

# Semantically-Prompted Language Models Improve Visual Descriptions

Anonymous ACL submission

## Abstract

Language-vision models like CLIP have made significant progress in zero-shot vision tasks, such as zero-shot image classification (ZSIC). However, generating specific and expressive visual descriptions remains a challenge as current methods produce descriptions that lack granularity and are ambiguous. To tackle these challenges, we propose V-GLOSS: Visual Glosses, a novel method that prompts language models with semantic knowledge to produce improved visual descriptions. We demonstrate that V-GLOSS can be used to achieve state-of-the-art results on benchmark ZSIC datasets, such as ImageNet and STL-10. In addition, we introduce a silver dataset with visual descriptions generated by V-GLOSS and demonstrate its utility for language-vision tasks.

## 1 Introduction

Language-vision models (Radford et al., 2021; Jia et al., 2021) have made significant progress in zero-shot vision tasks. However, we hypothesize that their accuracy is limited by a lack of visual concept descriptions that are both expressive and specific, that is, glosses that describe what images depicting a concept look like. In this work, we investigate this hypothesis by creating and testing a novel method for producing visual descriptions.

Improving visual descriptions is crucial for enhancing system performance in zero-shot vision tasks. Such descriptions facilitate the creation of more useful representations. Additionally, being able to describe a concept in terms of its appearance is essential for developing more robust and adaptable methods incorporating diverse visual information across various domains, without the need for extensive re-training.

Existing approaches to generating visual descriptions, such as those employed by CLIP (Radford et al., 2021) and CuPL (Pratt et al., 2022), involve directly plugging class labels into fixed templates




Class / Concept	WordNet Gloss	V-GLOSS (Ours)
CORKSCREW 	A bottle opener that pulls corks.	A <b>tool</b> with a <b>spiral blade</b> that is used to remove corks from bottles.
BRAMBLING 	Eurasian finch.	A <b>small brown</b> bird with a <b>black head</b> and a <b>white patch</b> on its chest.
BROCCOLI 	Branched green undeveloped flower heads.	A <b>green vegetable</b> with a <b>thick stalk</b> and <b>florets</b> that grow in a <b>dense head</b> .

Table 1: A qualitative comparison between WordNet concept glosses and V-GLOSS (Silver) class descriptions for some ImageNet classes. Our method describes what a class *looks like*, instead of what it *does* or *is*.

(e.g., *a photo of X*), or using large language models such as InstructGPT (Ouyang et al., 2022) to generate descriptions based on class labels (e.g., *what does X look like?*). These methods suffer from two main issues: class granularity and label ambiguity. Class granularity refers to the difficulty in distinguishing between visually similar classes, such as ALLIGATOR and CROCODILE. Label ambiguity is caused by using polysemous words as labels for distinct concepts. For example, CRANE can refer to either a bird or a construction machine. These issues limit the performance of existing models (Radford et al., 2021).

To address these challenges, we introduce V-GLOSS, a novel method that leverages language models (LMs) and semantic knowledge bases (SKBs) to generate improved visual descrip-

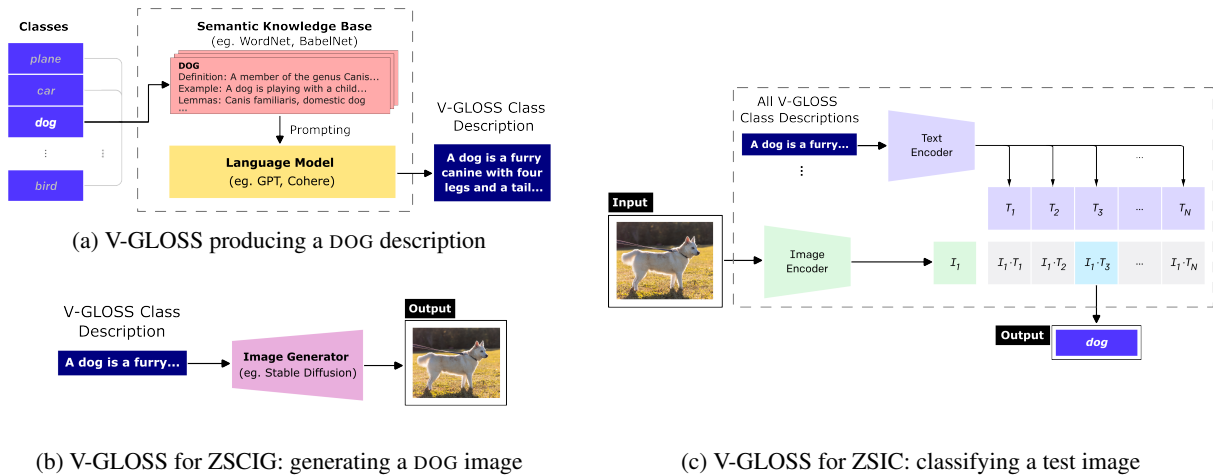


Figure 1: For the DOG class, we depict (a) V-GLOSS’s architecture (Section 4.2.1), along with adaptations: (b) ZSIC (Section 5.4.1) and (c) ZSCIG (Section 5.4.1)

tions – visual glosses. Table 1 shows some examples. By combining structured semantic information from SKBs such as WordNet (Miller, 1998), with a contrastive algorithm to distinguish similar classes, V-GLOSS is designed to mitigate the dual issues of granularity and ambiguity.

