# VLM-NCD:Novel Class Discovery with Vision-Based Large Language Models

Yuetong Su[1], Baoguo Wei[2], Xinyu Wang[3], Xu Li[4], and Lixin Li[5]

[1,2,3,4,5]School of Electronics and Information, Northwestern Polytechnical University, Xi'an, China
[1]syt15332639458@mail.nwpu.edu.cn
[2]wbg@nwpu.edu.cn
[3]xw18765@outlook.com
[4]lixu@nwpu.edu.cn
[5]lilixin@nwpu.edu.cn

December 12, 2025

## Abstract

Novel Class Discovery (NCD) aims to utilise prior knowledge of known classes to classify and discover unknown classes from unlabelled data. Existing NCD methods for images primarily rely on visual features, which suffer from limitations such as insufficient feature discriminability and the long-tail distribution of data. We propose LLM-NCD, a multimodal framework that breaks this bottleneck by fusing visual-textual semantics and prototype-guided clustering. Our key innovation lies in (1) modelling cluster centres and semantic prototypes of known classes by jointly optimising known class image and text features, and (2) a dual-phase discovery mechanism that dynamically separates known/novel samples via semantic affinity thresholds and adaptive clustering. Experiments on the CIFAR-100 dataset show that compared to the current state-of-the-art methods, this method achieves up to 25.3% improvement in accuracy for unknown classes. Notably, our method shows unique resilience to long-tail distributions—a first in NCD literature.

## 1 Introduction

With the widespread adoption of open-world intelligent applications, machine vision systems often face the need to identify unknown classes. Novel Class DiscoveryHan et al. [2019] aims to automatically identify and classify unknown classes from unlabelled samples using labelled data from known classes, and is one of the core problems in open-world recognitionBendale and Boult [2015]. Existing NCD methods are mainly based on deep visual feature clustering (such as DTCHan et al. [2019]), but face two major challenges:

**Insufficient feature discrimination.** Feature extractors trained on known classes in traditional models have difficulty generalising to unknown classes.

**Long-tail distribution bias.** Rare classes in real data are easily overwhelmed by high-frequency classes, leading to false negatives for novel classes.

In recent years, large-scale visual-language pre-training models (such as CLIPRadford et al. [2021]) have demonstrated powerful cross-modal semantic understanding capabilities. By training on a vast amount of image-text pairs, these models map images and text into a unified semantic space, making them naturally suited for open-world tasks involving multimodal associations. Currently, CLIP has achieved significant success in visual tasks such as image classificationZhou et al. [2022], Gao et al. [2024], Radford et al. [2021], Zhang et al. [2022], object detectionGu et al. [2021], Minderer et al. [2022], and video understandingWang et al. [2021], but its potential for application in the critical field of Novel Class Discovery has not yet been fully explored. This is primarily due to two important reasons:

**Insufficient task adaptability.** There is an inherent difference between the fixed class space of closed-set recognition tasks and the open-world characteristics required for Novel Class Discovery. Directly applying CLIP's zero-shot classification paradigm is difficult to adapt to the needs of a dynamically expanding class system.

**Lack of multimodal collaboration.** Vague descriptions of unknown classes (such as 'unknown objects') can easily lead to confusion in clustering boundaries. By fully utilising the complementarity of text and image features, it is possible to effectively develop the unique value of multimodal joint inference in Novel Class Discovery. In particular, the semantic a priori information contained in the text modality can provide key distinguishing clues for unknown classes.

For novel class discovery in open-world scenarios, we should leverage all available background knowledge or prior information beyond just the given labeled samples. Additional relevant knowledge can enhance classification confidence and thus improve model's novel class discovery capability. In this work, we employ Large Language Models as the representation of such background knowledge. There-
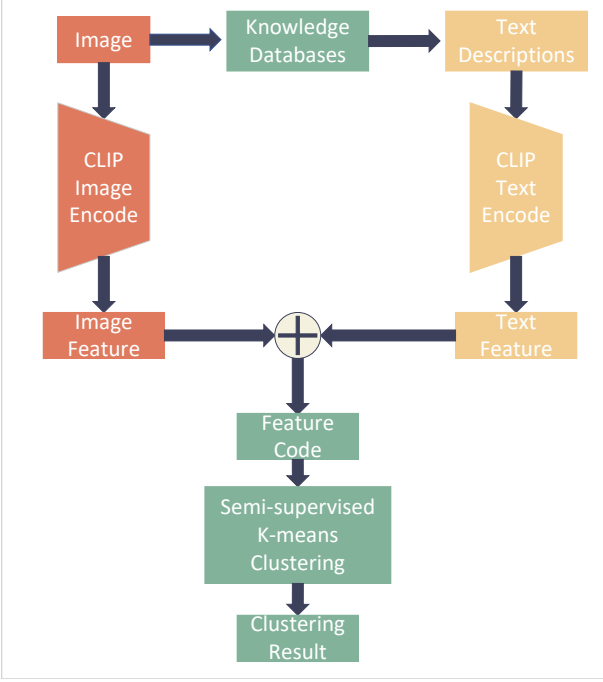
Figure 1: Methodology of this paper.

fore, to address the above issues, this paper proposes a Novel Class Discovery method based on Large Language Models, as shown in Figure 1. We propose a Novel Class Discovery framework based on multimodal semantic alignment. The core idea is to enhance class distinguishability through image-text feature fusion and to alleviate long-tail bias using a text-guided dynamic allocation strategy.

