

Waffling around for Performance: Visual Classification with Random Words and Broad Concepts

Karsten Roth¹, Jae Myung Kim¹, A. Sophia Koepke¹, Oriol Vinyals², Cordelia Schmid³, Zeynep Akata^{1,4}

¹University of Tübingen, ²Google DeepMind,

³Inria, Ecole normale supérieure, CNRS, PSL Research University, ⁴MPI for Intelligent Systems

Abstract

The visual classification performance of vision-language models such as CLIP can benefit from additional semantic knowledge, e.g. via large language models (LLMs) such as GPT-3. Further extending classnames with LLM-generated class descriptors, e.g. “waffle, which has a round shape”, or averaging retrieval scores over multiple such descriptors, has been shown to improve generalization performance. In this work, we study this behaviour in detail and propose WaffleCLIP, a framework for zero-shot visual classification which achieves similar performance gains on a large number of visual classification tasks by simply replacing LLM-generated descriptors with random character and word descriptors **without** querying external models. We extend these results with an extensive experimental study on the impact and shortcomings of additional semantics introduced via LLM-generated descriptors, and showcase how semantic context is better leveraged by automatically querying LLMs for high-level concepts, while jointly resolving potential class name ambiguities. Link to the codebase: <https://github.com/ExplainableML/WaffleCLIP>.

1. Introduction

Task-specific natural language prompts [31, 67, 8, 26] improve the performance of large vision-language models (VLMs) [47]. However, if the model does not have access to additional training data, i.e. in the zero-shot setting, prompt tuning is not an option. Instead, a promising alternative [42, 46, 36] is querying large language models (LLMs) to provide additional semantic context to enrich class representations. Extending classnames with fine-grained class descriptors generated by GPT-3 [5] via minimal human intervention boosts results [36, 46]. In particular, [36] use class-based descriptors on top of classnames, e.g. *a round shape* for *waffle*, and provide experimental evidence that additional semantic cues obtained this way are beneficial.

However, a closer inspection of GPT-3 generated de-

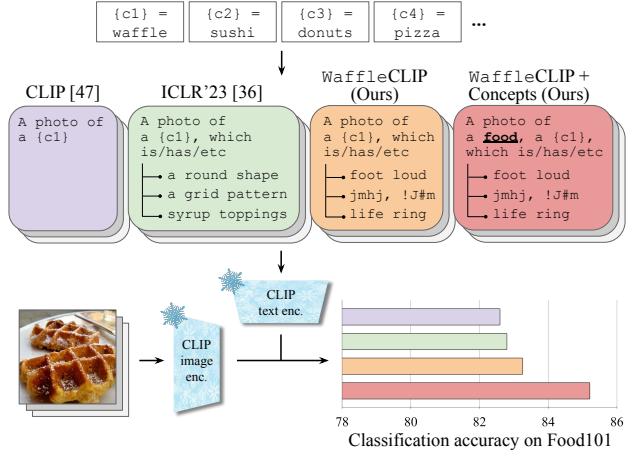


Figure 1: Substituting GPT-3 generated fine-grained descriptors with random word or character sequences yields competitive classification performance. Using high-level concepts in the prompt further allows classname ambiguities to be addressed for additional improvements.

scriptors indicates a high degree of diversity, limited visual relevance and ambiguity [36]. This means that multiple descriptors can get assigned to a class despite them likely not co-occurring, e.g. “steamed” and “fried”, can contain non-visual attributes, e.g. “a sour and spicy smell”, or can be associated with an ambiguous class interpretation, e.g. “webbed feet” for “Peking duck” as a food item. Hence, the underlying drivers of performance improvements when using generated fine-grained class descriptors are unclear.

To understand what is required to achieve these performance gains, we first evaluate a variant of [36] and show that each set of class-specific GPT-3 generated descriptors can be replaced with a fixed set of randomly selected, class-independent descriptors while still retaining similar benefits in performance. Motivated by this observation, we take this one step further and propose WaffleCLIP, named after *waffling around* the class name, that replaces the LLM-generated fine-grained descriptors, e.g. *a round shape*, *a grid pattern*, with random words (e.g. “foot loud”) or char-

acter lists (e.g. "jmhj, !J#m") based on average class name length and word counts (cf. Figure 1). As WaffleCLIP doesn't require access to LLMs for additional context (unlike e.g. [36, 42, 46, 56]), it remains *inherently zero-shot*.