Our results demonstrate the effectiveness of V-GLOSS in improving the performance of ZSIC systems. We achieve state-of-the-art (SOTA) results on benchmark datasets such as ImageNet (Deng et al., 2009), CIFAR-10, and CIFAR-100 (Krizhevsky et al., 2009) in the zero-shot setting, and STL-10 (Coates et al., 2011) in both the zero-shot and supervised settings. Additionally, we introduce V-GLOSS Silver, a silver dataset constructed by V-GLOSS, consisting of a visual gloss for each ImageNet class. We show that V-GLOSS Silver is useful for language-vision tasks such as ZSIC and ZSCIG, comparing favorably to WordNet glosses.

## 2 Tasks

Our main task is to generate a description for a given class or concept. For example, if an image classification dataset has the class DOG, we aim to produce a description such as “A dog is a furry, four-legged canine...” We consider such a description to be a specific kind of gloss.

We use two downstream tasks to compare methods of generating class descriptions: zero-shot image classification (ZSIC), and zero-shot class-conditional image generation (ZSCIG).

In ZSIC, the goal is to classify an image based on a set of classes, without having seen any la-

beled images belonging to those classes. The set of classes depends on the dataset. For example, given an image depicting a dog, we aim to predict the class DOG.

In ZSCIG, the goal is to generate an image that corresponds to a specific class, again without having seen any labeled examples. For example, given a class DOG, we aim to generate an image of a dog.

In short, ZSIC is the task of classifying a given image, while ZSCIG is the task of generating an image given a class. Both involve classes and images. Visual descriptions of classes provide useful information which can facilitate both tasks, by making it easier to either recognize or generate images of each class. Therefore, by developing a novel method to improve the generation of such descriptions, we hypothesize that performance on ZSIC and ZSCIG can be improved.

## 3 Related Work

**Language Models** The advent of transformer-based language models has revolutionized many natural language processing tasks (Radford et al., 2018; Devlin et al., 2018; Radford et al., 2019; Brown et al., 2020; Black et al., 2022; Ouyang et al., 2022). As these models are scaled up by their number of parameters and quantity of training data, they exhibit emergent abilities such as few-shot and zero-shot learning (Wei et al., 2022).

**Language-Vision Models** Significant strides have been made in the field of language-vision models such as CLIP (Radford et al., 2021) and ALIGN (Jia et al., 2021). These models apply contrastive pre-training approaches on large image-text

124 datasets, leading to improved representation learn-  
 125 ing for both text and images and enhanced perfor-  
 126 mance on several multi-modal tasks (Mokady et al.,  
 127 2021; Song et al., 2022). Further advancements  
 128 have been achieved by scaling up pre-training and  
 129 incorporating auxiliary training objectives (Pham  
 130 et al., 2021; Yu et al., 2022).

131 **Producing Descriptions & Prompting** The gen-  
 132 eration of descriptions and prompting has been  
 133 explored in various studies. Radford et al. (2021)  
 134 introduced the template ensemble (TE) method,  
 135 which uses a custom set of class labels and a fixed  
 136 set of templates. Each label is inserted into these  
 137 templates, and the completed templates for each  
 138 class are aggregated into a single representation of  
 139 the class. The CuPL method (Pratt et al., 2022)  
 140 utilizes InstructGPT (Brown et al., 2020; Ouyang  
 141 et al., 2022) to generate descriptions for ImageNet  
 142 classes. Both TE and CuPL can be used for zero-  
 143 shot image classification. Hao et al. (2022) fine-  
 144 tuned GPT models (Radford et al., 2018, 2019)  
 145 to rephrase image-generation prompts, resulting  
 146 in improved images. (Zhou et al., 2022) learned  
 147 soft prompts that improve performance, but are  
 148 intractable to humans. In this work, we prompt lan-  
 149 guage models with semantic knowledge to generate  
 150 visual descriptions.

## 151 4 Method

152 We begin by describing how we map classes to  
 153 concepts in a semantic knowledge base (SKB), in  
 154 order to leverage the concept-specific information  
 155 the SKB contains. We then introduce our novel  
 156 method V-GLOSS, which has two variants, *nor-*  
 157 *mal* and *contrastive*. We conclude by describing  
 158 the construction of V-GLOSS Silver, a set of class  
 159 descriptions produced using V-GLOSS.

### 160 4.1 Mapping Classes to WordNet Synsets

161 The ImageNet classes are already mapped to Word-  
 162 Net synsets by the dataset’s creators. For the  
 163 other datasets, we employ a heuristic that starts  
 164 by mapping each class to the most frequent sense  
 165 of the class label, as determined by WordNet<sup>1</sup>. For  
 166 CIFAR-10 and STL-10, this heuristic is sufficient.  
 167 For CIFAR-100, we manually re-map 18 classes.  
 168 For instance, we needed to re-map RAY from *light*  
 169 to *sea creature*, as the *light* sense is the most fre-  
 170 quent according to WordNet, but the RAY images  
 171 in the dataset depict sea creatures.

<sup>1</sup><https://www.nltk.org/>

What does a **platypus** look like?