In the top section, we first utilise CLIP's aligned visual-language representations, using the input image as a query to retrieve a set of highly relevant text descriptions from a large knowledge databases. To further leverage CLIP's large-scale pre-trained representations, the input image and its retrieved text are encoded into a set of feature encodings using a frozen CLIP image and text encoder. In the bottom section, given the connected text and image views, we employ a semi-supervised k-means clustering method to cluster the features into known and unknown classes. The specific contributions are as follows:

**Multimodal feature fusion mechanism.** We propose the first multimodal NCD framework based on CLIP, which utilises joint image and text features to enhance class distinguishability. The image and text features of CLIP are concatenated to generate a joint representation, while retaining visual granularity information and textual semantic constraints, thereby improving the discriminative power for known and unknown classes.

**Semantically guided two-stage discovery.** Design a semantically guided two-stage discovery mechanism that dynamically separates known classes from novel classes based on text similarity. In the known class modelling stage, construct text semantic prototypes. In the NCD stage, first filter known class samples based on text similarity, then perform k-means clustering on the remaining samples to reduce semantic drift.

**Lightweight solution.** High performance can be achieved with only a pre-trained CLIP model, without the need for fine-tuning or complex training processes, providing a lightweight solution for open-world recognition.

Experiments show that this method outperforms existing methods with novel class clustering accuracy rates of 85.6% and 78.2% on the CIFAR-100 and ImageNet-100 benchmarks, respectively. Ablation experiments validate the critical role of Large Language Models as feature extractors and text-guided strategies for long-tail data.

## 2 Related work

### 2.1 Novel Class Discovery

Novel Class Discovery (NCD) is a relatively emerging field, initially proposed as 'cross-task transfer', which involves transferring learning from labelled data to clustering unseen classes in unlabelled data. As an important research direction in open-world recognition, it has made significant progress in recent years. Hsu et al. [2017, 2019] use a paired conjoint network trained on labelled data and applies it to a clustering network trained on unlabelled data. Subsequent researchHan et al. [2019] improved upon this approach using specialized deep clustering methods. In RankStat[Han et al. [2020, 2021], a three-stage pipeline was deployed: first, the model was pre-trained on all data for representation learning, then fine-tuned on labelled data to capture higher-level semantic knowledge, and finally, ranking statistics were used to transfer knowledge from labelled to unlabelled data. Zhong et al. [2021] proposed a contrastive learning method that generates hard negatives by mixing labelled and unlabelled data in the latent space. UNOFini et al. [2021] introduced a unified cross-entropy loss that jointly trains the model on labelled and unlabelled data by swapping pseudo-labels in the classification head. The original GCDVaze et al. [2022] method involves k-means clustering of DINO embeddings, while XConFei et al. [2022] improves these results through additional contrastive training.

However, these methods have two main limitations. First, they typically train feature extractors on limited known class data, resulting in insufficient generalisation ability for unknown classes; second, these methods fail to effectively utilise semantic association information between classes, resulting in poor performance when processing data with complex semantic structures. Our work is based on a more realistic new environmentVaze et al. [2022], where unlabelled samples can come from either known or unknown classes. In our method, we focus on utilising multimodal models in multiple ways, and we demonstrate superior results compared to all previously published techniques.

### 2.2 Large Language Models

In recent years, large-scale pre-trained models have demonstrated tremendous potential in the field of computer vision. The CLIPRadford et al. [2021] model was pre-trained on 400 million image-text pairs, establishing cross-modal

semantic alignment capabilities. This contrastive learning-based training method enables the model to achieve outstanding zero-shot transferability. In image classification tasks, Gao et al.Gao et al. [2024] validated CLIP's performance in few-shot classification tasks, but their core experiments directly froze the CLIP visual encoder (without fine-tuning) and only adapted features through a simple linear layer. Radford et al.Radford et al. [2021] incorporated the CLIP visual and text encoders into a classifier based on image-text similarity; Zhang et al.Zhang et al. [2022] froze the CLIP visual encoder, and constructing a classifier using cached features, achieving a method that successfully classifies few-shot data without training. In object detection tasks, Gu et al.Gu et al. [2021] replaced the detector backbone with the CLIP visual encoder and distilled textual knowledge into the detection head to achieve open-word detection; Minderer et al.Minderer et al. [2022] designed a text-conditioned detection framework based on the CLIP visual encoder to achieve open-word detection. In segmentation tasks, Xu et al.Xu et al. [2022] extended the CLIP visual encoder and successfully achieved zero-shot semantic segmentation through a text-supervised learning grouping mechanism.

