# ADAM: Autonomous Discovery and Annotation Model using LLMs for Context-Aware Annotations

**Amirreza Rouhi**
Department of Electrical and Computer Engineering
Drexel University
Philadelphia, PA 19104
ar3755@drexel.edu

**Solmaz Arezoomandan**
Department of Electrical and Computer Engineering
Drexel University
Philadelphia, PA 19104
sa3747@drexel.edu

**Knut Peterson**
Department of Electrical and Computer Engineering
Drexel University
Philadelphia, PA 19104
kp3275@drexel.edu

**Joseph T. Woods**
Department of Electrical and Computer Engineering
Drexel University
Philadelphia, PA 19104
jw3897@drexel.edu

**David K. Han**
Department of Electrical and Computer Engineering
Drexel University
Philadelphia, PA 19104
dkh42@drexel.edu

## Abstract

Object detection models typically rely on predefined categories, limiting their ability to identify novel objects in open-world scenarios. To overcome this constraint, we introduce ADAM: Autonomous Discovery and Annotation Model, a training-free, self-refining framework for open-world object labeling. ADAM leverages Large Language Models (LLMs) to generate candidate labels for unknown objects based on contextual information from known entities within a scene. These labels are paired with visual embeddings from CLIP to construct an Embedding-Label Repository (ELR) that enables inference without category supervision. For a newly encountered unknown object, ADAM retrieves visually similar instances from the ELR and applies frequency-based voting and cross-modal reranking to assign a robust label. To further enhance consistency, we introduce a self-refinement loop that re-evaluates repository labels using visual cohesion analysis and kNN-based

majority relabeling. Experimental results on the COCO and PASCAL datasets demonstrate that ADAM effectively annotates novel categories using only visual and contextual signals—without requiring any fine-tuning or retraining.

# 1 Introduction

Humans often learn through context, as shown by their early language acquisition skills. If a child encounters a new word, such as "garnet," and hears, "The garnet shimmered brightly in the sunlight," and another time, "The garnet was like ruby in color," their understanding of meaning develops from the context in which that word is being presented. Over time, the child learns a "garnet" is an expensive, red stone—not by direct instruction but through inference [32, 47]. This principle is central to the ADAM methodology: open-world learning from context. ADAM aims to mimic this human capability—the learning of previously unfamiliar entities in an open-world environment.

Modern object detection models have achieved considerable success in identifying objects from predefined categories in controlled environments [26, 4, 5, 17, 37, 42]. However, their utility remains confined to closed-set settings, where all test-time categories are known at training time. These models fundamentally rely on supervised learning over fixed label spaces, rendering them ineffective in open-world environments where novel or long-tail object categories frequently appear. Recent advances in open-world recognition attempt to address these limitations by leveraging large-scale vision-language models (VLMs) [36, 24], prompt-driven classification [14], or knowledge-enhanced object detectors [50]. While effective, these methods often require costly fine-tuning, handcrafted prompts, or auxiliary annotations to adapt to unseen categories. Moreover, many approaches still assume the presence of training data for base categories or rely on model retraining, restricting real-time adaptability.

To address this challenge, we introduce **ADAM**: a *training-free*, *self-refining*, and *context-aware* framework for open-world object recognition. ADAM represents a paradigm shift in object detection by moving beyond category-constrained supervision. Instead, it mimics the human-like process of inference through context. ADAM applies this principle to visual recognition. Rather than relying on a fixed label set, ADAM constructs an **Embedding-Label Repository (ELR)** without any gradient-based updates or task-specific retraining. Starting with region proposals from a pretrained detector, ADAM uses a large language model (LLM) to assign provisional semantic labels based on contextual image captions. These labels and their associated embeddings populate the ELR. Crucially, the system includes a *self-refinement mechanism*: by examining the visual-semantic similarity among embeddings within each proposed class, ADAM revisits noisy predictions and reassigns labels using majority voting over the $k$-nearest neighbors in the embedding space.

**Our key contributions are as follows:**

- We propose a novel zero-shot framework for open-world object labeling that bypasses the need for supervised training on unknown categories.

- We design a context-aware prompt construction mechanism that integrates visual descriptors and spatial layouts to effectively guide large language models.

- We introduce an embedding-label memory module and a two-stage refinement process combining frequency-based aggregation and cross-modal re-ranking to improve robustness and disambiguate noisy predictions.

- We conduct extensive experiments on both COCO and PASCAL VOC datasets, demonstrating that ADAM can accurately annotate novel object categories using only contextual and visual signals, achieving competitive results without any fine-tuning of vision or language backbones.

# 2 Related Work

## 2.1 Open Vocabulary and Open World Object Detection

In recent years, object detection has moved beyond supervised learning, instead focusing on open-set detection. Many current approaches perform Open Vocabulary Detection (OVD), which takes a list of object labels for classification at inference time [51]. The limitations of these models such

as OWL-ViT, GLIP, Detic, and YOLO-World is that the user must know in advance what objects may appear in the scene [30, 24, 54, 9]. To circumvent the label list, Open World Object Detection (OWOD) methods such as ORE were developed [21]. A major limitation of these existing OWOD methods is that they require a human oracle to label unknown objects and iteratively retrain the model, which can be costly.

Later works, such as Maaz et. al. [29] and OW-DETR [15] focus on developing class-agnostic object detectors to identify unknown objects, but do not attempt to predict labels for them. OSODD [53] expands work in unknown object detection to cluster previously unseen classes, allowing for novel class grouping and discovery in an unsupervised manner, but they also do not predict labels for the new classes. Many of these works have made great progress in novel class detection and classification from provided sets of labels, however, **none of these works are capable of predicting novel labels for previously unseen object classes without a class list or oracle**.