A platypus looks like a beaver with a duck's bill

(a) CuPL Pratt et al. (2022)

...  
 Concept name: eagle  
 Hypernyms: bird or prey  
 Hyponyms: bald eagle, eaglet, golden eagle, harpy  
 Gloss: any of various large keen-sighted diurnal birds  
 of prey noted for their broad wings and strong...  
 Unique and expressive visual description: Eagles are  
 large birds of prey with dark brown bodies and wings...  
 ...  
 Concept name: **platypus**  
 Hypernyms: **duckbill, duckbilled platypus, ...**  
 Hyponyms: **egg-laying mammal**  
 Gloss: **small densely furred aquatic monotreme of**  
**Australia and Tasmania having a broad bill...**  
 Unique and expressive visual description:  
 Platypuses are water-dwelling mammals that have  
 broad duck-like bills and hind legs with a foot web that  
 has an intricate web of keratinised spongy hairs

(b) V-GLOSS

Figure 2: Class descriptions for PLATYPUS generated by two different methods that use LMs. Input prompts, output descriptions, and **plugged values** are shown.

## 172 4.2 V-GLOSS

173 We discuss the two variants of V-GLOSS below,  
 174 *normal* and *contrastive*. In both, for each class,  
 175 we produce multiple descriptions resulting in an  
 176 ensemble.

### 177 4.2.1 Normal V-GLOSS

178 We generate normal descriptions via in-context  
 179 learning with an LM, beginning by providing the  
 180 LM with a description of the task to be performed,  
 181 followed by multiple input-output examples. The  
 182 examples are fixed, involving the concepts EA-  
 183 GLE, BAT (animal), BAT (baseball), and TELEVI-  
 184 SION. We selected these to expose the model to am-  
 185 biguous class labels (*bat*), a natural object (*eagle*),  
 186 and an artificial object via (*television*). For each  
 187 class, we obtain from WordNet the hypernyms, hy-  
 188 ponyms, usage examples, synonyms, and glosses  
 189 of the sense to which the class is mapped, and pro-  
 190 vide this to the LM. Figure 2b shows a session with  
 191 the LM, beginning with the example of *eagle*, with  
 192 output generated for the class *platypus*. Table 1  
 193 compares our descriptions to WordNet glosses.

### 194 4.2.2 Contrastive V-GLOSS

195 During development, we observed that many er-  
 196 rors were caused by false positives involving vi-  
 197 sually similar classes. For example, the classes

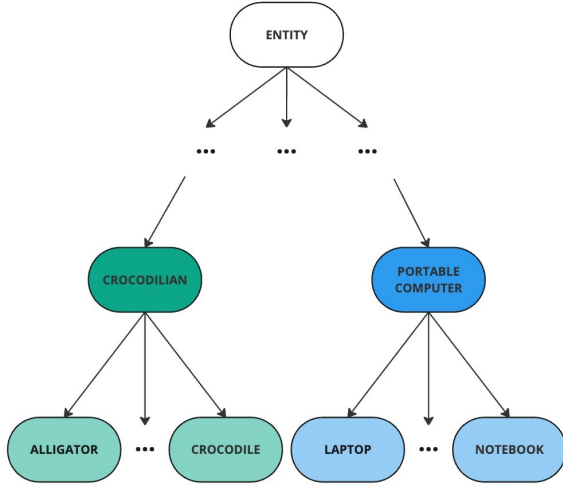


Figure 3: A sample of WordNet hypernym hierarchy. For *contrastive* prompting, we only distinguish classes that are semantically similar to the target class, like **ALLIGATOR** to **CROCODYLE**.

CROCODYLE for ALLIGATOR refer to similar-looking animals, and are often confused for one another. The contrastive variant of V-GLOSS is designed to address this by using semantic similarity between classes as a heuristic to estimate visual similarity. For each class, we search for other classes that are semantically similar, and if any are found, we add a negative instruction to the LM prompt, e.g. we generate a description for an ALLIGATOR *but not* a CROCODYLE, using the same prompt structure as for normal V-GLOSS.

We create a similarity matrix  $M$  as follows:

$$M_{i,j} = Sim(S[i], S[j]) \quad (1)$$

$Sim(s_1, s_2)$  is the Wu-Palmer path-similarity function (Wu and Palmer, 1994) comparing synsets  $s_1$  and  $s_2$ ; this similarity function uses the path between two concepts in the WordNet hypernym hierarchy (Figure 3) to measure semantic relatedness.  $S$  is the set of all classes in a dataset,  $\mathcal{D}$ , and  $i$  and  $j$  are indices ranging from 1 to  $|S|$ . Concisely, Equation 1 defines a similarity matrix containing similarity scores between all classes in a dataset.  $M$  is one of the inputs to our contrastive V-GLOSS variant, shown in Algorithm 1.

In Algorithm 1,  $\lambda$  is a threshold for minimum similarity. We only generate contrastive descriptions when classes have a similarity that exceeds or is equal to  $\lambda$ .  $N$  indicates the maximum number of classes to generate contrastive descriptions for.  $k$  is the number of distinct descriptions to gener-



Class / Concept	Normal	Contrastive
ALLIGATOR	 A large reptile with a long snout, a broad head, and a long tail.	A large, <b>dark-colored</b> reptile with a <b>rounded snout</b> , found in <b>freshwater</b> .
CROCODYLE	 A reptile with a broad, flat snout, a long tail, and a long, pointed snout.	A <b>grayish-green</b> reptile with a <b>v-shaped snout</b> , found in <b>brackish</b> or <b>saltwater</b> .

Table 2: Two similar classes with **key differences** between their *normal* and *contrastive* descriptions.

ate for a class pair.  $LM_c$  takes in the *target* class, a neighbor class, and  $k$ , then prompts the LM to generate  $k$  descriptions that distinguish the *target* and neighbor classes. In summary, for each class, Algorithm 1 identifies the classes most similar to it, excluding itself, and generates descriptions that distinguish them. Table 2 compares the normal and contrastive descriptions for ALLIGATOR and CROCODYLE; note that distinguishing features of the two classes are included in the LM’s output. Table 3 shows examples of classes with high false positive rates, and the classes they are contrasted with.