However, the application potential of CLIP in the critical field of Novel Class Discovery has not yet been fully explored. In our method, we not only fully leverage the capabilities of the CLIP visual-language large model but also simultaneously utilise both the visual encoder and text encoder, successfully validating CLIP's outstanding capabilities in the field of Novel Class Discovery and laying a solid foundation for the use of other large models.

## 2.3 Image-text multimodal

Multimodal learning offers new approaches for NCD tasks. Han et al. [2021] is the first to use image-text embedding to discover novel classes, and combining visual and linguistic features can significantly improve model performance. In the field of image classification, CLIPRadford et al. [2021] aligns image and text embeddings through contrastive learning, enabling zero-shot classification; Jia et al.Jia et al. [2021] trained multimodal models using noisy network data to enhance classification generalisation capabilities. In the field of object detection, Li et al.Li et al. [2022] jointly trained detection and language models to achieve semantically guided object detection; Zareian et al.Gu et al. [2021] utilised CLIP knowledge distillation to achieve open-word detection. In the field of object tracking, Yang et al.Yang et al. [2020] combined language descriptions with visual features for object tracking; Sha et al.Sha et al. [2023] focused on autonomous driving tasks and utilised language models to enhance tracking semantic understanding.

However, in NCD tasks, the issues of insufficient feature discriminability and the long-tail distribution of data severely limit feature extraction capabilities for novel classes. To address this challenge, this paper proposes a solution that combines visual and semantic information. By integrating image features with class text knowledge, this approach effectively enhances the feature representa-

tion capabilities for novel classes. Specifically, the method leverages the semantic prior information provided by the text modality to supplement the features of scarce samples while also providing additional supervisory signals for tail classes under long-tail distributions, thereby improving the model's generalisation performance.

# 3 Method

In this section, we first introduce the symbols and definitions of NCD. Then, we explain how to use CLIP in NCD and introduce our method for handling this task.

## 3.1 NCD Problem setting

The dataset $D$ consists of a labelled subset of dataset $D_{\mathcal{L}} = \{(X_i, y_i)\}_{i=1}^N \in \mathcal{X} \times \mathcal{Y}_{\mathcal{L}}$ and an unlabelled dataset $D_u = \{(X_i, Y_i)\}_{i=1}^M \subseteq \mathcal{X} \times \mathcal{Y}_u$, where $\mathcal{Y}_{\mathcal{L}} \subseteq \mathcal{Y}_u$. The goal is to learn a model that can group $D_u$ instances based on information from the $D_{\mathcal{L}}$ dataset. Leveraging advancements in visual transformers and their significant performance in self-supervised learningCaron et al. [2021], Vaze et al. Vaze et al. [2022] designed a two-stage training process for the NCD task: supervised contrastive learning on labelled data and unsupervised contrastive learning on all data.

Let $X_i$ and $X_i'$ be two randomly enhanced views of the same image in a small batch $B$. The unsupervised contrast loss is defined as follows:

$$\mathcal{L}_i^u = -\log \frac{\exp(Z_i \cdot Z_i'/\tau)}{\sum_n \mathbb{1}_{[n \neq i]} \exp(Z_i \cdot Z_n'/\tau)'}$$

where $Z_i = h(f(x_i))$ is the feature of the input image $x_i$ extracted by the backbone network $f(\cdot)$ and projected onto the embedding space through the projection head $h(\cdot)$, and $Z_i'$ is the feature of another view of the input image $x_i$.

The definition of supervised contrast loss is as follows:

$$L_i^s = -\frac{1}{|N(i)|} \sum_{q \in N(i)} \log \frac{\exp(Z_i \cdot Z_q/\tau)}{\sum_n \mathbb{1}_{[n \neq i]} \exp(Z_i \cdot Z_n/\tau)'}$$

where $N(i)$ indicates the index of other images with the same label as $X_i$ in the small batch $B$ . Then, the final objective function is a combination of these two losses:

$$\mathcal{L}^t = (1 - \lambda) \sum_{i \in \mathcal{B}_{\mathcal{L}} \cup \mathcal{B}_{\mathcal{U}}} \mathcal{L}_i^u + \lambda \sum_{i \in \mathcal{B}_{\mathcal{L}}} \mathcal{L}_i^s$$

where $\lambda$ is a weighting factor, $\mathcal{B}_{\mathcal{L}}$ and $\mathcal{B}_{\mathcal{U}}$ represent small batches of labelled and unlabelled images, respectively. For label assignment, the authors propose a semi-supervised k-means algorithm, which is similar to k-meansMcQueen [1967] but with an important difference: when calculating the clustering assignment at each step, the semi-supervised k-means algorithm considers the labelled data $\mathcal{D}_{\mathcal{L}}$. This means that labelled samples will always be assigned to the correct cluster, regardless of their distance from the nearest cluster centre.