Naturally, the convincing performance of WaffleCLIP across benchmarks raises questions regarding the true benefits of additional semantics introduced by LLM-generated descriptors. We provide answers with extensive experiments, showcasing that semantic descriptors produced by LLMs offer a *structurally different* and *complementary* impact on the classification behavior. However, we find this not to be fully driven by additionally introduced semantics, but rather a different form of structured noise ensembling. Instead, we show that actual semantic context is better introduced through coarse-grained, high-level concepts. Given access to external LLMs, we suggest a query mechanism for GPT-3 to automatically generate these (e.g. *food* for *waffle, peking duck*), while jointly resolving issues of context-dependent class label ambiguity for further gains.

In summary, our contributions are: **1)** We motivate and propose WaffleCLIP to use random character and word descriptors to enhance the semantic retrieval process in VLMs (particularly CLIP); **2)** we demonstrate that WaffleCLIP yields similar or better zero-shot image classification performances compared to methods reliant on external LLM-generated descriptors; **3)** we extensively study the semantic context introduced through LLM-generated descriptors and propose (automatically extracted) high-level LLM-generated concepts as an alternative for better use of semantics while tackling classname ambiguities.

2. Related Work

Image classification with VLMs such as CLIP [47] has gained popularity particularly in low-data regimes. As input prompts have a significant impact on the performance, recent research has focused on the exploration of learnable prompts for the text encoder [67, 66, 33, 54], the visual encoder [1, 7, 61, 32] or for both encoders jointly [62]. Alternatively synthetic images generated from theclassnames can support image classification [56, 2, 20]. In contrast, we do not tune prompts or query image generation methods, but propose to use prompts containing random characters or words to enhance the zero-shot capabilities of VLMs.

Adding external knowledge to language prompts. Recently, multiple works have leveraged LLMs to obtain more effective prompts. [46, 38, 35] utilized GPT-3 [5] to produce and study lengthy, descriptive sentences that articulate the visual concepts of each category, while [42] generated semantic hierarchies to identify subclasses of categories for zero-shot class prediction. [36] used multiple fine-grained LLM-generated class descriptors, which enhance accuracy and appear to provide interpretability by assigning weights to each descriptor. Similarly, different kinds of descriptions

have been used for image classification, by manually crafting descriptions [48, 21], or by utilizing external databases based on Wikipedia [16, 45, 39, 11], the WordNet hierarchy [37, 52, 49], or the ImageNet-Wiki [6]. Whilst external knowledge from LLMs can be valuable, we can match the image classification performance of using fine-grained LLM-generated descriptors with randomly sampled characters and words as class descriptors. In addition, we find that if semantic context is available through LLMs, it is better integrated through high-level context (c.f. e.g. also [15]), for which we provide an automatic extraction mechanism.

Noise augmentation. Data augmentation through noise is known to enhance the performance and robustness of model training for a variety of tasks and domains [53, 17]. In the language domain, noise can be incorporated in the embedding or input space. For instance, [55, 9, 19] used linguistic embedding space augmentations inspired by mixup [64], and [10] added Gaussian embedding space noise. Augmentation through input space noise has been performed at the word- [28, 60], token- [60] or character-level [51, 22, 3, 41]. For character-level noise augmentation, characters are randomly substituted, added or removed [22, 3, 41]. In all cases, these augmentation are used to prevent overfitting *during training*. Instead, our approach utilizes character- and word-level language augmentation to perturb the class prompts for improved zero-shot image classification.

3. Method

We first describe image classification using class descriptors following [36] (§3.1), before motivating and explaining our LLM-free, random semantic descriptor alternative WaffleCLIP (§3.2). Finally, if LLMs are available, we highlight a simple extension to incorporate semantics while jointly resolving ambiguities via automatic high-level semantic concepts extraction for additional benefits (§3.3).