## 2.2 Vision-Text Features and Contextual Reasoning

When it comes to predicting novel text labels for new classes, most models are limited to the open-vocabulary approach by using long lists of class labels and matching them to images, such as in CLIP [35] or RAM [52]. Other more recent works have begun exploring ways to surpass this limitation by using LLMs, such as in RAM++ [18] which uses LLMs to add additional descriptors based on ground truth tag information, expanding the set of tags that can be matched to images. DVDet [19] uses a different approach, by matching visual embeddings with fine-grained text descriptors of object parts to improve detections. This approach has also been explored to extract text features of objects, such as texture with CLIP [48], or general descriptors in OvarNet [6]. While extracting text features of objects does not directly provide new object class labels, additional textual information can be used to aid in novel label generation.

Other methods combine image and language more effectively into single architectures, such as BLIP [23] and LLaVa [28], both of which excel at tasks such as visual question answering. These models do not rely on class lists for image captioning, allowing them to generate novel captions of image contents from scratch. However, both models require massive amounts of supervised training data to achieve effective results. Another approach to leveraging LLMs is to use their contextual knowledge to help make predictions where existing labels or visual features are not sufficient. Rouhi et. al. [39] explored this approach using LLaMA [43] to enhance object detection by generating contextually aware labels for occluded or poorly visible objects, achieving significant improvements.

## 2.3 Label Refinement and Consistency in Embedding Spaces

Clustering is a core unsupervised technique for grouping data in high-dimensional spaces without ground truth labels [34, 40, 49], with applications in data organization [10], image retrieval [46, 2], and anomaly detection [22, 3]. Contrastive learning methods like SimCLR [8], SCAN [44], and SPICE [33] have shown improved clusterability of learned representations. Advanced clustering models such as DeepDPM [38], ClusterGAN [31], DINO-ViT [1], and ClusterFormer [25] leverage pretrained or generative models. However, clustering often struggles with scalability and requires predefined cluster counts [16], while methods like the elbow criterion are costly at scale [41].

To overcome these limitations, Nearest Neighbor Search (NNS) is a flexible alternative as it retrieves similar data based on distance metrics without requiring global structure. Works such as [13, 1] have applied NNS for embedding-based retrieval using cosine similarity. Tools like FAISS [11, 20] and SPANN [7] offer scalable and memory-efficient NNS solutions.

## 3 Methodology

Our method of labeling unknown objects in open-world settings leverages contextual information and visual cues to infer object identities without relying on predefined categories. ADAM addresses this challenge through a training-free, self-refining framework that associates unknown object regions with candidate labels based on contextual and visual similarity.

## 3.1 Theoretical Motivation

In open-world labeling, inferring a semantic label $Y$ benefits from contextual variables $X_1, \ldots, X_n$ by reducing its conditional entropy:

$$H(Y \mid X_1, \ldots, X_n) \leq H(Y \mid X_1, \ldots, X_{n-1}).$$

This principle—central to information theory—states that each additional relevant context lowers the uncertainty in $Y$. It underpins our context-aware refinement process (see Supplementary Material).

## 3.2 System Overview

ADAM is composed of two primary stages: (1) building an Embedding-Label Repository (ELR) and (2) predicting labels for novel unknown objects using similarity-based retrieval and refinement (Figure 1 and Figure 2). In the first stage, ADAM constructs the ELR by associating visual embeddings of masked unknown regions with candidate labels predicted by a large language model (LLM) prompted with contextual information. In the second stage, when a new unknown object is encountered, its embedding is compared against the repository to retrieve similar entries. The associated label predictions are then refined using frequency-based aggregation and cross-modal reranking.

## 3.3 Generating the Embedding-Label Repository (ELR)

The embedding-label repository serves as the foundational knowledge base for ADAM, storing visual embeddings of unknown objects and their corresponding predicted labels generated by the LLM. This process is divided into two main components:

### 3.3.1 Context-Aware Object Prediction (COP) Module

The main goal of the COP module is to generate a list of predicted labels for an unknown object. Given an input image, objects are categorized into two groups: a list of detected objects with known labels and a separate list of detected, yet unknown objects.

**Textual Characterization of Object Features.** The cropped image $I_u$ of the unknown object $o_u$ is extracted from the input image $\mathcal{I}$ as follows:

$$I_u = \text{Crop}(\mathcal{I}, \mathcal{B}_u). \tag{1}$$

Multiple characteristics are then selected by creating a different caption for each text feature within a given characteristic, and then finding the similarities between embeddings of potential captions and the cropped image embedding using CLIP to select the most relevant features:

$$\mathbf{C}_{u,\text{texture}} = \text{CLIP}_{\text{visual}}^{\text{texture}}(I_u), \mathbf{C}_{u,\text{color}} = \text{CLIP}_{\text{visual}}^{\text{color}}(I_u), \mathbf{C}_{u,\text{material}} = \text{CLIP}_{\text{visual}}^{\text{material}}(I_u). \tag{2}$$

The aggregated description of the unknown object is:

$$\mathbf{C}_u = \{\mathbf{C}_{u,\text{texture}}, \mathbf{C}_{u,\text{color}}, \mathbf{C}_{u,\text{material}}, \ldots\}. \tag{3}$$

Details of these extractions and a full list of text characteristics (e.g. shape, pattern, ...) are detailed in the supplementary materials.