## 3.2 Our method

By combining text and image information, language-image models can achieve better performance in various tasks. Therefore, we suggested utilising CLIP's zero-shot learning capabilities and multimodal alignment encoders to address two major issues in NCD, and propose a retrieval-based data augmentation method, with the model architecture shown in Figure 2.

### 3.2.1 CLIP in NCD

We address the NCD task by leveraging CLIP's multimodal joint embeddings. The CLIP model has two branches, the image branch CLIP-Image and the text branch CLIP-Text, which encode images and text into global feature representations, respectively. CLIP is trained on a large-scale image-text pair dataset, ensuring that paired images and texts are close in the embedding space, while unpaired ones are pushed apart.

To improve data representation, especially for labelled and unlabelled data, we optimise representations by combining supervised contrastive learning on labelled data and unsupervised contrastive learning on all data. We achieve this by fine-tuning the model on the target data simultaneously. CLIP learns image representations by comparing the representation of an image with its textual description, such as 'a photo of class name'. The text description is called a prompt, and its design is crucial for improving CLIP's performance.

However, the unlabelled data contains latent novel classes that are entirely unknown — we possess no prior knowledge (e.g., class names, quantity, or semantic descriptions) to even provisionally define them, rendering prompt-based methods fundamentally inapplicable due to the absence of a valid reference for prompt construction. Therefore, inspired by research in the field of image caption generation, we generate a set of text descriptions for both labelled and unlabelled data, providing supplementary information for the input image, as shown in Figure 3, which includes details and information about the input image to map it to the feature space. Training a separate caption generation model to generate text descriptions may be expensive and challenging, so for each labelled and unlabelled image, we retrieve the k most relevant descriptions from a text corpus. Current mainstream description databases include Conceptual Captions (3M)Sharma et al. [2018], Conceptual Captions (12M)Changpinyo et al. [2021], MS CocoLin et al. [2014], and LIONSchuhmann et al. [2021], among others.

**Text description retrieval.** Given an image query, the goal is to retrieve the top k most relevant text descriptions from the description database. To achieve this goal, we propose using cross-modal joint embedding from CLIPRadford et al. [2021] for cross-modal retrieval tasks. Specifically, we use CLIP-Text to encode all descriptions in the description database as search terms. Images are encoded as queries by CLIP-Image. We then search the description database for the top k text descriptions with the highest cosine similarity scores. Examples of the top 4

results are shown in Figure 3.

**Multi-view generation.** Our feature vector extraction and multi-view generation framework is illustrated in Figure 2 (left). Given an image and a set of text descriptions, we first encode the image using the CLIP image encoder, then encode the text descriptions using the CLIP text encoder, and generate a view (sentence embedding) through mean pooling. Finally, the feature vectors of the image and text views are concatenated and projected into the CLIP latent space, where clustering is performed directly.

**Semi-supervised k-means clustering.** Given image views and text views, we concatenate the feature vectors and apply semi-supervised k-means clustering according to the method in Vaze et al. [2022] to group unlabelled data into known and unknown classes. Semi-supervised k-means clustering is a constrained version of the traditional k-means method, where the number of clusters k is assumed to be known. This requires $\mathcal{D}_{\mathcal{L}}$ data instances to be assigned to the correct cluster centres based on their true class labels. First, use real class labels to obtain $\mathcal{D}_{\mathcal{L}}$ data in semi-supervised k-means clustering centres. Second, for novel classes, use k-means++Arthur and Vassilvitskii [2006] to obtain cluster centres from the data. During the process of updating and assigning cluster centres, instances of the same class are always grouped together, while instances can be assigned to any cluster based on their distance to each cluster centre. Once the algorithm converges, each instance can be assigned a cluster label.

## 4 Experiment

### 4.1 Model Architecture Details

CLIPRadford et al. [2021] consists of two encoders, CLIP-Image and CLIP-Text, which are pre-trained Transformer models for images and text, respectively. CLIP-Text is a basic Transformer model comprising 12 layers with a hidden size of 768, and the final linear projection layer generates a representation vector of size 512. CLIP-Image is a hybrid ViT-Base model (the same as the DINO training model, for fair comparison), consisting of 12 stacked layers, starting with a convolutional layer for feature extraction. For a given image, 49 embeddings with a hidden size of 768 are generated, and the output hidden state is projected from 768 to 512 dimensions to match the output of the CLIP-Text encoder. We fine-tuned the last block of the visual transformer with an initial learning rate of $5e^{-5}$ and decayed it over time using cosine annealing scheduling. We trained the model for 100 epochs using batches of size 128 and set the $\lambda$ value in the loss function to 0.25. The optimal hyperparameters were selected by tuning and testing on a separate validation set.

