers.

We prompt the VLM to provide three possible answers for each example, since we find that this results in the VLM supplying more comprehensive descriptions. For an empirical analysis of prompt, please refer to Appendix A.3. For these experiments, we simply treat each word in the description as a keyword. Afterwards, we apply standard language post-processing, i.e. stop word removal and lemmatization, on the keywords. The novelty score of each image is computed using the average frequency of the keywords.

For ImageNet-LT, we use a ResNet-50 (He et al. 2016) backbone and the optimization protocol from (Du et al. 2024). (Du et al. 2024) uses a contrastive loss as an auxiliary loss and achieves the state-of-the-art in long-tail learning. Furthermore, we found the logit-adjustment loss (Menon et al. 2020) to be more stable than the softmax loss when training on the sub-sampled dataset. The model is trained for 90 epochs with a batch size of 256, a cosine-decayed learning rate of 0.1, and SGD optimizer with 0.9 momentum. For Places-LT, we follow the protocol from (Liu et al. 2019) and use an ImageNet pre-trained ResNet-152 (He et al. 2016) as the backbone. We fine-tune the model for 30 epochs with a batch size of 256, a cosine-decayed learning rate of 0.02, a SGD optimizer with 0.9 momentum, and standard data augmentations, i.e. randomized image crops and flips. We also apply the logit-adjustment method (Menon et al. 2020) to focal loss to improve learning.

Results. Figure 3 shows the classification accuracy on the test split for ImageNet-LT and Places-LT when different numbers of examples are mined. The number of examples mined from the “unlabeled” pool varies from 10% to 50%. The left-hand side of Figure 3 plots the top-1 accuracy for all classes while the right-hand side plots the accuracy for only the tail classes. The model trained with data obtained from VLMine consistently outperforms the baselines. Note that classifying the tail classes is quite challenging as the proportion of those classes is small compared with head classes.

To provide more insight, we evaluate the mining algorithms by inspecting the long-tailness of the mined examples. We define the rareness of each mined example by the number of images corresponding to its class in the original “labeled” pool. Thus, a technique is better at identifying long-tail examples if the mining frequency of rarer examples is higher. Figure 4 shows the distribution of the 2000 examples mined by VLMine and the predictive entropy uncertainty-based approach. As we can see, VLMine mines significantly more of the rarist classes than the uncertainty-based method. Demonstrating that a VLM is more capability at identifying long-tail examples in this case.

Furthermore, we evaluate our proposed Pareto mining with and without the novelty score from VLMine and the results are shown in Figure 3. For ImageNet-LT, Pareto mining with VLMine gives a boost in performance, especially for tail classes, when more examples are mined. In comparison, without VLMine and only using the model-based al-

gorithms, Pareto mining does not give as significant of an improvement. In Figure 5, we plot the signals from each algorithm for each example in ImageNet-LT. We observe that the signals from the two uncertainty-based approaches are highly correlated while the signal from VLMine is orthogonal to the model-based approaches. This illustrates why combining model-agnostic and model-based signals improves data mining.

4.2 Long-tail 3D Object Detection

We evaluate our proposed method on the real-world task of 3D object detection for autonomous driving. For this challenging setting, VLMine needs to identify long-tail objects within LiDAR point clouds using the corresponding camera images.

Setup. For our experiments, we use the large-scale Waymo Open Dataset (Sun et al. 2020). The dataset consists of 1150 sequences with each sequence containing 20 seconds of driving data captured at 10 Hz. Each frame consists of a LiDAR point cloud and 5 camera images. We follow the active learning setup from (Jiang et al. 2022), which achieves the state-of-the-art for this task. We randomly split the sequences within the training set into a 10% “labeled” pool and a 90% “unlabeled” pool. The large “unlabeled” pool can be viewed as the stream of data coming from the autonomous fleet, from which we want to mine long-tail examples. For these experiments, the goal is to identify a small portion of the “unlabeled” pool (only 3% following (Jiang et al. 2022)) that adds the most value when added to the training set.

Existing self-driving benchmarks lack fine-grained inter- or intra-class labels. For that reason, (Jiang et al. 2022) uses the shape of an object as a proxy for intra-class long-tailness. Their assumption is that a vehicle with a larger size and a pedestrian with a smaller height are part of the long-tail. Hence, they evaluate the mining techniques by decomposing the Average Precision (AP) metrics for the vehicle and pedestrian classes. Specifically, they report detection performance on larger vehicles (longer than 8 and 10 meters) and smaller pedestrians (shorter than 1.4 and 1.2 meters). Although these length/height-based decompositions are far from ideal for evaluating long-tail performance, according to the analysis in (Jiang et al. 2022), these road users are uncommon in the Waymo Open Dataset. Our baseline for comparisons is D-REM (Jiang et al. 2022), which is an uncertainty-based approach using an ensemble of models to infer the novelty score.