**Prompt Construction.** To generate a prompt for the LLM, we include the labels and bounding box coordinates of the known objects, represented as $\{L_i, \mathcal{B}_i \mid o_i \in \mathcal{O}_{\text{known}}\}$, the aggregated characteristics of the unknown object extracted from the CLIP image encoder, $\mathbf{C}_u$, and the bounding box $\mathcal{B}_u$ of the unknown object. The resulting textual prompt $\mathcal{P}$ is structured as:

$$\mathcal{P} = \text{GeneratePrompt}(\{L_i, \mathcal{B}_i\}, \mathbf{C}_u, \mathcal{B}_u). \tag{4}$$

An example of a complete prompt can be seen in Figure 1.

**Candidate Label Prediction.** The prompt $\mathcal{P}$ is input to the LLM, which generates a list of candidate labels for the unknown object:

$$L_{u,j} = \text{LLM}(\mathcal{P}), \tag{5}$$

where $L_{u,j} = \{\text{label}_1, \text{label}_2, \ldots, \text{label}_m\}$ is a set of $m$ candidate labels. For this study, the output of the LLM for each unknown object is a list of 50 potential object labels ($m = 50$).

### 3.3.2 Visual Embedding Generation with Image Encoder

The second component of embedding-label repository generation focuses on extracting visual embeddings for the unknown objects using an image encoder. In our framework, we utilize the CLIP image encoder for this purpose.

**Visual Embedding Computation.** The cropped image $I_u$ of the unknown object is input to the CLIP image encoder to generate its visual embedding:

$$\vec{v}_u = \text{CLIP}(I_u). \tag{6}$$

The resulting embedding $\vec{v}_u$ represents the unknown object in a high-dimensional feature space and is used in subsequent similarity searches and label assignments.

The embedding-label repository is maintained as a collection of instance visual embeddings $\{\vec{v}_i \in \mathbb{R}^d\}$, represented as a matrix:

$$V = \begin{bmatrix} \vec{v}_1 \\ \vec{v}_2 \\ \vdots \\ \vec{v}_N \end{bmatrix}, \quad V \in \mathbb{R}^{N \times d}, \tag{7}$$

where $N$ is the total number of objects for which embeddings have been extracted, and $d$ is the dimensionality of each embedding. For each visual embedding, the corresponding LLM-predicted label list is stored in the embedding-label repository as shown in Figure 1.

### 3.4 Label Prediction for New Unknown Objects

Assume we have a new unknown object $o_{u'}$ that we want to label. To begin, the object is cropped from the image, and its visual embedding, $\vec{v}_{u'}$, is calculated based on equations 1 and 8.

**Similarity Search.** We do not directly apply a clustering algorithm, and instead rely on the natural formation of distinct clusters within the CLIP embeddings and use nearest neighbor search to identify similar visual embeddings. Given the query image embedding $\vec{v}_{u'}$, the objective is to identify the top $k$ most similar image indices within the repository. Using FAISS and cosine similarity, this search is defined as:

$$\{i_1, i_2, \ldots, i_k\} = \arg \text{top-}k(S(\vec{v}_{u'}, V)), \tag{8}$$

where $S$ denotes the cosine similarity function. The indices $\{i_1, i_2, \ldots, i_k\}$ correspond to the $k$ most similar embeddings in $V$.

**Retrieving Corresponding Labels.** For the query embedding $\vec{v}_{u'}$, the top $k$ most similar embeddings are identified from the repository, indexed as $\{i_1, i_2, \ldots, i_k\}$. Each $\vec{v}_i$ represents the visual embedding of a single object stored in the repository. The corresponding label lists associated with these embeddings are then retrieved as:

$$\mathcal{L}_{\text{similar}} = \{L_{i_1,j}, L_{i_2,j}, \ldots, L_{i_k,j}\}, \tag{9}$$

where $L_{i,j}$ is the list of candidate labels generated by the LLM for the specific object represented by the embedding $\vec{v}_i$.

**Frequency-Based Ranking.** The extracted label lists $\mathcal{L}_{\text{similar}}$ are combined and sorted based on the frequency of occurrence of each label across all retrieved lists. Labels that occur in less than 50% of the label lists are removed. The result is a ranked list of labels $L_{\text{sorted}}$, with the most frequent label appearing first:

$$L_{\text{sorted}} = \text{RankByFrequency}(\mathcal{L}_{\text{similar}}). \tag{10}$$

To ensure robustness, we retain only those labels that appear in more than half of the retrieved lists.

**Cross-Modal Reranking** We apply CLIP to perform cross-modal reranking [45] of the candidate labels in $L_{\text{sorted}}$ to improve classification results. For each label $l \in L_{\text{sorted}}$, the CLIP text encoder computes its text embedding:

$$\vec{v}_{\text{text}}(l) = \text{CLIP}_{\text{text}}(l). \tag{11}$$

The cross-modal similarity between the previously computed visual embedding $\vec{v}_{u'}$ of the unknown object and each text embedding $\vec{v}_{\text{text}}(l)$ is calculated as:

$$\text{score}(l) = \cos(\vec{v}_{u'}, \vec{v}_{\text{text}}(l)). \tag{12}$$
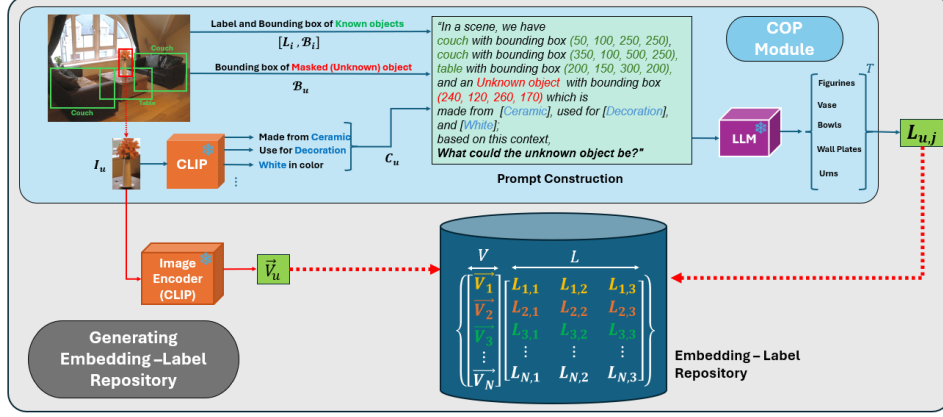
Figure 1: **Generating the Embedding-Label Repository:** This figure illustrates the pipeline for generating the embedding-label repository in ADAM. The Context-Aware Object Prediction (COP) module uses known objects' labels and locations, along with characteristics of the unknown object extracted using CLIP, to construct a context-rich prompt. The prompt is input to an LLM, which generates candidate labels ($L_{u,j}$) for the unknown object. Simultaneously, the unknown object's visual embedding ($\vec{v}_u$) is computed using the CLIP image encoder. The embedding and corresponding labels are stored in the embedding-label repository for future retrieval and refinement.
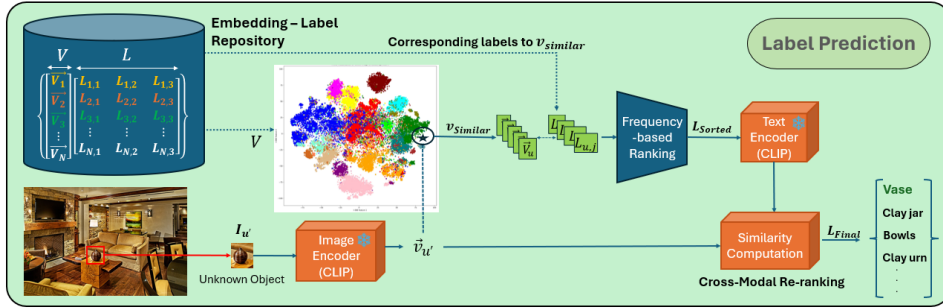


Figure 2: **Schematic of the Label Prediction Process for New Unknown Objects:** The unknown object $I_{u'}$ is processed through a CLIP image encoder to extract its visual embedding $\vec{v}_{u'}$. Similar embeddings $\vec{v}_{\text{similar}}$ are retrieved from the embedding-label repository, and their associated candidate labels are ranked using frequency-based ranking. Finally, CLIP's text encoder performs cross-modal reranking to refine the labels, yielding the top prediction $L_{\text{final}}$.

The label with the highest similarity score is then selected as the final prediction:

$$L_{\text{final}} = \arg \max_{l \in L_{\text{sorted}}} \text{score}(l). \tag{13}$$

This step ensures that the final label aligns strongly with the visual embedding of the unknown object.

### 3.5 Self-Refining Embedding-Label Repository

To enhance the consistency of label assignments, we extend ADAM with a self-refining mechanism that detects and corrects noisy predictions in the ELR. This process is entirely training-free and operates through unsupervised visual validation and neighborhood-based relabeling.

**Intra-Class Cohesion Analysis.** After the initial construction of the repository, we evaluate the visual consistency within each predicted class by computing intra-class cohesion. For each predicted label $l$, we collect all associated embeddings $\mathcal{V}_l = \{\vec{v}_i \mid l \in L_{i,j}\}$ and calculate the average pairwise cosine similarity:

$$\text{Cohesion}(l) = \frac{1}{|\mathcal{V}_l|^2} \sum_{\vec{v}_i, \vec{v}_j \in \mathcal{V}_l} \cos(\vec{v}_i, \vec{v}_j), \tag{14}$$

which quantifies the internal visual agreement among instances of class $l$.

**Outlier Detection and Flagging.** For each embedding $\vec{v}_u \in \mathcal{V}_l$, we compute its average similarity to all other embeddings in the same class. If this value falls below the class-level cohesion score, the embedding is flagged for potential relabeling:

$$\mathcal{V}_{\text{flagged}} = \{\vec{v}_u \in \mathcal{V}_l \mid \text{avg\_sim}(\vec{v}_u) < \text{Cohesion}(l)\}. \tag{15}$$

**Local Relabeling via K-NN Voting.** Each flagged embedding $\vec{v}_u \in \mathcal{V}_{\text{flagged}}$ is re-evaluated by querying its $k_{\text{refinement}}$ nearest neighbors in the embedding space using cosine similarity:

$$\{i_1, i_2, \ldots, i_{k_{\text{refinement}}}\} = \arg \text{top-}k(S(\vec{v}_u, V)). \tag{16}$$

We aggregate the associated label lists of the retrieved neighbors and apply majority voting to assign a new label:

$$L_{\text{new}} = \text{MajorityVote}(\{L_{i_1,j}, L_{i_2,j}, \ldots, L_{i_{k_{\text{refinement}}},j}\}). \tag{17}$$

If $L_{\text{new}} \neq L_{\text{old}}$, we update the label in the repository accordingly.