### 4.2 Datasets and Evaluation metrics

**Datasets.** We evaluate our method on general image classification tasks. Following the approach outlined in Vaze et al. [2022], we selected CIFAR-10, CIFAR-100, and ImageNet-100 as general image classification datasets. We
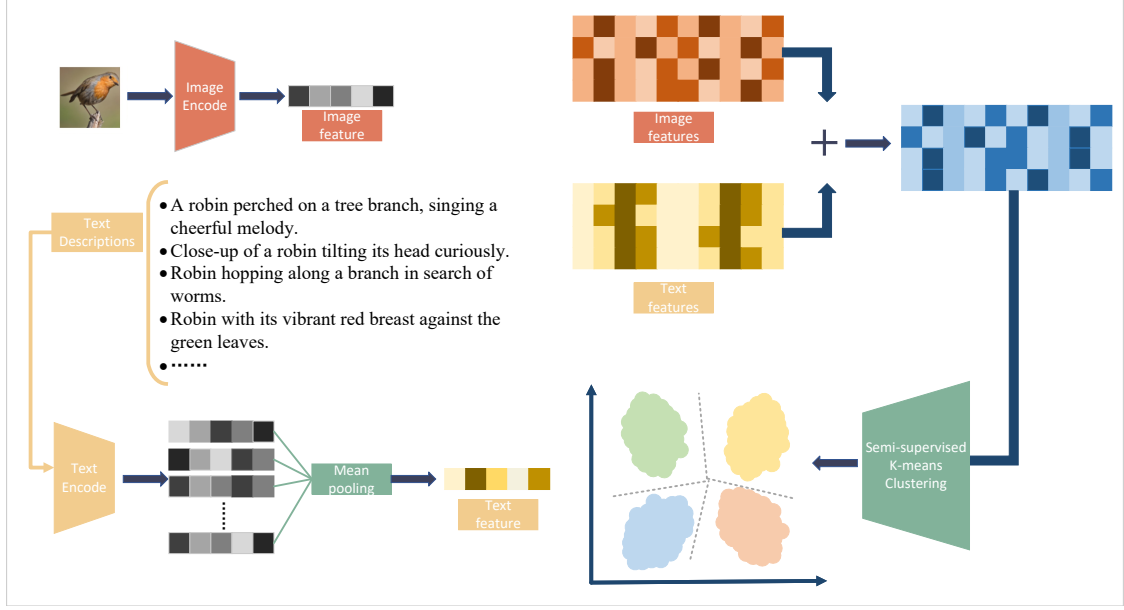
Figure 2: Model architecture. First, we propose a cross-modal retrieval module to retrieve context-aware text descriptions as supplementary information for clustering, complementing labelled and unlabelled data. Second, we concatenate image views and text views and use a semi-supervised k-means clustering algorithm to group known and unknown classes.
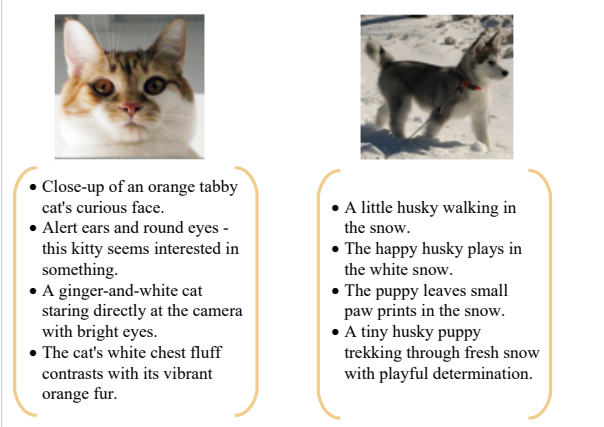


Figure 3: Shows the top 4 text descriptions most relevant to an image in the ImageNet dataset, retrieved from the Conceptual Captions (3M) database.

|  | CIFAR-10 | CIFAR-100 | ImageNet-100 |
|---|---|---|---|
| $|y_{\mathcal{L}}|$ | 5 | 80 | 50 |
| $|y_u|$ | 10 | 100 | 100 |
| $|\mathcal{D}_{\mathcal{L}}|$ | 12.5k | 20k | 31.9k |
| $|\mathcal{D}_u|$ | 37.5k | 30k | 95.3k |

Table 1: Dataset division. $(|y_{\mathcal{L}}|, |y_u|)$ represents the number of classes in the labelled and unlabelled sets. $(|\mathcal{D}_{\mathcal{L}}|, |\mathcal{D}_u|)$ represents the number of images in each set.

divided the training data into labelled datasets and unlabelled datasets. We equally divided all classes into known classes and unknown classes, then randomly sampled 50% of the images from the known classes as unlabelled data to ensure that the unlabelled dataset $\mathcal{D}_{\sqcap}$ contains images from both known and unknown classes, while the labelled dataset only contains images from known classes. These divisions are summarised in Table 1.