Prompt and Implementation Details. For these experiments, we again use LLaVA v1.5-7B (Liu et al. 2024) as the VLM, but generate keywords using a LLM, GPT3.5-turbo (OpenAI 2023). We prompt the VLM to describe the camera images with the following: Describe the uncommon or abnormal vehicles, pedestrians, and cyclists related to traffic in this image. Afterwards, we prompt the LLM to summarize the description into a set of keywords using: Please return keywords for each image description in this format: keyword1, keyword2, and

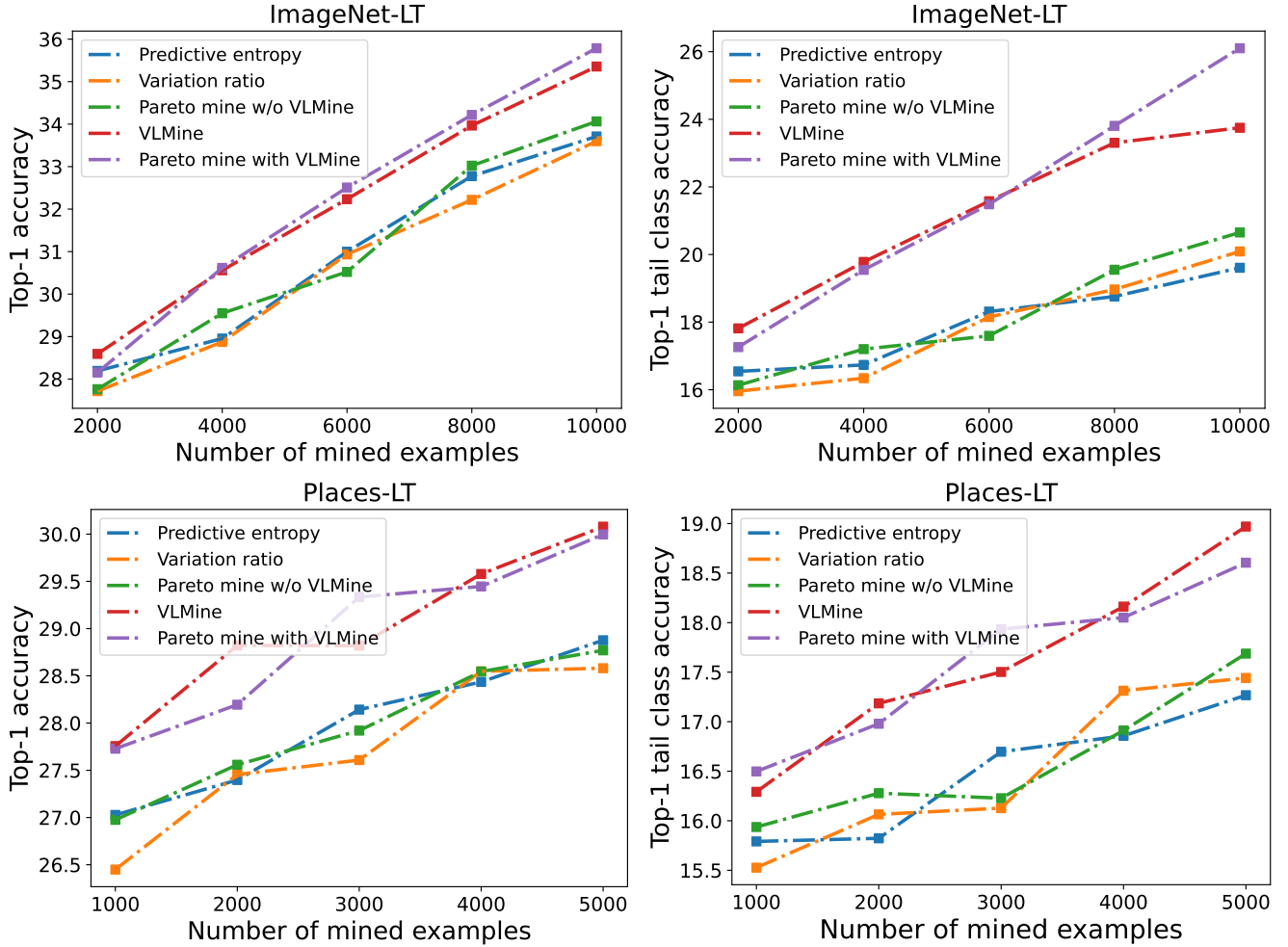


Figure 3: Data mining experiments on ImageNet-LT and Places-LT.

etc. The main reason we use a LLM for keyword extraction in this case is that the situations are more complicated, and we need to ensure the keywords well represent the intra-class variations. For example, in the following description, “There is a construction truck on the street,” we want the algorithm to output keywords like “construction truck” rather than “construction” and “truck”. For that reason, the heuristic approach from the previous experiment is insufficient.

We use a two-step approach to generate keywords by first asking the VLM to provide a description and then extracting keywords based on that description. Alternatively, we can directly ask the VLM to describe the example using only keywords. However, we identify some issues with that approach and refer readers to Appendix A.4 for more details.

The evaluation only considers long-tail vehicles and pedestrians, but the description from the VLM and the keywords from the LLM might contain information irrelevant to those types of objects. To address this, we filter the keywords by asking the LLM what type of road user a keyword is related to using the prompt: For each

of the following words separated by ; , identify whether they are related to a description of vehicle, pedestrian, or cyclist. Example answers are: child - pedestrian; bus - vehicle; dish - not related.

Since each frame contains several camera images, we simply concatenate the keywords from all the images. In Section 4.1, the images are either object-centric or scene-centric, so all the keywords are relevant to decide the rareness of the example. However, for 3D object detection, each frame may contain many objects, and the keywords describing common objects may not be relevant to the potential long-tail objects in the example. Therefore, we use the minimum frequency to compute the novelty score instead of the average frequency. This approach allows us to focus on the presence of less common keywords to determine the long-tailness of the example.

For the experiments, our detector is based on CenterPoint (Yin, Zhou, and Krahenbuhl 2021). For D-REM, we train 5 models with different random initializations while keeping

Method	Vehicle (AP/APH)			Pedestrian (AP/APH)		
	All	$\geq 8m$	$\geq 10m$	All	$\leq 1.4m$	$\leq 1.2m$
D-REM	62.5/61.9	27.1/25.9	16.7/14.8	66.4/50.8	15.2/9.8	10.4/5.9
VLMine	63.0/62.3	38.0/37.2	25.2/22.0	65.2/49.6	18.7/13.4	9.6/5.5
Pareto Mining (VLMine + D-REM)	62.9/62.1	39.3/37.6	25.1/24.3	65.4/50.3	21.6/15.9	11.6/7.2

Table 1: Data mining experiment on the Waymo Open Dataset.