Throughout all experiments involving repository self-refinement, we set $k_{\text{refinement}} = 8$ for the nearest neighbor voting step. This value was empirically selected to balance label stability and correction strength, and further discussion of this choice is provided in the supplementary material.

**Iterative Refinement.** This self-refinement process is performed iteratively to ensure semantic consistency across the repository. At each iteration:

1. Recompute class-level cohesion scores.

2. Identify new outliers using updated class assignments.

3. Relabel flagged instances via $k_{\text{refinement}}$-NN majority voting.

The process terminates when either no label changes occur or at a fixed number of iterations $T$.


# 4 Experiments

## 4.1 Dataset and Evaluation Protocol

To test the ADAM framework, we conducted a comprehensive experimental study using both the COCO 2017 [27] and PASCAL VOC 2012 [12] datasets.

The **Emedding-Label Repository** was built using COCO training set. Ground truths were chosen to limit the interference from false positives that occur using an object detector. For an object class $o_u$, we mask its annotations across all images, treating it as the unknown. We then use the remaining set of annotations as the known objects and follow the repository generation process described in Section 3.3. To complete the repository, this process is repeated for every $o_u$ in the dataset and the results are concatenated together. For generating predicted labels, we utilized the LLaMA v3.2 LLM.

With 860,001 annotations in the training set, the Embedding-Label Repository size is $860,001 \times (50 + 768)$, where each annotation is represented by $m = 50$ LLM-predicted labels and a $d = 768$-dimensional visual embedding. For the **Label Prediction** step, we set $k = 250$ to retrieve the top $k$ visually similar embeddings during the similarity search. This value was empirically chosen according to section 4.4 (Ablation Study).

During evaluation, we used the validation sets of COCO and PASCAL VOC.

It is important to note that our results were obtained without any fine-tuning of the pre-trained CLIP or LLaMA models. For experimentation, we utilized three NVIDIA RTX 3090 GPUs. The average response time from LLaMA for each prompt input was 7.67 seconds.

A key challenge in evaluation is the mismatch between the closed-set evaluation labels and the open-set outputs generated by the LLM, which can include synonyms or semantically similar terms (e.g., "Television" and "TV"). To address this, predicted labels were mapped to ground truth labels using CLIP's text encoder and cosine similarity. Labels with a similarity of 0.7 or higher to the ground truth were considered correct.
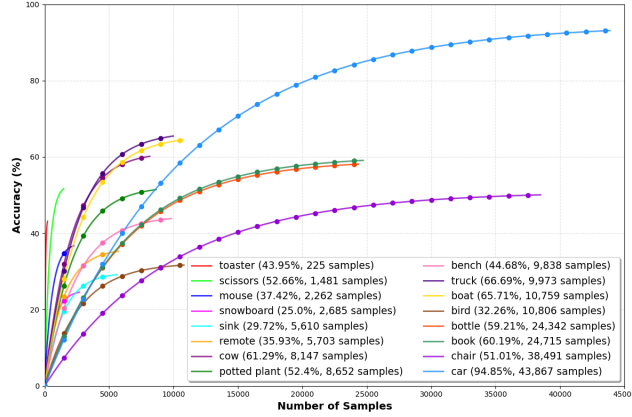
Figure 3: Top-1 accuracy of ADAM for different object categories as a function of the number of same-class samples in the embedding-label repository (**Before Self-Refining ELR**).

## 4.2 Results and Performance Analysis

The evaluation results of ADAM are presented in Table 1, which reports the accuracy of the framework for Top-1, Top-3 and Top-5 predictions. Information on accuracy by class is available in the supplemental materials. The mean accuracy for all object categories demonstrates that ADAM achieves promising results in labeling unknown objects, with a mean Top-1 accuracy of 61.30%. For comparison, we apply either CLIP or BLIP-VQA to the ground truth bounding boxes, providing the list of all COCO labels as input. CLIP significantly underperforms ADAM with a Top-1 accuracy of only 45.65%. BLIPachieves the best performance in Top-3 and Top-5, but still falls behind ADAM in Top-1. However, both CLIP and BLIP-VQA require an explicit list of candidate labels or answers at inference time, which limits their applicability in truly open-world scenarios. In contrast, ADAM does not rely on any label list input. Instead, it generates candidate labels dynamically using contextual reasoning and LLM predictions, making it well-suited for environments where novel categories emerge without prior specification. These results validate the effectiveness of the proposed label prediction framework in open-world object detection tasks.

Table 1: Mean accuracy on the classification of COCO val boxes. Note that ADAM's refinement iterations only apply to its Top-1 accuracy (61.3%).

| Method | Top-1 Accuracy | Top-3 Accuracy | Top-5 Accuracy |
|--------|----------------|----------------|----------------|
| CLIP | 45.65 | 64.23 | 71.55 |
| BLIP-VQA | 60.18 | **92.12** | **98.16** |
| ADAM | **61.30** | 70.64 | 75.95 |

### 4.2.1 Effect of Repository Size on Performance

Figure 3 presents the performance of ADAM as a function of the number of same-class samples in the original embedding-label repository, with each graph line representing a unique object. The results clearly indicate that increasing the number of samples improves the performance of the model.

This observation aligns with the hypothesis that a larger repository provides more representative embeddings and reduces contextual ambiguity for label prediction.

This trend underscores the importance of repository size in the proposed framework in open-world object detection tasks.