**Algorithm 1** Generate Contrastive Descriptions: We generate contrastive descriptions to help distinguish the most similar classes.

**Require:**  $M$ : Equation 1 result

**Require:**  $\lambda, N, k$ : Hyperparameters

**Require:**  $S$ : All classes in dataset,  $\mathcal{D}$

**Require:**  $LM_c$ : LM prompted contrastively

- 1:  $G \leftarrow$  empty  $|S|$ -list for class descriptions
- 2: **for**  $i \leftarrow 0$  to  $|S| - 1$  **do**
- 3:    $target \leftarrow S[i]$
- 4:    $S^* \leftarrow$  top  $N$  classes :  $\lambda \leq M_{i,*} \leq 1$
- 5:   **for**  $s^*$  in  $S^*$  **do**
- 6:      $samples \leftarrow LM_c(target, s^*, k)$
- 7:      $G[i].insert(samples)$
- 8: **return**  $G$

## 5 Evaluation

Toward evaluating V-GLOSS, we describe our datasets, evaluation metrics, baselines, previous methods, and experiments.

Class	False Positives	Contrastives
AFRICAN ELEPHANT	TUSKER (44), ASIAN ELEPHANT (6)	TUSKER, ASIAN ELEPHANT
NOTEBOOK	LAPTOP (22), DESKTOP (10), SPACE BAR (2)	LAPTOP, DESKTOP, SPACE BAR

Table 3: False positives and their counts vs. classes selected by the contrastive algorithm (see Equation 1 and Algorithm 1). Hits and misses are shown.

## 5.1 Datasets

We evaluate our method on the test splits of four widely used benchmark datasets: ImageNet (Deng et al., 2009) consists of 50,000 images equally distributed across 1,000 classes, and serves as our *primary* benchmark. CIFAR-10 and CIFAR-100 (Krizhevsky et al., 2009) both comprise 10,000 test samples across 10 and 100 classes, respectively. Finally, STL-10 (Coates et al., 2011) comprises 100,000 test samples and is designed for unsupervised learning. For CIFAR-10, CIFAR-100, and STL-10, which are not pre-mapped to WordNet, we employ the two-step process detailed in Section 4.1 to map each class to a WordNet synset.

Experiment 1 (Section 5.4) involves ImageNet alone and covers both the ZSCIG and ZSIC tasks. In contrast, Experiment 2 (Section 5.5), which is our main experiment, tests the impact of various class description methods on the ZSIC task and uses all datasets. In Experiment 2, we allow methods to use ensembles of descriptions of each class, while in Experiment 1, we experiment with only a single description.

The datasets we selected to evaluate the following properties of V-GLOSS:

- Performance on benchmark datasets with varying numbers of classes.** Each dataset has its own set of classes, ranging from ImageNet with 1,000 classes, to CIFAR-100 with 100 classes, to CIFAR-10 and STL-10, each with 10 classes.
- Ability to represent diverse concepts at varying levels of granularity.** The datasets we use contain a wide range of concepts across various domains, rather than those targeting specific subareas such as pets (Parkhi et al., 2012), foods (Bossard et al., 2014), cars (Krause et al., 2013), scenes (Xiao et al., 2010), or airplanes (Maji et al., 2013).

## 5.2 Evaluation Metrics

**Top-1 Accuracy** In ZSIC, this metric is the frequency with which the model’s top prediction for an image matches the gold label.

**Fréchet Inception Distance (FID)** For ZSCIG, FID (Heusel et al., 2017) quantifies the divergence between ground truth and generated images, with lower scores signifying a better ability to produce images similar to the ground truth.

**Inception Score** Also for ZSCIG, the inception score (Salimans et al., 2016) uses an Inception model’s (Szegedy et al., 2015) output probability distribution to assess the diversity and realism of generated images, with higher scores indicating more diverse and convincing images. Unlike the above metrics, this does not require ground-truth images for comparison.

## 5.3 Baseline & Previous Methods

In this section, we describe the methods to which we compare V-GLOSS. For methods that produce ensembles of class descriptions (i.e. multiple descriptions per class), a single representation of the class is obtained by averaging individual representations for each description.

First, the **1-Template baseline** inserts a class label into a *single* specific template. For example, given the class DOG, the baseline produces “A *photo of a dog.*”

**Template Ensemble** (Radford et al., 2021) generates an ensemble of descriptions for a class by inserting the class label into each of a set of templates. For example, some descriptions for DOG are: “A *photo of a dog.*”, “A *blurry photo of a dog.*”, and “An *origami dog.*” This method uses a modified list of class labels<sup>2</sup> designed to reduce ambiguity.

**CuPL** (Pratt et al., 2022) also generates an ensemble of descriptions for each class. The descriptions are generated by prompting a LLM, Instruct-GPT (Ouyang et al., 2022), with questions such as: “What does a dog look like?” and “Describe an image of a dog from the internet.” CuPL uses the same class labels as Template Ensemble.

The authors of CuPL also combined their method with Template Ensemble. The resulting method, **CuPL + Template Ensemble**, combines the class descriptions from both methods.