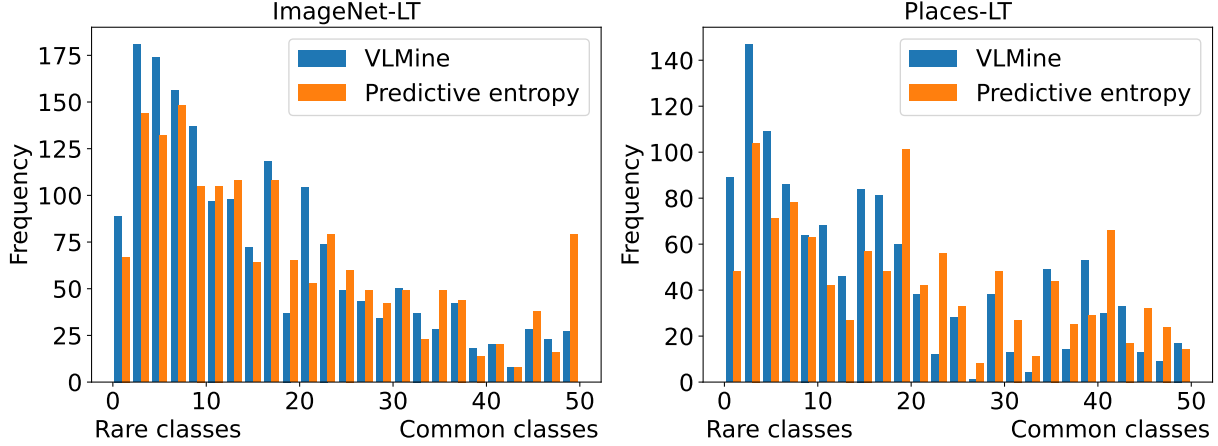


Figure 4: Distribution of the mined data sorted by rareness (the bars further to the left represent rare classes while the bars on the right correspond to more common classes). For readability, we only show classes that have less than 50 images in the original “labeled” pool. The rareness of each class is quantified by the number of images for the class in the original “labeled” pool. We plot the frequency of mined examples for different rarenesses.

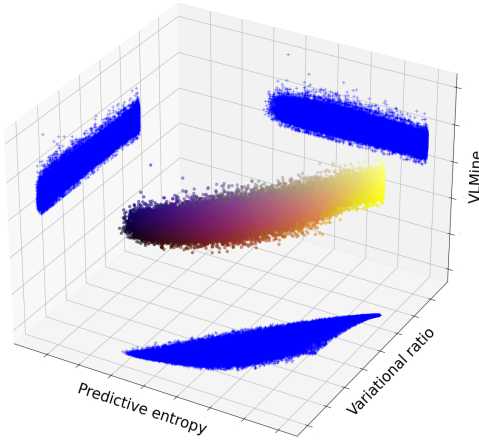


Figure 5: Correlation of the novelty scores from different algorithms on ImageNet-LT. We plot the scores from three different algorithms for each example in the unlabeled pool and project the scores to show the correlation between each pair of algorithms. As we can see, the scores between predictive entropy and variational ratio are highly correlated, while VLMine provides orthogonal signals.

Figure 6: Visualizations of Pareto mining.

all other training configurations the same. We also evaluate the integration of the two approaches using the proposed Pareto mining.

Results. The results of these experiments are summarized in Table 1. When considering all objects, unsurprisingly all of the methods have similar performance. Our proposed VLMine significantly outperforms D-REM on all tail objects except for pedestrians under 1.2 meters where the two approaches perform comparably. We observe a further boost in performance on tail objects when employing the proposed Pareto mining technique, which demonstrates the value of combining model-agnostic and model-based novelty scores. We refer readers to Appendix A.1 for a visualization of the mined examples, and to Appendix A.2 for an illustration of Pareto examples.

5 Conclusion

This work proposes VLMine, a simple yet effective long-tail data mining algorithm that leverages the knowledge from a vision language model. VLMine provides a strong signal to detect long-tail examples which is complementary to the signals given by traditional uncertainty-based approaches. Furthermore, we show how our proposed Pareto mining can integrate multiple mining signals to further enhance performance.