### 4.2.2 Effect of Contextual Information on Accuracy

The availability of contextual information significantly influences the accuracy of ADAM's label predictions for unknown objects. Table 2 illustrates the relationship between the number of known

Table 2: Accuracy of ADAM vs. Number of Known Objects in a Scene (**Before Self-Refining ELR**).

| Number of Known Objects | 0 | 1–2 | 3–4 | 5–8 | 9+ |
|---|---|---|---|---|---|
| **Average Accuracy (%)** | 6.2 | 30.4 | 38.2 | 56.7 | **58.2** |

objects in a scene and the model's prediction accuracy. Scenes with no known objects achieve a Top-1 accuracy of only 6%, as in this case, the method relies solely on the visual characteristics of the unknown object. However, when one to two known objects are present, the accuracy improves to 30%. As the number of known objects increases beyond eight, the accuracy further rises to 58%, highlighting the critical role of contextual relationships in enhancing label prediction performance.

## 4.3 Evaluating with Region Proposals from Faster R-CNN

To evaluate ADAM in a practical detection pipeline and assess its generalization to a new dataset, we apply it to the PASCAL VOC dataset using region proposals from a pretrained Faster R-CNN model [37]. The embedding-label repository built from COCO is reused here. More details on this process are provided in the supplementary material. For comparison, we apply CLIP and BLIP to the same RPN outputs and provide them with the full COCO label list (which includes all PASCAL categories). In contrast, ADAM operates without any fixed label list.

As shown in Table 3, ADAM outperforms CLIP, yet BLIP has a slight edge in average precision. However, ADAM has better performance than BLIP on most classes, surpassing BLIP in 14 out of 20 classes. This demonstrates ADAM's strong transferability and its effectiveness as a training-free, open-world labeling solution.

Table 3: Precision Comparison by Class with FasterRCNN on the PASCAL validation set

| Metric | person | cat | dog | car | sofa | mbike | train | plane | bus | horse | bicycle | chair | bird | table | boat | cow | plant | sheep | tv | bottle | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| FasterRCNN+CLIP | **94.8** | 80.0 | 73.1 | **95.2** | 62.8 | 90.1 | 95.8 | 90.1 | 96.9 | 90.3 | 77.9 | 10.1 | 85.2 | 47.8 | 83.4 | 92.5 | **93.3** | 86.0 | 17.2 | 47.0 | 75.4 |
| FasterRCNN+BLIP | 92.3 | 79.5 | 69.9 | 92.8 | **85.4** | 74.6 | 78.8 | 95.1 | 45.1 | 89.3 | 75.6 | **74.2** | 88.6 | **61.6** | 80.0 | 88.8 | 88.6 | **95.1** | **71.1** | 60.5 | **79.3** |
| FasterRCNN+ADAM | 92.8 | **82.3** | **75.1** | 87.5 | 30.2 | **90.7** | **99.4** | **97.7** | **97.4** | **91.3** | **78.1** | 44.6 | **96.5** | 49.4 | **94.0** | **98.5** | 89.2 | 77.0 | 49.3 | **62.5** | 79.2 |

## 4.4 Ablation Study

We conducted an ablation study, summarized in Table 4, to analyze the contribution of key components in the ADAM framework. The study evaluates the accuracy under various settings across five values of $k$ (0, 50, 150, 250, and 500). Here, $k = 0$ represents the baseline where the embedding-label repository is not utilized, effectively disabling the similarity search step.

**Prompt Changes.** Excluding key elements from the COP module prompt results in significant performance degradation. Removing the bounding box of the unknown object reduces accuracy by 7.7% at $k = 250$, emphasizing the importance of spatial information. The absence of CLIP-derived characteristics, such as color and texture, results in a similar drop of over 10%. The most severe decline occurs when the known object context is excluded, reducing accuracy from 30.13% to 3.13% at $k = 0$ and from 42.14% to 8.24% at $k = 250$, highlighting its critical role in contextual reasoning.

**Model Structure.** Disabling majority voting and relying solely on the top prediction leads to a sharp performance decline, achieving only 32.23% accuracy at $k = 250$. This shows the importance of frequency-based ranking in aggregating predicted labels. Removing cross-modal re-ranking also reduces accuracy, particularly for larger $k$ values, with a drop from 53.60% to 45.90% at $k = 250$, underscoring the need for CLIP-based refinement of predicted labels. With the embedding-label repository, we see significant improvement from no repository search at $k = 0$ (33.95%) to $k = 250$ (57.64%). After $k = 250$, performance declines to 52.98%, suggesting that larger neighborhoods introduce noise from less relevant embeddings. The vast improvement with the addition of the repository demonstrates its critical role in label assignment. Lastly, the self-refinement of the repository improves performance by 3.66% at $k = 250$, highlighting the benefits of the process.