<sup>2</sup><https://github.com/anishathalye/imagenet-simple-labels>

## 5.4 Experiment 1: V-GLOSS Silver

This experiment evaluates V-GLOSS’s ability to generate a *single* description for each class, without relying on ensembling. We then evaluate the V-GLOSS description of each class against its WordNet gloss.

To construct this set of class descriptions, which we view as a silver dataset of such descriptions, we generate a *single, normal* description for each ImageNet class via greedy decoding. We generate only *normal* descriptions because they outperform *contrastive* ones when only a single description is used. We call the resulting dataset *V-GLOSS Silver*.

We extrinsically evaluate V-GLOSS Silver by using it for the ZSIC and ZSCIG tasks, and comparing the results to those achieved using the 1-Template baseline, and WordNet glosses. We do not compare V-GLOSS Silver to CuPL or other previous methods which do not produce a single description for each class.

### 5.4.1 Technical Details

**ZSIC** We employ CLIP (Radford et al., 2021), which comprises an image encoder and a text encoder, as the ZSIC backbone model. Our procedure consists of three steps: First, we use the CLIP text encoder to create an aggregate representation for each class based on its description(s). Then, at test time, we employ the CLIP image encoder to generate a representation of the input image. Finally, we predict the class which maximizes the cosine similarity between the representation of its description(s), and the image representation (see Figure 1c). We evaluate the predictions using top-1 accuracy.

**ZSCIG** For ZSCIG (see Figure 1b), we condition Stable Diffusion (Rombach et al., 2022) on each class description before generating an image. We use a guidance scale of 7.5 and run 50 diffusion steps. We evaluate the generated images using Inception and FID scores.

### 5.4.2 Results

The results of Experiment 1 are shown in Table 4. Based on our extrinsic evaluation in the ZSIC and ZSCIG tasks, *V-GLOSS Silver* descriptions yield better performance compared to baseline and WordNet Glosses. On ZSIC, we improve accuracy by 1.3%; on ZSCIG, we improve Inception and FID scores by 9.9 and 5.7, respectively. This demonstrates the effectiveness and utility of V-GLOSS:

	ZSIC		ZSCIG	
	Accuracy ↑	Inception ↑	FID ↓	
Baseline (1-Template)	71.0	99.7	25.7	
WordNet Glosses	44.7	58.5	30.0	
V-GLOSS Silver	<b>72.3</b>	<b>109.6</b>	<b>20.0</b>	

Table 4: Extrinsic evaluation on the tasks of ZSIC and ZSCIG. ↓ means that lower is better.

our visual descriptions yield better results on ZSIC and ZSCIG.

### 5.4.3 Analysis

V-GLOSS Silver descriptions are considerably more detailed, more expressive, and better visually grounded than their WordNet gloss counterparts (see Figure 1). Specifically, we observe that V-GLOSS descriptions make greater use of descriptive words and phrases, e.g. *spiral, brown, green, thick, small*, etc.

## 5.5 Experiment 2: ZSIC

Our second experiment assesses the effectiveness of V-GLOSS descriptions in facilitating ZSIC. The details for the ZSIC pipeline are largely similar to those described in Experiment 1 (Section 5.4), except that we generate an ensemble of descriptions per class, as opposed to only one description. We also experiment with two image encoder variants: ViT (Dosovitskiy et al., 2020) and RN50 (He et al., 2016). For all baselines and methods (Section 5.3, Section 4.2.1), we follow the same evaluation procedure after generating class descriptions.

### 5.5.1 Technical Details

We generate class descriptions using the 6.1B-parameter Cohere LM<sup>3</sup>. We choose Cohere over alternatives due to its extensive cost-free availability, reducing the cost of our experiments. Cohere has comparable performance to the similarly-sized InstructGPT (Brown et al., 2020; Ouyang et al., 2022) variant, as demonstrated by Liang et al. (2022) across various benchmarks. Therefore, we do not gain any advantage by using Cohere instead of InstructGPT.

When generating class descriptions with *normal* V-GLOSS, we use a temperature of 2.5 to produce an ensemble of 50 descriptions per class. When generating *contrastively*, we use a temperature of 1.5 to generate an ensemble of 20 descrip-

<sup>3</sup><https://docs.cohere.com/docs/models>

Method	Model	Accuracy (%) on Datasets				# LM Parameters
		ImageNet	CIFAR-100	CIFAR-10	STL-10	
Baseline (1-Template)	ViT	72.4	77.3	95.2	99.5	0
	RN50	68.7	57.7	81.0	98.4	
Template Ensemble	ViT	76.2	77.9	96.2	99.4	0
	RN50	73.2	61.3	86.8	98.3	
CuPL	ViT	76.7	-	-	-	175B
CuPL + Template Ensemble	ViT	77.6	-	-	-	175B
	RN50	75.1	-	-	-	
V-GLOSS ( <i>Normal-Only</i> )	ViT	77.3	77.5	95.6	99.4	6.1B
	RN50	73.3	63.5	86.8	98.3	
V-GLOSS ( <i>Normal + Contrastive</i> )	ViT	<b>78.5</b>	<b>78.2</b>	<b>97.0</b>	<b>99.6</b>	6.1B
	RN50	74.5	64.6	87.8	98.8	

Table 5: Top-1 accuracy on ZSIC. ViT and RN are Transformer- and ResNet-based CLIP variants.

tions per class. Like Pratt et al. (2022), we observe that performance saturates around 50 descriptions for *normal* V-GLOSS, but we also observe saturation at around 20 descriptions for *contrastive* V-GLOSS. Based on tuning on development data, we set  $N = 5$ ,  $\lambda = 0.5$ , and  $k = 4$  (see Algorithm 1). In total, we obtain 70 class descriptions. During generation, we set the maximum number of tokens to 35, but also terminate generation when the *boundary parameter* or *newline* token is reached.