**Evaluation metrics.** We use clustering accuracy (ACC) to measure the performance of the model, which is defined as follows:

$$\text{ACC} = \max_{p \in P(y_U)} \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}\{y_i = p(\hat{y}_i)\}$$

where ACC is the clustering accuracy, which is cal-

culated using the Hungarian method [32] to match the model's predictions $\hat{y}_i$ with the true labels $y_i$. N is the total number of images in the unlabelled dataset. Following the approach in Vaze et al. [2022], we use this metric on three different sets, namely:

**All.** Refers to the entire unlabeled set $\mathcal{D}_u$. **Old.** Refers to images in the unlabelled dataset $\mathcal{D}_u$ that belong to $\mathcal{Y}_{\mathcal{L}}$. **New.** Refers to images in the unlabelled dataset $\mathcal{D}_u$ that belong to $y_u \setminus y_{\mathcal{L}}$.
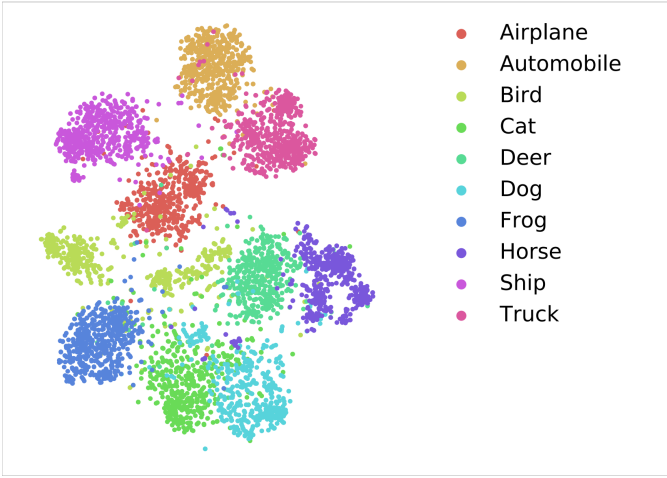
## 4.3 Comparison results

Our method is first compared with leading techniques methods on image classification. RankStats+Han et al. [2021] and UNO+Fini et al. [2021] are two competitive baseline methods for NCD, modified for the GCD setting. XConFei et al. [2022] is a method tailored for the GCD setting. The results in the General Image Recognition benchmark are shown in Table 2. Across all datasets in our experiments, our method achieves the best performance on most classes, often improving significantly over previous methods. On ImageNet-100 and CIFAR-100, our method outperforms other methods on all subsets, further confirming that the dual use of multimodal models enhances performance compared to using visual models alone.

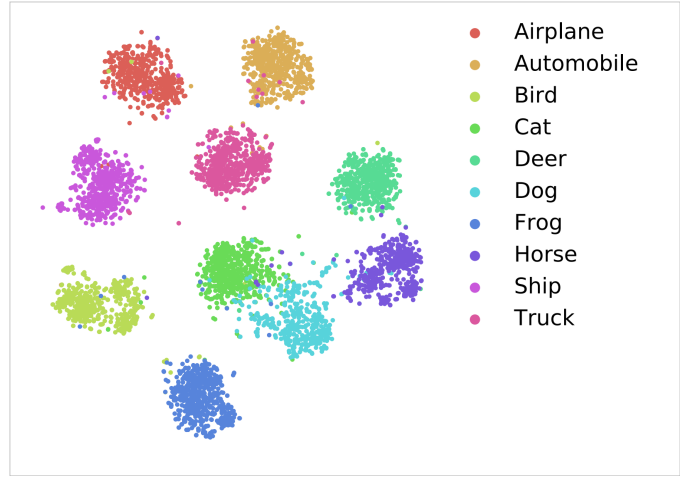| | | CIFAR-10 | | | CIFAR-100 | | | ImageNet-100 | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Class | All | Old | New | All | Old | New | All | Old | New |
| RankStats+Han et al. [2021] | 46.8 | 19.2 | 60.5 | 58.2 | 77.6 | 19.3 | 37.1 | 61.6 | 24.8 |
| UNO+Fini et al. [2021] | 68.6 | **98.3** | 53.8 | 69.5 | 80.6 | 47.2 | 70.3 | 95.0 | 57.9 |
| GCDVaze et al. [2022] | 91.5 | 97.9 | 88.2 | 73.0 | 76.2 | 66.5 | 74.1 | 89.8 | 66.3 |
| XConFei et al. [2022] | 96.0 | 97.3 | 95.4 | 74.2 | 81.2 | 60.3 | 77.6 | 93.5 | 69.7 |
| Ours | **96.6** | 97.2 | **96.4** | **85.2** | **85.0** | **85.6** | **84.0** | **95.5** | **78.2** |