References

- Alayrac, J.-B.; Donahue, J.; Luc, P.; Miech, A.; Barr, I.; Hasson, Y.; Lenc, K.; Mensch, A.; Millican, K.; Reynolds, M.; et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35: 23716–23736.
- Alshammari, S.; Wang, Y.-X.; Ramanan, D.; and Kong, S. 2022. Long-tailed recognition via weight balancing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6897–6907.
- Awadalla, A.; Gao, I.; Gardner, J.; Hessel, J.; Hanafy, Y.; Zhu, W.; Marathe, K.; Bitton, Y.; Gadre, S.; Sagawa, S.; et al. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*.
- Beluch, W. H.; Genewein, T.; Nürnberger, A.; and Köhler, J. M. 2018. The power of ensembles for active learning in image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 9368–9377.
- Cao, K.; Wei, C.; Gaidon, A.; Arechiga, N.; and Ma, T. 2019. Learning imbalanced datasets with label-distribution-aware margin loss. *Advances in neural information processing systems*, 32.
- Chawla, N. V.; Bowyer, K. W.; Hall, L. O.; and Kegelmeyer, W. P. 2002. SMOTE: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16: 321–357.
- Chiang, W.-L.; Li, Z.; Lin, Z.; Sheng, Y.; Wu, Z.; Zhang, H.; Zheng, L.; Zhuang, S.; Zhuang, Y.; Gonzalez, J. E.; et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023).
- Choi, J.; Elezi, I.; Lee, H.-J.; Farabet, C.; and Alvarez, J. M. 2021. Active learning for deep object detection via probabilistic modeling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 10264–10273.
- Chowdhery, A.; Narang, S.; Devlin, J.; Bosma, M.; Mishra, G.; Roberts, A.; Barham, P.; Chung, H. W.; Sutton, C.; Gehrmann, S.; et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113.
- Chu, P.; Bian, X.; Liu, S.; and Ling, H. 2020. Feature space augmentation for long-tailed data. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX 16*, 694–710. Springer.
- Cui, J.; Zhong, Z.; Liu, S.; Yu, B.; and Jia, J. 2021. Parametric contrastive learning. In *Proceedings of the IEEE/CVF international conference on computer vision*, 715–724.
- Cui, Y.; Jia, M.; Lin, T.-Y.; Song, Y.; and Belongie, S. 2019a. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 9268–9277.
- Cui, Y.; Jia, M.; Lin, T.-Y.; Song, Y.; and Belongie, S. 2019b. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 9268–9277.
- Dai, W.; Li, J.; Li, D.; Tiong, A.; Zhao, J.; Wang, W.; Li, B.; Fung, P.; and Hoi, S. 2023. InstructBLIP: Towards General-purpose Vision-Language Models with Instruction Tuning. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Dong, B.; Zhou, P.; Yan, S.; and Zuo, W. 2022. LPT: long-tailed prompt tuning for image classification. In *The Eleventh International Conference on Learning Representations*.
- Du, C.; Wang, Y.; Song, S.; and Huang, G. 2024. Probabilistic Contrastive Learning for Long-Tailed Visual Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Freeman, L. C. 1965. Elementary applied statistics: for students in behavioral science. (No Title).
- Gal, Y.; Islam, R.; and Ghahramani, Z. 2017. Deep bayesian active learning with image data. In *International conference on machine learning*, 1183–1192. PMLR.
- Gudovskiy, D.; Hodgkinson, A.; Yamaguchi, T.; and Tsukizawa, S. 2020. Deep active learning for biased datasets via fisher kernel self-supervision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 9041–9049.
- Han, H.; Wang, W.-Y.; and Mao, B.-H. 2005. Borderline-SMOTE: a new over-sampling method in imbalanced data sets learning. In *International conference on intelligent computing*, 878–887. Springer.
- Harakeh, A.; Smart, M.; and Waslander, S. L. 2020. Bayesod: A bayesian approach for uncertainty estimation in deep object detectors. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, 87–93. IEEE.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- Holub, A.; Perona, P.; and Burl, M. C. 2008. Entropy-based active learning for object recognition. In *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, 1–8. IEEE.
- Hsieh, T.-I.; Robb, E.; Chen, H.-T.; and Huang, J.-B. 2021. Droploss for long-tail instance segmentation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, 1549–1557.
- Huang, C.; Li, Y.; Loy, C. C.; and Tang, X. 2019. Deep imbalanced learning for face recognition and attribute prediction. *IEEE transactions on pattern analysis and machine intelligence*, 42(11): 2781–2794.
- Hyun Cho, J.; and Krähenbühl, P. 2022. Long-tail detection with effective class-margins. In *European Conference on Computer Vision*, 698–714. Springer.
- Jiang, C. M.; Najibi, M.; Qi, C. R.; Zhou, Y.; and Anguelov, D. 2022. Improving the intra-class long-tail in 3d detection via rare example mining. In *European Conference on Computer Vision*, 158–175. Springer.
- Jiang, X.; Liu, F.; Fang, Z.; Chen, H.; Liu, T.; Zheng, F.; and Han, B. 2024. Negative Label Guided OOD Detection