### 5.5.2 Results

The results of Experiment 2 are shown in Table 5. V-GLOSS yields better accuracy than the baseline by an average of 3.60% overall (2.22% with ViT and 4.98% with RN50). V-GLOSS also outperforms Template Ensemble and CuPL + Template Ensemble, by 1.21% and 0.15% respectively. This improvement is especially notable since the top 15 results on the ImageNet benchmark differ by less than 1% accuracy.<sup>4</sup>

In addition, we make the following observations. (1) V-GLOSS (*Normal + Contrastive*) surpasses V-GLOSS (*Normal-Only*), by an average of 0.91% accuracy. (2) We outperform *CuPL + Template Ensemble* using an LLM with 28.7x fewer parameters. (3) The RN backbone (He et al., 2016), which is generally less capable than ViT (Dosovitskiy et al., 2020), sees a more significant benefit from the V-GLOSS method, on average 3.8%. (4) For STL-10, V-GLOSS matches the top-performing supervised system (Gesmundo, 2022) with a score of 99.6%.

We also note that the *contrastive* component is more helpful on the larger datasets: CIFAR-100

<sup>4</sup><https://paperswithcode.com/sota/image-classification-on-imagenet>

and ImageNet, which have more opportunities for mutual ambiguity between different classes, than on the smaller ones: CIFAR-10 and STL-10. Concretely, this improvement is 1.05%, on average. Later, in Section 6, we discuss these results and their implications more extensively.

### 5.5.3 Analysis

In Section 1, we pointed out several problems in previous methods. Here, we carefully analyze how our V-GLOSS method addresses these issues.

**Label Ambiguity:** Without adequate context, text models may fail to grasp the intended meaning of a polysemous word. *Crane* is a polysemous word, and ImageNet (Deng et al., 2009) has two classes that refer to different senses of the word: *construction machine* and *wading bird*, but use the same label. Thus, in *1-Template*, for example, both classes have the same description. This point highlights an important benefit of linking classes to WordNet, which resolves such ambiguity. Empirically, when compared with a ViT backbone to the *Lex Baseline*, our accuracy on **CRANE** (machine) and **CRANE** (bird) increase from 0% and 46% to 76% and 78%, respectively.

### Relationship Between Performance & Context:

When comparing the baselines to the other methods, we observe that accuracy generally improves as the amount of surrounding context increases. On one hand, if a sentence consists of “*my crane.*” alone, the sense of *crane* is unclear. On the other, if the sentence is “*my construction crane,*” the meaning of *crane* becomes clear. We see that providing additional context helps to disambiguate words. When a description provides more useful context,



Figure 4: V-GLOSS Attention Map



Figure 5: WordNet Gloss Attention Map

models can form better representations of specific classes. By comparing V-GLOSS to the baselines (see Table 5), we can observe that the benefits of additional context extend to the vision-language setting. Concretely, providing visually-grounded context in the description improves performance.

**Class Granularity:** We consider pairs of classes that are similar enough to be mistaken, such as ALLIGATOR and CROCODILE. In WordNet, relationships between synsets are modeled through *is-a* (hyponymy-hypernymy) and *part-of* (meronymy-holonymy) relationships. For example, CROCODILIAN is a hypernym of both ALLIGATOR and CROCODILE, while only ALLIGATOR is a holonym of SNOUT, since alligators have snouts while crocodiles do not. Using our contrastive algorithm, we generate descriptions that highlight how images of a CROCODILE should depict a greener animal with a rounded snout. Empirically, using ViT, the average accuracy of V-GLOSS across these two classes jumps from 36% to 68% when contrastive glosses are used. This improvement highlights the effectiveness of our contrastive V-GLOSS variant in reducing false positives between visually similar classes.

**Attention To Relevant Context:** We analyze the model’s attention maps to better understand V-GLOSS’s impact. Figure 4 shows the attention map for V-GLOSS (see Table 1 for descriptions), indicating effective utilization of visually-relevant context. Conversely, Figure 5 shows the attention map for the WordNet glosses (baseline), where the attention score on *bottle* is 3.5x higher, implying less distraction in V-GLOSS. These maps demonstrate success in steering the model’s attention toward relevant context, thus improving classification accuracy across different classes and descriptions.

## 6 Discussion

When looking at our results, a pertinent question arises: Why does an SKB, such as WordNet, help us do better on tasks related to vision? In this section, we formulate two insights on how the synergy between SKBs and LMs supports our improvements.

### Insight #1: SKBs represent concepts precisely

When LMs are prompted with better information, they produce better output (Borgeaud et al., 2022). WordNet provides a precise representation of a class and its relationship to other classes, leaving minimal room for ambiguity. Afterward, we can prompt an LM with this precise information to produce unambiguous and high-quality class descriptions.

### Insight #2: Semantic similarity is a useful proxy for visual similarity

WordNet models lexical semantics as a graph (see Figure 3), with synsets as nodes and *is-a* relationships as directed edges. The distance between different nodes reflects the level of semantic similarity and is by extension an indicator of the level of visual similarity between synsets. ALLIGATOR and CROCODILE are semantically similar because they are both kinds of CROCODILIAN, but they are visually similar as well (see Table 2). Semantic similarity informs what classes we distinguish with our contrastive descriptions and why they work (see Table 3). This is because semantic and visual similarity are highly correlated.