3.1. Image classification with class descriptors

Given target categories C and a query image x , the zero-shot image classification protocol used in CLIP [47] defines the classification problem as nearest neighbour retrieval:

$$\tilde{c} = \arg \max_{c \in C} s(\phi_I(x), \phi_L(f(c))), \quad (1)$$

with prompt $f(c) = \text{"A photo of a }\{c\}\text{."}$ and image and language encoder ϕ_I and ϕ_L . To improve the retrieval process, [36] converts the simple class-embedding retrieval to a dictionary-based one, where a class c is associated with a set of descriptors D_c via " $\{c\}$ which (is/has/etc) {descriptor}." with e.g. $c = \text{"waffle"}$ and descriptor = "a round shape". Given D_c for classes c , classification is reformulated as

$$\arg \max_{c \in C} \frac{1}{|D_c|} \sum_{d \in D_c} s(\phi_I(x), \phi_L(d)), \quad (2)$$

ViT-B/32	ImageNetV2	ImageNet	CUB200	EuroSAT	Places365	Food101	Oxford Pets	DTD	Avg
CLIP [47]	54.71	62.01	51.28	40.78	39.12	82.59	85.06	43.18	57.34
DCLIP [36]	55.82	63.12	52.47	43.29	40.47	82.79	86.54	43.99	58.56
DCLIP (same, 1x)	55.47 \pm 0.24	62.89 \pm 0.19	52.64 \pm 0.28	39.74 \pm 2.69	40.29 \pm 0.47	83.82 \pm 0.48	87.04 \pm 0.27	43.35 \pm 0.41	58.16 \pm 1.01
DCLIP (same, 2x)	55.75 \pm 0.21	63.10 \pm 0.19	52.72 \pm 0.23	39.73 \pm 1.66	40.61 \pm 0.22	84.01 \pm 0.23	87.10 \pm 0.14	43.29 \pm 0.22	58.29 \pm 0.62

Table 1: **Motivating random class descriptors.** Comparing CLIP [47] and GPT-descriptor-extended CLIP [36] (DCLIP) with the same set of randomly sampled descriptors used for each class, where the set size is either the average number of descriptions per class in DCLIP (*same, 1x*), or twice that (*same, 2x*). A random set of descriptors per class can match or even outperform DCLIP across backbone architectures confirming that randomized prompt averaging leads to higher performance.

which defines the similarity between image x and class c as the average similarity to all its descriptor variants. We abbreviate this descriptor-based extension of CLIP as *DCLIP*.

3.2. WaffleCLIP

DCLIP [36]¹ requires external LLMs for descriptors that convert the single-class matching problem to one over an ensemble of fine-grained class representations.

Motivation. However, we observe that various such LLM-generated descriptors reveal high diversity, limited visual relevance, and ambiguity. From a conceptual perspective, this makes it hard to pin down the exact benefits of generated class descriptors used e.g. in [46] or [36]. To understand a possible driver of performance improvements, we conduct a simple experimental study, shown in Tab. 1. We take all available LLM-generated descriptors for a dataset from [36], sample a small set of descriptors where the cardinality of the set is the average number of descriptors per class used in DCLIP, and assign this same set of random descriptors to every class, i.e. *DCLIP (same, 1x)*. This shows a close match to DCLIP (e.g. 58.56% and 58.16% for ViT-B/32, 69.14% and 68.80% for ViT-L/14, and 54.77% and 54.71% for ResNet50 in total average) and in parts even better performance (e.g. 0.83%, 0.34%, 1.49% improvement in Food101 for ViT-B/32, ViT-L/14, ResNet50, respectively). This reveals averaging over descriptor variations as one of the key drivers for performance. The results further improve when increasing the number of random LLM-generated descriptors for each class (*DCLIP (same, 2x)*, e.g. 58.16% \rightarrow 58.29% on ViT-B/32 or 68.80% \rightarrow 69.12% on ViT-L/14). This indicates that the role of additional descriptor semantics is likely overestimated, especially when uncurated descriptors are used. Building on the benefits of averaging over various prompt variants to extract a better semantic representation estimate of an associated class, we investigate whether fully randomized prompt descriptors can provide similar benefits, **without** querying external LLMs.