- with Pretrained Vision-Language Models. In *The Twelfth International Conference on Learning Representations*.
- Kang, B.; Xie, S.; Rohrbach, M.; Yan, Z.; Gordo, A.; Feng, J.; and Kalantidis, Y. 2019. Decoupling representation and classifier for long-tailed recognition. *arXiv preprint arXiv:1910.09217*.
- Khan, S.; Hayat, M.; Zamir, S. W.; Shen, J.; and Shao, L. 2019. Striking the right balance with uncertainty. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 103–112.
- Kim, J.; Jeong, J.; and Shin, J. 2020. M2m: Imbalanced classification via major-to-minor translation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 13896–13905.
- Kwon, S.; Park, J.; Kim, M.; Cho, J.; Ryu, E. K.; and Lee, K. 2024. Image Clustering Conditioned on Text Criteria. In *The Twelfth International Conference on Learning Representations*.
- Li, T.; Wang, L.; and Wu, G. 2021. Self supervision to distillation for long-tailed visual recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, 630–639.
- Li, Y.; Wang, T.; Kang, B.; Tang, S.; Wang, C.; Li, J.; and Feng, J. 2020. Overcoming classifier imbalance for long-tail object detection with balanced group softmax. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10991–11000.
- Liu, B.; Li, H.; Kang, H.; Hua, G.; and Vasconcelos, N. 2021. Gistnet: a geometric structure transfer network for long-tailed recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 8209–8218.
- Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2024. Visual instruction tuning. *Advances in neural information processing systems*, 36.
- Liu, Z.; Miao, Z.; Zhan, X.; Wang, J.; Gong, B.; and Yu, S. X. 2019. Large-scale long-tailed recognition in an open world. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2537–2546.
- Long, A.; Yin, W.; Ajanthan, T.; Nguyen, V.; Purkait, P.; Garg, R.; Blair, A.; Shen, C.; and van den Hengel, A. 2022. Retrieval augmented classification for long-tail visual recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 6959–6969.
- Ma, T.; Geng, S.; Wang, M.; Shao, J.; Lu, J.; Li, H.; Gao, P.; and Qiao, Y. 2021. A simple long-tailed recognition baseline via vision-language model. *arXiv preprint arXiv:2111.14745*.
- Ma, Y.; Peri, N.; Wei, S.; Hua, W.; Ramanan, D.; Li, Y.; and Kong, S. 2023. Long-Tailed 3D Detection via 2D Late Fusion. *arXiv preprint arXiv:2312.10986*.
- Menon, A. K.; Jayasumana, S.; Rawat, A. S.; Jain, H.; Veit, A.; and Kumar, S. 2020. Long-tail learning via logit adjustment. *arXiv preprint arXiv:2007.07314*.
- OpenAI, R. 2023. Gpt-4 technical report. arxiv 2303.08774. *View in Article*, 2: 13.
- Peri, N.; Dave, A.; Ramanan, D.; and Kong, S. 2023. Towards long-tailed 3d detection. In *Conference on Robot Learning*, 1904–1915. PMLR.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.
- Ren, J.; Yu, C.; Ma, X.; Zhao, H.; Yi, S.; et al. 2020. Balanced meta-softmax for long-tailed visual recognition. *Advances in neural information processing systems*, 33: 4175–4186.
- Segal, S.; Kumar, N.; Casas, S.; Zeng, W.; Ren, M.; Wang, J.; and Urtasun, R. 2022. Just label what you need: Fine-grained active selection for p&p through partially labeled scenes. In *Conference on Robot Learning*, 816–826. PMLR.
- Sener, O.; and Savarese, S. 2017. Active learning for convolutional neural networks: A core-set approach. *arXiv preprint arXiv:1708.00489*.
- Settles, B. 2009. Active learning literature survey.
- Shannon, C. E. 1948. A mathematical theory of communication. *The Bell system technical journal*, 27(3): 379–423.
- Shi, J.-X.; Wei, T.; Zhou, Z.; Han, X.-Y.; Shao, J.-J.; and Li, Y.-F. 2023. Parameter-Efficient Long-Tailed Recognition. *arXiv preprint arXiv:2309.10019*.
- Sick, B.; Walter, M.; and Abhau, J. 2023. Adaptive-Shape: Solving Shape Variability for 3D Object Detection with Geometry Aware Anchor Distributions. *arXiv preprint arXiv:2302.14522*.
- Sinha, S.; Ebrahimi, S.; and Darrell, T. 2019. Variational adversarial active learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 5972–5981.
- Sun, P.; Kretschmar, H.; Dotiwalla, X.; Chouard, A.; Patnaik, V.; Tsui, P.; Guo, J.; Zhou, Y.; Chai, Y.; Caine, B.; et al. 2020. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2446–2454.
- Tan, J.; Lu, X.; Zhang, G.; Yin, C.; and Li, Q. 2021a. Equalization loss v2: A new gradient balance approach for long-tailed object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 1685–1694.
- Tan, J.; Lu, X.; Zhang, G.; Yin, C.; and Li, Q. 2021b. Equalization loss v2: A new gradient balance approach for long-tailed object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 1685–1694.
- Tan, J.; Wang, C.; Li, B.; Li, Q.; Ouyang, W.; Yin, C.; and Yan, J. 2020. Equalization loss for long-tailed object recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 11662–11671.
- Tang, K.; Huang, J.; and Zhang, H. 2020. Long-tailed classification by keeping the good and removing the bad momentum causal effect. *Advances in Neural Information Processing Systems*, 33: 1513–1524.

Tian, C.; Wang, W.; Zhu, X.; Dai, J.; and Qiao, Y. 2022. Vl-ltr: Learning class-wise visual-linguistic representation for long-tailed visual recognition. In *European Conference on Computer Vision*, 73–91. Springer.

Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; Azhar, F.; et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Wang, J.; Zhang, W.; Zang, Y.; Cao, Y.; Pang, J.; Gong, T.; Chen, K.; Liu, Z.; Loy, C. C.; and Lin, D. 2021. Seesaw loss for long-tailed instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 9695–9704.

Wang, X.; Lian, L.; Miao, Z.; Liu, Z.; and Yu, S. X. 2020. Long-tailed recognition by routing diverse distribution-aware experts. *arXiv preprint arXiv:2010.01809*.

Wang, Y.-X.; Ramanan, D.; and Hebert, M. 2017. Learning to model the tail. *Advances in neural information processing systems*, 30.

Wu, J.; Song, L.; Wang, T.; Zhang, Q.; and Yuan, J. 2020. Forest r-cnn: Large-vocabulary long-tailed object detection and instance segmentation. In *Proceedings of the 28th ACM international conference on multimedia*, 1570–1578.

Xiang, L.; Ding, G.; and Han, J. 2020. Learning from multiple experts: Self-paced knowledge distillation for long-tailed classification. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, 247–263. Springer.

Yin, T.; Zhou, X.; and Krahenbuhl, P. 2021. Center-based 3d object detection and tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 11784–11793.

Zhang, S.; Li, Z.; Yan, S.; He, X.; and Sun, J. 2021. Distribution alignment: A unified framework for long-tail visual recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2361–2370.

Zhong, Z.; Cui, J.; Liu, S.; and Jia, J. 2021. Improving calibration for long-tailed recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 16489–16498.

Zhu, D.; Chen, J.; Shen, X.; Li, X.; and Elhoseiny, M. 2024. MiniGPT-4: Enhancing Vision-Language Understanding with Advanced Large Language Models. In *The Twelfth International Conference on Learning Representations*.