Table 2: Comparison results of general image recognition datasets.

| Datasets | Image Encoders | Knowledge Databases | All | Old | New |
|---|---|---|---|---|---|
| | DINO | N | 73.0 | 76.2 | 66.5 |
| | DINO | Y | 75.9 | 79.7 | 67.3 |
| CIFAR-100 | CLIP | N | 84.2 | 83.1 | 82.3 |
| | CLIP | Y | **85.2** | **85.0** | **85.6** |

Table 3: Clustering accuracy on whether using knowledge datasets and different image encoders.



(a) t-SNE results using only image features on the CIFAR-10 dataset

(b) t-SNE results using both image and text features on the CIFAR-10 dataset

Figure 4: t-SNE results showing whether text features were used in the CIFAR-10 dataset.

| Datasets | Knowledge Databases | All | Old | New |
|---|---|---|---|---|
| | CC-12M | **85.2** | 85.0 | **85.6** |
| | CC-3M | 82.8 | 82.6 | 83.2 |
| CIFAR-100 | MSCOCO | 85.1 | **85.5** | 84.2 |
| | LAION-400M | 82.0 | 82.6 | 80.8 |
| | LAION-5B | 82.5 | 83.4 | 80.6 |

Table 4: Accuracy of models using different knowledge databases as text description sources.



Figure 5: The impact of different Top-k values on the results.

## 4.4 Analysis

We analysed the contribution of certain aspects of our method through ablation experiments. Specifically, we emphasised the importance of the impact of different image encoders, the impact of whether to use text descriptions and different description databases, and the impact of the number of texts k for each image on model accuracy.

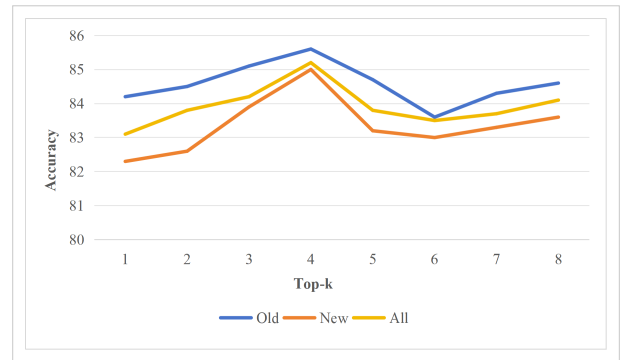As shown in Table 3, the table compares the clustering accuracy of models using different image encoders and whether they use only visual features or both visual and text features. DINO refers to the use of pre-trained weights from GCD's DINOCaron et al. [2021], while CLIP refers to the use of a pre-trained ViT-B/16 backbone network. The results show that models using both image and text features with CLIP significantly outperform those using only

image features, confirming that incorporating language information can significantly improve performance compared to image-only models. Additionally, the retrieval mechanism, which uses CLIP's cross-modal joint embeddings for text retrieval, further enhances performance.

Overall, in NCD settings, semantic information plays a crucial role in improving model representation and performance, and retrieval-based augmentation mechanisms further improve clustering accuracy compared to using visual features alone. Descriptive retrieval text is essential to our method. In typical datasets, text descriptions may vary in their association with images. Ideally, we aim to encode salient objects and meaningful details within images to enhance representation learning for object recognition tasks.

For contrastive models, the learned representations are influenced by the text encoder, suggesting that textual descriptions of image scene content can enhance the transferability of CLIP models. We validated this hypothesis and quantified the descriptiveness of textual descriptions using multiple text datasets. We conducted the first four cross-modal retrieval experiments on multiple data sources, including Conceptual Captions (3M)Sharma et al. [2018], Conceptual Captions (12M)Changpinyo et al. [2021], MS CocoLin et al. [2014], and LIONSchuhmann et al. [2021], and recorded the model's accuracy on All, Old, and New subsets for each dataset. Table 4 shows the model's results on the CIFAR100 dataset.

## 4.5   Qualitative results

We further demonstrate the results of t-SNEKuhn [1955] projections using only image features and simultaneously using image and text features on the CIFAR-10 dataset, as shown in Figure 4, where we present the clustering characteristics of unlabelled data. When using only image features, data points of the same class are typically projected near each other, forming distinct clusters, but there is some overlap between classes. In contrast, when using both image and text features, the image-text features form distinct clusters with more pronounced separation, further confirming the utility of language in this setting.