## 7 Conclusion

This study concentrates on generating visual class descriptions for ZSIC and ZSCIG tasks. We utilize a unique method that merges Semantic Knowledge Bases (SKBs) and Language Models (LMs) to create high-quality descriptions. Our findings reveal that the semantic information from SKBs can condition an LM to generate accurate, expressive, and visually grounded descriptions. Furthermore, we observe that LMs, although pre-trained solely on text, contain latent knowledge about the visual properties of concepts. This knowledge can be harnessed using our novel V-GLOSS method, thus improving the accuracy of zero-shot image classification and generation models. This underscores the strong relationship between language and vision, suggesting potential for LMs in future multi-modal tasks.



## 570 Limitations

571 **The dataset must be mapped to an SKB.** As  
572 described earlier, mapping the dataset to WordNet,  
573 although a one-time step, is not fully automatic. In  
574 future work, we look to fully automate this step,  
575 possibly by selecting a synset based on the simi-  
576 larity between sample class images and potential  
577 senses of the class label.

578 **We are limited in terms of language, dataset**  
579 **class count, and our SKB’s size.** First, our  
580 English-focused stance may prove a limiting factor  
581 in our method being applied to ZSIC or ZSCIG  
582 tasks based in other languages. Some classes are  
583 strongly related to non-English languages.

584 Second, our largest evaluation dataset, ImageNet  
585 (Deng et al., 2009), has 1,000 classes, representing  
586 just 0.64% coverage of WordNet. We look forward  
587 to evaluating our methods on a larger ImageNet  
588 set: ImageNet-21k, which would cover 14.06% of  
589 WordNet.

590 Third, although our method can be applied to  
591 BabelNet (Navigli and Ponzetto, 2012), which has  
592 over 1.5 billion synsets, we focus on WordNet,  
593 which has 155,287. We look to explore alternative  
594 SKBs such as BabelNet, or non-English wordnets,  
595 both of which offer the benefit of being multilin-  
596 gual.

## 597 Ethics Statement

598 In normal use, we discover no direct ethical issues  
599 with our method. Note, however, that we may  
600 inherit ethical problems from the components used  
601 by our method. Both CLIP (Agarwal et al., 2021)  
602 and LMs (Liang et al., 2021) have independently  
603 been shown to exhibit some level of bias. Also,  
604 semantic resources such as WordNet (Miller, 1995)  
605 tend to focus on formalized concepts. This poses  
606 a problem if our method’s use concerns people on  
607 the fringes of society.

608 We noted earlier that our method is mostly  
609 English-focused. This could be a source of bias  
610 if our method is applied in a multilingual context.  
611 We ask that people do not apply our method to real-  
612 world problems where multilingual knowledge is  
613 required. There is also the issue of semantic re-  
614 sources for low-resource languages not being ex-  
615 tensive enough (Magueresse et al., 2020).

## References

- Sandhini Agarwal, Gretchen Krueger, Jack Clark, Alec Radford, Jong Wook Kim, and Miles Brundage. 2021. Evaluating clip: towards characterization of broader capabilities and downstream implications. *arXiv preprint arXiv:2108.02818*. 617-621
- Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, et al. 2022. Gpt-neox-20b: An open-source autoregressive language model. *arXiv preprint arXiv:2204.06745*. 622-626
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*, pages 2206–2240. PMLR. 627-633
- Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. 2014. Food-101—mining discriminative components with random forests. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part VI 13*, pages 446–461. Springer. 634-639
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901. 640-645
- Adam Coates, Andrew Ng, and Honglak Lee. 2011. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223. JMLR Workshop and Conference Proceedings. 646-651
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee. 652-656
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*. 657-660
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*. 661-667
- Andrea Gesmundo. 2022. A continual development methodology for large-scale multitask dynamic ml systems. *arXiv preprint arXiv:2209.07326*. 668-670