WaffleCLIP. This motivates WaffleCLIP, an *LLM-free* descriptor alternative that uses simple randomized de-

¹DCLIP [36] reports improvements over CLIP by using the phrase " $\{c\}$, which (is/has/etc) {descriptor}." instead of "A photo of $\{c\}$, which (is/has/etc) {descriptor}." as suggested in the original CLIP paper. For fair comparison with CLIP, we utilize the latter.

scriptors. In particular, we populate D_c with class-independent, random word sequences or random character lists, with fixed number of characters per word, l_w , and fixed number of words, n_w . For example, $l_w = 4$ and $n_w = 2$ for `char_seq_1 = "aks@, pg2f"` in Figure 2. To avoid introducing hyperparameters, we leverage a simple heuristic where the average number of words and average number of characters per word in the provided class labels determines l_w and n_w . As a result, this converts the standard CLIP input prompt "A photo of a $\{c\}$." into "A photo of a $\{c\}$, which (is/has/etc) {random_sequence}.", where we follow the extension structure used in [36].

3.3. Better semantics and reduced ambiguity via high-level concepts

Due to the limited impact of additional semantics introduced by fine-grained descriptors (c.f. §3.2), we propose an alternative way of querying LLMs, that does not require averaging across multiple descriptors and jointly addresses the issue of class ambiguities. Therefore, we suggest taking a step back and to search not for additional class details, but instead for higher-level commonalities *between* the classes, akin to the use of class hierarchies in image classification [18]. Understanding commonalities between multiple target classes can help resolve ambiguities. If the class "boxer" is seen in the context of animal classification, it likely refers to the animal instead of a human athlete. We propose to automatically produce such high-level concepts by using available class names (or subsets if the class count exceeds the maximum LLM input sequence length) $C_{\mathcal{D}}$ for a dataset \mathcal{D} and querying GPT-3 [5] with:

```
"Q: Tell me in five words or less what
{list_of_classes} have in common. It may be
nothing. A: They are all ".
```

After extracting the shared concept, a simple concept filtering is attached that checks if generated concepts fall into non-specific categories, namely "Object", "Thing", "Verb", "Adjective", "Noun" or "Word". If so, high-level concept guidance is omitted (only the case for three out of eleven visual classification benchmarks, see also §4).