In addition, we examined the model's sensitivity to the number of images retrieved (top-k) and demonstrated its accuracy using Conceptual Captions (12M)Changpinyo et al. [2021] on the CIFAR-100 dataset. Figure 5 shows that some text descriptions may not contain useful information, causing the model's accuracy to plateau or even decline after a certain number.

## 5   Conclusion

In this paper, we propose a new Large Language Models-based approach to address the Novel Class Discovery problem, leveraging the multimodal capabilities of the CLIP model. First, we utilise the image encoder of CLIP to establish a robust baseline for NCD. Second, we introduce a supplementary retrieval-based augmentation method, specifically retrieving text descriptions from a text corpus

and jointly clustering image and text embeddings. We conducted a rigorous analysis to demonstrate that our method is highly suitable for the NCD setting.

We present quantitative results on three general classification datasets, achieving significant performance improvements over previous methods. Additionally, we demonstrate that the two components of CLIP are complementary and necessary, working together to achieve strong performance. However, there are still some limitations and future work, including improving the retrieval process to enhance the quality of retrieved contextual knowledge and conducting experimental validation on other fine-grained datasets.

## References

David Arthur and Sergei Vassilvitskii. k-means++: The advantages of careful seeding. Technical report, Stanford, 2006.

Abhijit Bendale and Terrance Boult. Towards open world recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1893–1902, 2015.

Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.

Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3558–3568, 2021.

Yixin Fei, Zhongkai Zhao, Siwei Yang, and Bingchen Zhao. Xcon: Learning with experts for fine-grained category discovery. *arXiv preprint arXiv:2208.01898*, 2022.

Enrico Fini, Enver Sangineto, Stéphane Lathuilière, Zhun Zhong, Moin Nabi, and Elisa Ricci. A unified objective for novel class discovery. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9284–9292, 2021.

Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *International Journal of Computer Vision*, 132(2):581–595, 2024.

Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. *arXiv preprint arXiv:2104.13921*, 2021.

Kai Han, Andrea Vedaldi, and Andrew Zisserman. Learning to discover novel visual categories via deep transfer clustering. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8401–8409, 2019.

Kai Han, Sylvestre-Alvise Rebuffi, Sebastien Ehrhardt, Andrea Vedaldi, and Andrew Zisserman. Automatically discovering and learning new visual categories with ranking statistics. *arXiv preprint arXiv:2002.05714*, 2020.

Kai Han, Sylvestre-Alvise Rebuffi, Sebastien Ehrhardt, Andrea Vedaldi, and Andrew Zisserman. Autonovel: Automatically discovering and learning novel visual categories. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):6767–6781, 2021.

Yen-Chang Hsu, Zhaoyang Lv, and Zsolt Kira. Learning to cluster in order to transfer across domains and tasks. *arXiv preprint arXiv:1711.10125*, 2017.

Yen-Chang Hsu, Zhaoyang Lv, Joel Schlosser, Phillip Odom, and Zsolt Kira. Multi-class classification without multi-class labels. *arXiv preprint arXiv:1901.00544*, 2019.

Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pages 4904–4916. PMLR, 2021.

Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.

Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10965–10975, 2022.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.

James B McQueen. Some methods of classification and analysis of multivariate observations. In *Proc. of 5th Berkeley Symposium on Math. Stat. and Prob.*, pages 281–297, 1967.

Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection. In *European conference on computer vision*, pages 728–755. Springer, 2022.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021.

Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.

Hao Sha, Yao Mu, Yuxuan Jiang, Guojian Zhan, Li Chen, Chenfeng Xu, Ping Luo, Shengbo Eben Li, Masayoshi Tomizuka, Wei Zhan, et al. Large language models as decision makers for autonomous driving. 2023.

Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018.

Sagar Vaze, Kai Han, Andrea Vedaldi, and Andrew Zisserman. Generalized category discovery. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7492–7501, 2022.

Mengmeng Wang, Jiazheng Xing, and Yong Liu. Actionclip: A new paradigm for video action recognition. *arXiv preprint arXiv:2109.08472*, 2021.

Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit: Semantic segmentation emerges from text supervision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18134–18144, 2022.

Zhengyuan Yang, Tushar Kumar, Tianlang Chen, Jingsong Su, and Jiebo Luo. Grounding-tracking-integration. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(9):3433–3443, 2020.

Renrui Zhang, Wei Zhang, Rongyao Fang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. Tip-adapter: Training-free adaption of clip for few-shot classification. In *European conference on computer vision*, pages 493–510. Springer, 2022.

Zhun Zhong, Enrico Fini, Subhankar Roy, Zhiming Luo, Elisa Ricci, and Nicu Sebe. Neighborhood contrastive learning for novel class discovery. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10867–10875, 2021.

Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022.