A Appendix

A.1 Analysis of Mined Examples

We visualize the data mined by VLMine in Section 4.2 for the Waymo Open Dataset. Figure 8 shows the images and keywords for mined long-tail vehicles. Figure 9 depicts some typical failure modes and/or less accurate keywords from VLMine. In Figure 9a-d, since the objects are hard to describe, VLMine fails to give accurate keywords. However, such failures may not degrade the quality of the mined data. The content of the images is often part of the long-tail, and although the keywords are not accurate, their frequencies are lower resulting in a higher novelty score. Figure 9e illustrates a case where VLMine detects an improperly parked van but misses the fire truck. Figure 9f shows a case where VLMine wrongly recognizes a yellow wall as a yellow truck.

A.2 Pareto Mining Visualization

Figure 10 illustrates how the signal from VLMine differs from traditional uncertainty-based approaches like D-REM. We depict four types of examples corresponding to long-tail vehicles. Three are Pareto examples where at least one of the mining methods gives the example a high score. The other is a dominated example where both techniques give the example a low score. The examples that have a high novelty score for both D-REM and VLMine are indeed long-tail, such as the trolley train and supply truck. D-REM will often score challenging yet common vehicles highly when VLMine gives them a small novelty score. We observe that large vehicles tend to be given a higher uncertainty by D-REM. As a result, typical examples with a low D-REM scores but a high VLMine scores are regular-sized but uncommon vehicles, such as a small-sized dump truck or a detached trailer. The examples where both scores are low are common low-complexity scenes with few objects.

A.3 Ablation Studies

Prompt Design We study how the design of the prompt influences the final results. For this study, we use ImageNet-LT. For our experiments in Section 4.1, we prompt the VLM to provide a description of the image with What are the possible classes for this image? Give three possible answers. We find such a prompt, which we refer to as Prompt 1, will force the model to produce more informative descriptions of the images and thus benefits the data mining. To illustrate this, we conduct an experiment with the following prompt: Describe the object in the image, which we refer to as Prompt 2. Here are the outputs of the VLM when the two prompts are used for the image shown in Figure 7.

- Prompt 1: “1. Sports and Recreation: The image features a man holding a large fish, which could be a trophy catch from a fishing trip. This suggests that the man is likely an avid fisherman or enjoys spending time outdoors engaging in recreational activities. 2. Nature and Wildlife: The image showcases a man holding a fish, possibly a walleye, in a natural setting. This could be a representation of the beauty of nature and the importance of preserving and protecting aquatic ecosystems. 3. Fishing and Angling:



Figure 7: Example of an image from ImageNet-LT.

The man is holding a fish, which is a common subject in fishing and angling-related images. This image could be used to promote fishing trips, fishing equipment, or simply to showcase the thrill of catching a large fish.”

- Prompt 2: “The object in the image is a large fish, specifically a muskie, which is being held by a man.”

Figure 11 shows the classification accuracy when the different prompts are used. We observe consistent improvement when using Prompt 1, and the improvements are significant for the tail classes.

Integrating Novelty Scores The proposed Pareto mining provides a simple and principled mechanism to integrate the long-tail signals from multiple algorithms. A linear combination of the novelty scores is another way to leverage multiple signals. In this section, we compare Pareto mining with linear combination. Since the scores from different algorithms might have different statistics, we first normalized all the scores to have a mean of zero and a standard deviation of one. We combine all novelty scores for a particular example by averaging the normalized scores. Figure 12 summarizes our results. Overall, we observe that Pareto mining consistently outperforms the linear combination.

A.4 Remark on One-step Keywords Generation

It is noteworthy that we use a two-step approach to generate keywords for each example. We start by obtaining a description of the example by prompting the VLM, and then in the second step, we extract the keywords from the description using a rule-based approach or a LLM. Alternatively, we can directly ask the VLM to describe the example using keywords. However, we find that the VLM might output non-informative keywords, miss important objects, or generate inconsistent text. For example, using the following prompt: List the uncommon or abnormal vehicles, pedestrians and cyclists related to traffic in this image. Return only the keywords in the following format: keyword1, keyword2, etc. gives the subsequent results. For Figure 8a, the output is repetitive and inconsistent: “Concrete truck, car, cement mixer, cement truck, white car, concrete truck, cement



(a) concrete mixer truck



(b) double-length subway train



(c) open-top sightseeing bus



(d) excavator



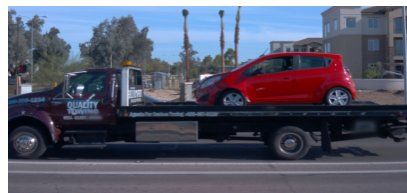
(e) construction vehicle



(f) tailgate open



(g) truck with open tailgate



(h) tow truck



(i) truck carrying advertisement

Figure 8: Example images from the Waymo Open Dataset and their keywords.