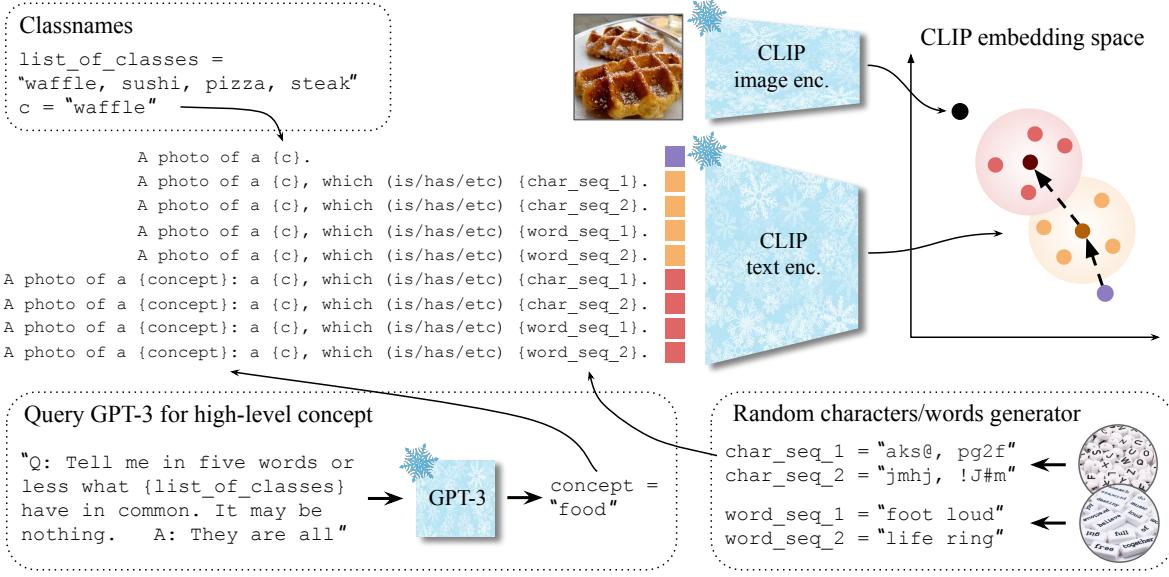


Figure 2: **Visual classification with WaffleCLIP using random characters/words.** Introducing character-level or word-level noise following the classname increases the similarity between the ground-truth text features and image features (orange). WaffleCLIP can be further enhanced by adding a high-level concept descriptor in the prompt (red).

We then augment the default CLIP prompt to "A photo of a {concept}: a {c}." and for WaffleCLIP, the prompt is then extended to "A photo of a {concept}: a {c}, which (is/has/etc) {random_sequence}.". While the prompt style can likely be improved, this naive extension already offers remarkable benefits.

4. Experiments

We start with implementation details before comparing WaffleCLIP to DCLIP in §4.1. Extending our observational experiments in Tab. 1, we study the source of performance gains via LLM-generated descriptors (§4.2) and show a better way to introduce semantics into the retrieval process while tackling semantic ambiguities via automatic high-level concept extraction (§4.3). Finally, §4.4 provides additional insights on additional (OOD) benchmarks and a comparison to prompt ensembles and latent space noise. For additional details, see Supplementary.

Implementation details. We utilize CLIP [47] as the underlying VLM for WaffleCLIP. As there is no direct cost associated with generating random character or word sequences, their number is only bounded by inference speed requirements (which is minimal as all respective language embeddings can be computed *a priori* [36]). However, we find diminishing returns for very high numbers (see also §4.4), and use 30 random descriptors per class (or 15 random character and word descriptor pairs) if not mentioned otherwise, with similar performance for both half or double the descriptor count (c.f. §4.4). All experiments use PyTorch [44] and are conducted on a single NVIDIA

3090Ti. Wherever necessary, fine-grained LLM-generated descriptors are either taken from or generated following the codebase provided by [36], which we build on. If not mentioned explicitly, every result involving WaffleCLIP is computed over at least seven random seeds.

Benchmarks. The datasets considered are (mostly from [36]) ImageNet [14] and ImageNetV2 [29], CUB200-2011 [57] (fine-grained bird classification data), EuroSAT [23] (satellite image recognition data), Places365 [65] with scene imagery, Food101 [4] with different food classes, Oxford IIIT Pets [43], DTD (Describable Textures Dataset, [12]), Flowers102 [40], FGVC Aircraft [34] and Stanford Cars [30].

High-level concepts. Following §3.3, the GPT-3 generated high-level concept for CUB200-2011 is "Bird", "Land Use" for EuroSAT, "Place" for Places365, "Food" for Food101 and "Breed" for Oxford Pets. For additional benchmarks, extracted concepts are noted in the resp. section §4.4. For ImageNet (V2) and DTD, the concepts are too generic and thus filtered out ("Object", "Noun" or "Adjective"), with high-level guidance omitted.

4.1. WaffleCLIP vs LLM-generated descriptors

We start by analyzing the impact of randomization beyond fixed, randomized sets of fine-grained LLM-generated descriptors as done in Tab. 1, by instead using randomized character or word descriptors through our proposed WaffleCLIP. For that, we investigate visual classification accuracies across the eight diverse benchmarks studied in [36] in Tab. 2, where we compare WaffleCLIP, which does