671	Yaru Hao, Zewen Chi, Li Dong, and Furu Wei. 2022.	Roberto Navigli and Simone Paolo Ponzetto. 2012. Ba-	722
672	Optimizing prompts for text-to-image generation.	belnet: The automatic construction, evaluation and	723
673	<i>arXiv preprint arXiv:2212.09611</i> .	application of a wide-coverage multilingual semantic	724
674	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian	network. <i>Artificial intelligence</i> , 193:217–250.	725
675	Sun. 2016. Deep residual learning for image recog-		
676	nition. In <i>Proceedings of the IEEE conference on</i>	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Car-	726
677	<i>computer vision and pattern recognition</i> , pages 770–	roll L Wainwright, Pamela Mishkin, Chong Zhang,	727
678	778.	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	728
679	Martin Heusel, Hubert Ramsauer, Thomas Unterthiner,	2022. Training language models to follow in-	729
680	Bernhard Nessler, and Sepp Hochreiter. 2017. Gans	structions with human feedback. <i>arXiv preprint</i>	730
681	trained by a two time-scale update rule converge to a	<i>arXiv:2203.02155</i> .	731
682	local nash equilibrium. <i>Advances in neural informa-</i>	Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman,	732
683	<i>tion processing systems</i> , 30.	and CV Jawahar. 2012. Cats and dogs. In <i>2012</i>	733
684	Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana	<i>IEEE conference on computer vision and pattern</i>	734
685	Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen	<i>recognition</i> , pages 3498–3505. IEEE.	735
686	Li, and Tom Duerig. 2021. Scaling up visual and	Hieu Pham, Zihang Dai, Golnaz Ghiasi, Kenji	736
687	vision-language representation learning with noisy	Kawaguchi, Hanxiao Liu, Adams Wei Yu, Jiahui	737
688	text supervision. In <i>International Conference on</i>	Yu, Yi-Ting Chen, Minh-Thang Luong, Yonghui	738
689	<i>Machine Learning</i> , pages 4904–4916. PMLR.	Wu, et al. 2021. Combined scaling for open-	739
690	Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-	vocabulary image classification. <i>arXiv preprint</i>	740
691	Fei. 2013. 3d object representations for fine-grained	<i>arXiv: 2111.10050</i> .	741
692	categorization. In <i>Proceedings of the IEEE inter-</i>	Sarah Pratt, Rosanne Liu, and Ali Farhadi. 2022. What	742
693	<i>national conference on computer vision workshops</i> ,	does a platypus look like? generating customized	743
694	pages 554–561.	prompts for zero-shot image classification. <i>arXiv</i>	744
695	Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learn-	<i>preprint arXiv:2209.03320</i> .	745
696	ing multiple layers of features from tiny images.	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	746
697	Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	747
698	Ruslan Salakhutdinov. 2021. Towards understand-	try, Amanda Askell, Pamela Mishkin, Jack Clark,	748
699	ing and mitigating social biases in language models.	Gretchen Krueger, and Ilya Sutskever. 2021. <b>Learn-</b>	749
700	In <i>International Conference on Machine Learning</i> ,	<b>ing transferable visual models from natural language</b>	750
701	pages 6565–6576. PMLR.	<b>supervision</b> .	751
702	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris	Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya	752
703	Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian	Sutskever, et al. 2018. Improving language under-	753
704	Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar,	standing by generative pre-training.	754
705	et al. 2022. Holistic evaluation of language	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	755
706	models. <i>arXiv preprint arXiv:2211.09110</i> .	Dario Amodei, Ilya Sutskever, et al. 2019. Language	756
707	Alexandre Magueresse, Vincent Carles, and Evan Heet-	models are unsupervised multitask learners. <i>OpenAI</i>	757
708	derks. 2020. Low-resource languages: A review	<i>blog</i> , 1(8):9.	758
709	of past work and future challenges. <i>arXiv preprint</i>	Robin Rombach, Andreas Blattmann, Dominik Lorenz,	759
710	<i>arXiv:2006.07264</i> .	Patrick Esser, and Björn Ommer. 2022. High-	760
711	Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew	resolution image synthesis with latent diffusion mod-	761
712	Blaschko, and Andrea Vedaldi. 2013. Fine-grained	els. In <i>Proceedings of the IEEE/CVF Conference</i>	762
713	visual classification of aircraft. <i>arXiv preprint</i>	<i>on Computer Vision and Pattern Recognition</i> , pages	763
714	<i>arXiv:1306.5151</i> .	10684–10695.	764
715	George A Miller. 1995. Wordnet: a lexical database for	Tim Salimans, Ian Goodfellow, Wojciech Zaremba,	765
716	english. <i>Communications of the ACM</i> , 38(11):39–41.	Vicki Cheung, Alec Radford, and Xi Chen. 2016.	766
717	George A Miller. 1998. <i>WordNet: An electronic lexical</i>	Improved techniques for training gans. <i>Advances in</i>	767
718	<i>database</i> . MIT press.	<i>neural information processing systems</i> , 29.	768
719	Ron Mokady, Amir Hertz, and Amit H Bermano. 2021.	Haoyu Song, Li Dong, Wei-Nan Zhang, Ting Liu, and	769
720	Clipcap: Clip prefix for image captioning. <i>arXiv</i>	Furu Wei. 2022. Clip models are few-shot learners:	770
721	<i>preprint arXiv:2111.09734</i> .	Empirical studies on vqa and visual entailment. <i>arXiv</i>	771
		<i>preprint arXiv:2203.07190</i> .	772
		Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Ser-	773
		manet, Scott Reed, Dragomir Anguelov, Dumitru	774
		Erhan, Vincent Vanhoucke, and Andrew Rabinovich.	775

776 2015. Going deeper with convolutions. In *Proceed-*  
777 *ings of the IEEE conference on computer vision and*  
778 *pattern recognition*, pages 1–9.

779 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,  
780 Barret Zoph, Sebastian Borgeaud, Dani Yogatama,  
781 Maarten Bosma, Denny Zhou, Donald Metzler, et al.  
782 2022. Emergent abilities of large language models.  
783 *arXiv preprint arXiv:2206.07682*.

784 Zhibiao Wu and Martha Palmer. 1994. Verb seman-  
785 tics and lexical selection. *arXiv preprint cmp-*  
786 *lg/9406033*.

787 Jianxiong Xiao, James Hays, Krista A Ehinger, Aude  
788 Oliva, and Antonio Torralba. 2010. Sun database:  
789 Large-scale scene recognition from abbey to zoo. In  
790 *2010 IEEE computer society conference on computer*  
791 *vision and pattern recognition*, pages 3485–3492.  
792 IEEE.

793 Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Ye-  
794 ung, Mojtaba Seyedhosseini, and Yonghui Wu. 2022.  
795 Coca: Contrastive captioners are image-text founda-  
796 tion models. *arXiv preprint arXiv:2205.01917*.

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