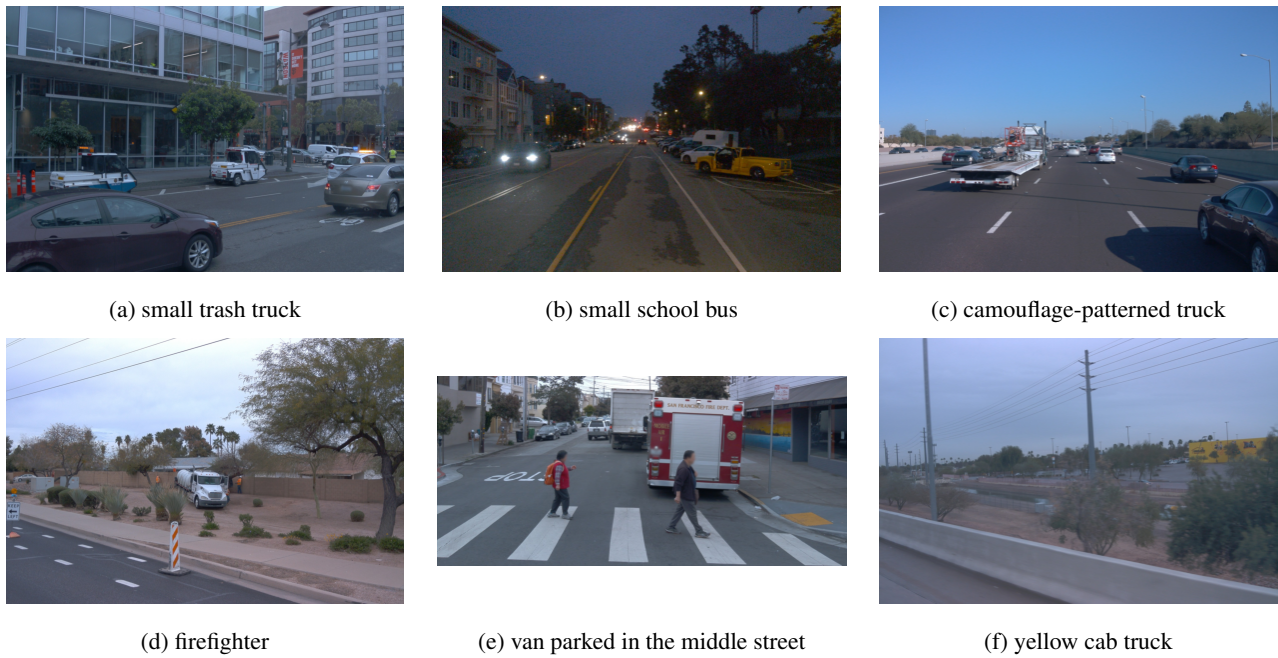


Figure 9: Examples of inaccurate keywords from the Waymo Open Dataset.

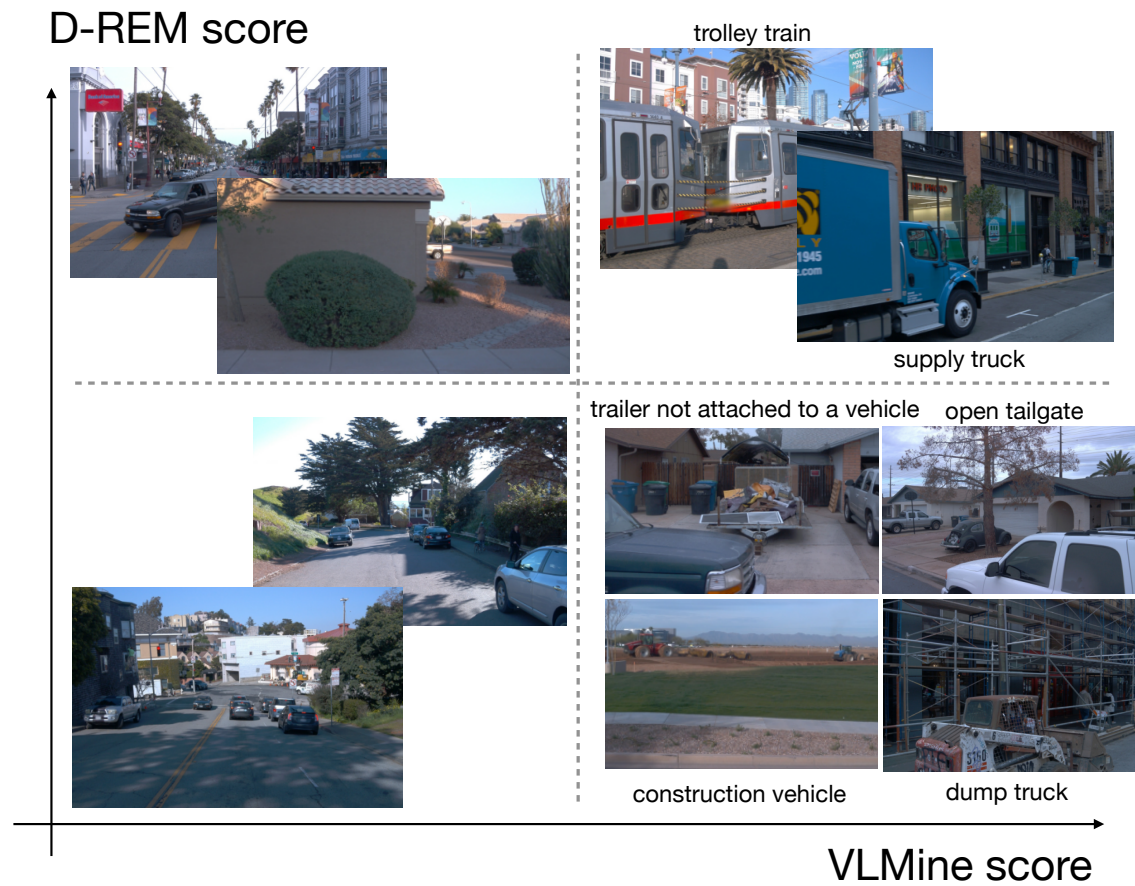


Figure 10: Visualization of Pareto and dominated examples from the Waymo Open Dataset. Images are cropped for readability.