Table 4: Ablation study for ADAM. The table reports accuracy (%) for different ablation settings across five $k$ values, where $k = 0$ denotes no repository search.

| Method | k Values | | | | |
|---|---|---|---|---|---|
| | $k=0$ | $k=50$ | $k=150$ | $k=250$ | $k=500$ |
| *Prompt Changes* | | | | | |
| No CLIP Characteristics Descriptor | 30.13 | 38.20 | 40.62 | 42.14 | 43.48 |
| No Bounding Box for Unknown Object | 28.45 | 42.88 | 46.32 | 49.91 | 47.76 |
| No Known Object Context | 3.13 | 4.47 | 8.07 | 8.24 | 10.80 |
| *Model Structure* | | | | | |
| No Majority Voting (Top Prediction Only) | – | 28.11 | 30.09 | 32.23 | 34.32 |
| No Cross Modal Reranking | 29.95 | 32.11 | 43.73 | 45.90 | 45.20 |
| With Embedding-Label Repository (ADAM) | 33.95 | 45.95 | 53.21 | **57.64** | 52.98 |
| **+ Self-Refining (Full Model)** | **36.10** | **47.30** | **55.30** | **61.30** | **54.20** |

## 4.5 Limitations

While ADAM achieves high accuracy across diverse object categories, it faces challenges in scenarios where contextual cues are limited. There is also room to test the method on additional datasets, especially those which have more complete labeling of every object in the scene. Additionally, the use of an LLM for label generation, while improving the set of objects that can be classified, comes at the cost of increased computation.

## 5 Conclusion

ADAM introduces a completely novel approach to open-world object detection that operates in a fully zero-shot setting, utilizing an open-world vocabulary without the need for predefined categories or labeled training data for unknown objects. This framework bridges the gap between vision and language by integrating LLMs with visual embedding methods to infer labels based solely on contextual and visual features. This study establishes a foundation for future research in open-world object detection, paving the way for more adaptive and intelligent systems that may reduce the burden of human annotation for computer vision datasets.

## References

[1] Shir Amir, Yossi Gandelsman, Shai Bagon, and Tali Dekel. Deep vit features as dense visual descriptors, 2022. URL https://arxiv.org/abs/2112.05814.

[2] J Anju and R Shreelekshmi. A faster secure content-based image retrieval using clustering for cloud. *Expert Systems with Applications*, 189:116070, 2022.

[3] Riyaz Ahamed Ariyaluran Habeeb, Fariza Nasaruddin, Abdullah Gani, Mohamed Ahzam Amanullah, Ibrahim Abaker Targio Hashem, Ejaz Ahmed, and Muhammad Imran. Clustering-based real-time anomaly detection—a breakthrough in big data technologies. *Transactions on Emerging Telecommunications Technologies*, 33(8):e3647, 2022.

[4] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162, 2018.

[5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020.

[6] Keyan Chen, Xiaolong Jiang, Yao Hu, Xu Tang, Yan Gao, Jianqi Chen, and Weidi Xie. Ovarnet: Towards open-vocabulary object attribute recognition. In *CVPR*, 2023.

[7] Qi Chen, Bing Zhao, Haidong Wang, Mingqin Li, Chuanjie Liu, Zengzhong Li, Mao Yang, and Jingdong Wang. Spann: Highly-efficient billion-scale approximate nearest neighborhood search. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan,

editors, *Advances in Neural Information Processing Systems*, volume 34, pages 5199–5212. Curran Associates, Inc., 2021. URL `https://proceedings.neurips.cc/paper_files/paper/2021/file/299dc35e747eb77177d9cea10a802da2-Paper.pdf`.

[8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.

[9] Tianheng Cheng, Lin Song, Yixiao Ge, Wenyu Liu, Xinggang Wang, and Ying Shan. Yolo-world: Real-time open-vocabulary object detection. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, 2024.

[10] Zhihua Cui, Xuechun Jing, Peng Zhao, Wensheng Zhang, and Jinjun Chen. A new subspace clustering strategy for ai-based data analysis in iot system. *IEEE Internet of Things Journal*, 8 (16):12540–12549, 2021.

[11] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. *arXiv preprint arXiv:2401.08281*, 2024.

[12] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html, 2012.

[13] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *CVPR*, 2023.

[14] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. In *International Conference on Learning Representations*, 2022. URL `https://openreview.net/forum?id=lL3lnMbR4WU`.

[15] Akshita Gupta, Sanath Narayan, KJ Joseph, Salman Khan, Fahad Shahbaz Khan, and Mubarak Shah. Ow-detr: Open-world detection transformer. In *CVPR*, 2022.

[16] Jiawei Han, Jian Pei, and Hanghang Tong. *Data mining: concepts and techniques*. Morgan kaufmann, 2022.

[17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.

[18] Xinyu Huang, Yi-Jie Huang, Youcai Zhang, Weiwei Tian, Rui Feng, Yuejie Zhang, Yanchun Xie, Yaqian Li, and Lei Zhang. Open-set image tagging with multi-grained text supervision. *arXiv e-prints*, pages arXiv–2310, 2023.

[19] Sheng Jin, Xueying Jiang, Jiaxing Huang, Lewei Lu, and Shijian Lu. LLMs meet VLMs: Boost open vocabulary object detection with fine-grained descriptors. In *The Twelfth International Conference on Learning Representations*, 2024. URL `https://openreview.net/forum?id=usrChqw6yK`.

[20] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547, 2019.

[21] K. J. Joseph, Salman H. Khan, Fahad Shahbaz Khan, and Vineeth N. Balasubramanian. Towards open world object detection. *CoRR*, abs/2103.02603, 2021. URL `https://arxiv.org/abs/2103.02603`.

[22] Jinbo Li, Hesam Izakian, Witold Pedrycz, and Iqbal Jamal. Clustering-based anomaly detection in multivariate time series data. *Applied Soft Computing*, 100:106919, 2021.

[23] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, 2022.

[24] Liunian Harold Li*, Pengchuan Zhang*, Haotian Zhang*, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, Kai-Wei Chang, and Jianfeng Gao. Grounded language-image pre-training. In *CVPR*, 2022.