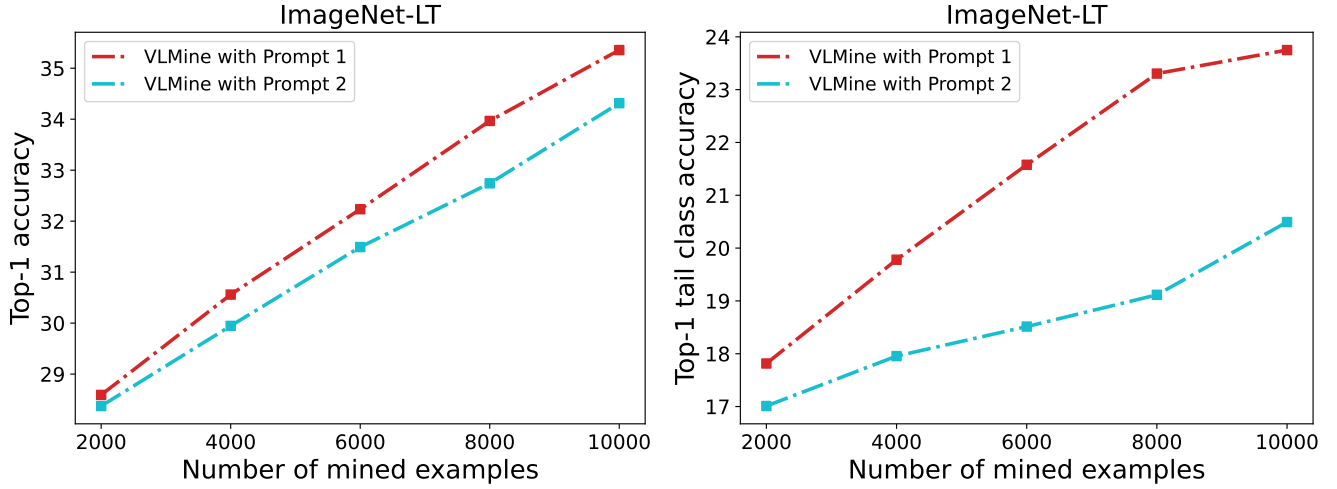


Figure 11: Ablation study on prompt engineering with ImageNet-LT.

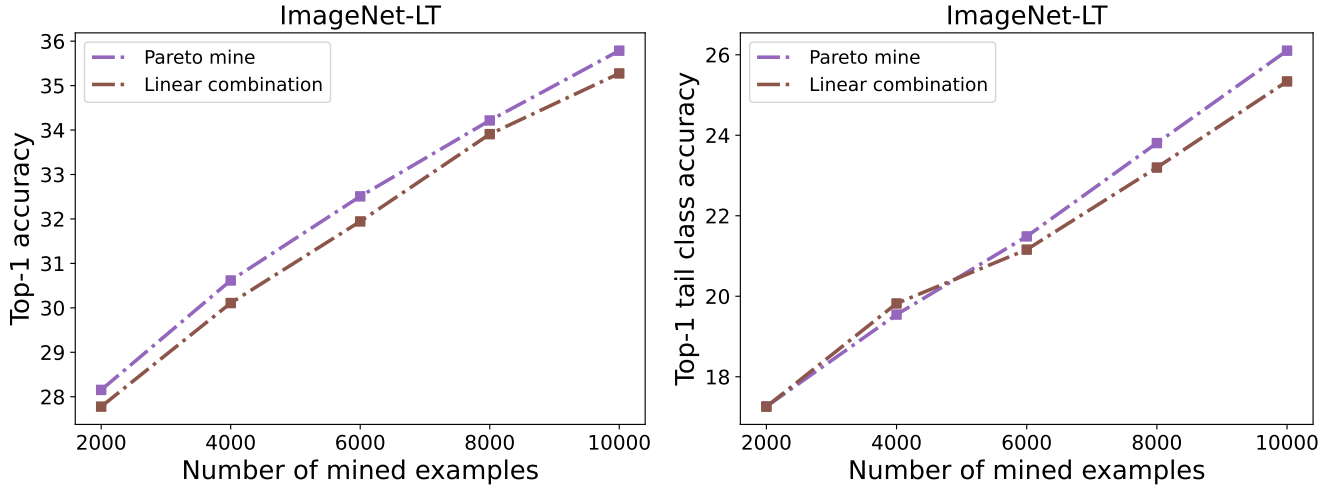


Figure 12: Comparing Pareto mining to a linear combination on ImageNet-LT.

truck, white car, concrete truck, cement truck, ...” For Figure 8b, important objects are missed: “None,” and for Figure 8c, keywords are not descriptive: “Bus.” It is possible that with a better VLM, we might reconsider one-step keyword generation, but we will leave this for future work.

A.5 Remarks on Hallucination

The VLM may hallucinate on some input images, but we find that hallucinations often have repeating patterns. For example, a typical case is to hallucinate the presence of a dog when the image shows a busy city street. Below is an example of hallucination when Figure 13 is used as the input image: “In the image, there are several cars and a bus, which are common vehicles on the road. However, there is also a dog crossing the street along with the group of pedestrians, which is an unusual sight. Additionally, there are multiple bicycles in the scene, including one with a person riding it, which is another uncommon mode of transportation in this



Figure 13: Example of a camera image from the Waymo Open Dataset.

particular scene. The presence of these unconventional vehicles and pedestrians creates a unique and lively atmosphere in the busy city street.”

Since such hallucinations share common patterns, they will appear a considerable amount of times when describing the images of a typical autonomous driving dataset. We argue that our approach is robust to these types of hallucinations for two reasons: if the hallucination is not related to the traffic object we are interested in, it will be filtered by the LLM; even if it is related to traffic objects, since it is repeated multiple times, the keywords that come from hallucination will have a relatively high frequency. Of course hallucination is not ideally, but we expect VLMine to improve as the underlining VLM improves.