[25] James Liang, Yiming Cui, Qifan Wang, Tong Geng, Wenguan Wang, and Dongfang Liu. Clusterfomer: clustering as a universal visual learner. *Advances in neural information processing systems*, 36, 2024.

[26] Chuang Lin, Yi Jiang, Lizhen Qu, Zehuan Yuan, and Jianfei Cai. Generative region-language pretraining for open-ended object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13958–13968, 2024.

[27] 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 *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.

[28] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023.

[29] Muhammad Maaz, Hanoona Abdul Rasheed, Salman H. Khan, Fahad Shahbaz Khan, Rao Muhammad Anwer, and Ming-Hsuan Yang. Multi-modal transformers excel at class-agnostic object detection. *CoRR*, abs/2111.11430, 2021. URL https://arxiv.org/abs/2111.11430.

[30] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, Xiao Wang, Xiaohua Zhai, Thomas Kipf, and Neil Houlsby. Simple open-vocabulary object detection. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Computer Vision – ECCV 2022*, pages 728–755, Cham, 2022. Springer Nature Switzerland. ISBN 978-3-031-20080-9.

[31] Sudipto Mukherjee, Himanshu Asnani, Eugene Lin, and Sreeram Kannan. Clustergan: Latent space clustering in generative adversarial networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 4610–4617, 2019.

[32] William E. Nagy, Patricia A. Herman, and Richard C. Anderson. Learning words from context. *Reading Research Quarterly*, 20(2):233–253, 1985. ISSN 00340553. URL http://www.jstor.org/stable/747758.

[33] Chuang Niu, Hongming Shan, and Ge Wang. Spice: Semantic pseudo-labeling for image clustering. *IEEE Transactions on Image Processing*, 31:7264–7278, 2022.

[34] Gbeminiyi John Oyewole and George Alex Thopil. Data clustering: application and trends. *Artificial Intelligence Review*, 56(7):6439–6475, 2023.

[35] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *CoRR*, abs/2103.00020, 2021. URL https://arxiv.org/abs/2103.00020.

[36] 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.

[37] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE transactions on pattern analysis and machine intelligence*, 39(6):1137–1149, 2016.

[38] Meitar Ronen, Shahaf E Finder, and Oren Freifeld. Deepdpm: Deep clustering with an unknown number of clusters. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9861–9870, 2022.

[39] Amirreza Rouhi, Diego Patiño, and David K Han. Enhancing object detection by leveraging large language models for contextual knowledge. In *International Conference on Pattern Recognition*, pages 299–314. Springer, 2025.

[40] Amit Saxena, Mukesh Prasad, Akshansh Gupta, Neha Bharill, Om Prakash Patel, Aruna Tiwari, Meng Joo Er, Weiping Ding, and Chin-Teng Lin. A review of clustering techniques and developments. *Neurocomputing*, 267:664–681, 2017.

[41] Erich Schubert. Stop using the elbow criterion for k-means and how to choose the number of clusters instead. *ACM SIGKDD Explorations Newsletter*, 25(1):36–42, 2023.

[42] Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chenfeng Xu, Wei Zhan, Masayoshi Tomizuka, Lei Li, Zehuan Yuan, Changhu Wang, et al. Sparse r-cnn: End-to-end object detection with learnable proposals. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14454–14463, 2021.

[43] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[44] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc Van Gool. Scan: Learning to classify images without labels. In *European conference on computer vision*, pages 268–285. Springer, 2020.

[45] Xiaohua Wang, Zhenghua Wang, Xuan Gao, Feiran Zhang, Yixin Wu, Zhibo Xu, Tianyuan Shi, Zhengyuan Wang, Shizheng Li, Qi Qian, Ruicheng Yin, Changze Lv, Xiaoqing Zheng, and Xuanjing Huang. Searching for best practices in retrieval-augmented generation. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 17716–17736, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.981. URL https://aclanthology.org/2024.emnlp-main.981/.

[46] Xiyue Wang, Yuexi Du, Sen Yang, Jun Zhang, Minghui Wang, Jing Zhang, Wei Yang, Junzhou Huang, and Xiao Han. Retccl: Clustering-guided contrastive learning for whole-slide image retrieval. *Medical image analysis*, 83:102645, 2023.

[47] Stuart Webb. The effects of context on incidental vocabulary learning. *University of Hawaii National Foreign Language Resource Center*, 20(2):232–245, 2008.

[48] Chenyun Wu and Subhransu Maji. How well does clip understand texture? In *ECCV 2022 CVinW Workshop*, 2022. URL https://arxiv.org/abs/2203.11449.

[49] Wen-Bo Xie, Yan-Li Lee, Cong Wang, Duan-Bing Chen, and Tao Zhou. Hierarchical clustering supported by reciprocal nearest neighbors. *Information Sciences*, 527:279–292, 2020.

[50] Yuhang Zang, Wei Li, Kaiyang Zhou, Chen Huang, and Chen Change Loy. Open-vocabulary detr with conditional matching. In *European Conference on Computer Vision*, 2022.

[51] Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14393–14402, June 2021.

[52] Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li, Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin, Tong Luo, Yaqian Li, Shilong Liu, et al. Recognize anything: A strong image tagging model. *arXiv preprint arXiv:2306.03514*, 2023.

[53] Jiyang Zheng, Weihao Li, Jie Hong, Lars Petersson, and Nick Barnes. Towards open-set object detection and discovery. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 3960–3969, 2022. doi: 10.1109/CVPRW56347.2022.00441.

[54] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *ECCV*, 2022.