ViT-B/32	ImageNetV2	ImageNet	CUB200	EuroSAT	Places365	Food101	Oxford Pets	DTD	Avg
CLIP [47]	54.71	62.01	51.28	40.78	39.12	82.59	85.06	43.18	57.34
+ Concepts	↓	↓	52.23	48.86	39.31	84.66	86.73	↓	58.96
DCLIP [36]	55.82	63.12	52.47	43.29	40.47	82.79	86.54	43.99	58.56
WaffleCLIP (ours)	55.92 ±0.08	63.31 ±0.09	52.38 ±0.12	44.31 ±1.07	40.56 ±0.07	83.25 ±0.21	85.70 ±0.25	43.16 ±0.25	58.57 ±0.41
+ Concepts	↓	↓	52.83 ±0.19	48.51 ±0.70	40.97 ±0.08	85.21 ±0.06	87.52 ±0.10	↓	59.47 ±0.42
+ GPT descr. + Concepts	↓	↓	52.77 ±0.26	51.64 ±0.25	41.35 ±0.09	84.87 ±0.05	87.71 ±0.18	↓	60.21 ±0.20

Table 2: **Image classification performance with WaffleCLIP.** WaffleCLIP extends input prompts with random word and character sequences and matches the performance of DCLIP [36] which utilizes GPT-generated class-descriptors. We additionally find that random word and character noise complements the use of GPT-generated descriptors (+ *GPT descr.*). Introducing additional semantic context through coarse, high-level concepts (+ *Concepts*) can offer an additional boost in performance, particularly on benchmarks where some classnames may be generic or ambiguous. We use (↓) to denote same results as previous lines where high-level concept guidance is not applicable. For results on ViT-L/14 and RN50, see Supp.

not use any external LLMs, with DCLIP. We find that averaging over randomized descriptors yields performances comparable to or better than those obtained with LLM-generated fine-grained descriptors over a majority of studied datasets, with average performance similarly matching: 58.56% using DCLIP versus 58.57% for WaffleCLIP with a ViT-B/32 backbone and (see Supp.) 69.14% → 68.95% for ViT-L/14, and 54.77% → 54.20% for ResNet50. Beyond the inherently zero-shot nature of WaffleCLIP and ease of use, these results highlight that improved visual classification with pretrained VLMs does not require external LLMs, and further cements prompt averaging as a potential key driver behind DCLIP.

4.2. Are descriptors from LLMs obsolete?

Our results above question the benefits of LLM-generated fine-grained semantics, as averaging over fully randomized character and word sequences achieves comparable performance. But does that mean that there is no benefit in leveraging descriptors produced by LLMs?

Impact of Averaging. To better understand this, we extend our motivational experiments from Tab. 1. First, we look at what happens when not performing averaging over all image-descriptor distances as in DCLIP, but instead choosing the maximum. If additional fine-grained semantics were indeed beneficial, selecting the most suitable one should similarly raise the performance. However, as Tab. 3 reveals, performance actually drops, highlighting that the VLM can not leverage the additional semantics to improve visual classification performance². Instead, it again points to descriptor ensembling *as the main driver in performance*.

We further support this by studying additional descriptor randomization variants beyond those in §3.2. In particular, instead of swapping specific descriptors, we interchange full class-specific descriptor lists (*interchanged*). As descriptions often contain class-specific keywords, this models a systematic semantic shift away from the actual class. Additionally, we evaluate shuffling words within a

²This is potentially influenced by bag-of-words behaviour of CLIP-like VLMs [63], which we leave to future research to study in more detail.

ViT-B/32	ImageNetV2	ImageNet	CUB200	EuroSAT	Places365
CLIP [47]	54.71	62.01	51.28	40.78	39.12
DCLIP [36] (<i>mean</i>)	55.82	63.12	52.47	43.29	40.47
DCLIP [36] (<i>max</i>)	54.41	61.67	52.40	37.11	37.21
Food101	Oxford Pets	DTD	Flowers102	FGVCAircraft	Stanford Cars
82.59	85.06	43.18	62.89	24.99	58.54
82.79	86.54	43.99	64.01	26.94	57.08
82.37	88.03	43.35	63.62	25.77	56.21
					Avg
					55.01
					56.05
					54.74

Table 3: **Importance of semantics in DCLIP.** We compare DCLIP with similarity score averaging (*mean*) and maximum similarity selection (*max*). As can be seen, simply taking the most similar entry even underperforms CLIP on average. This points to the limited impact of additional LLM-generated semantics on the improved visual classification.

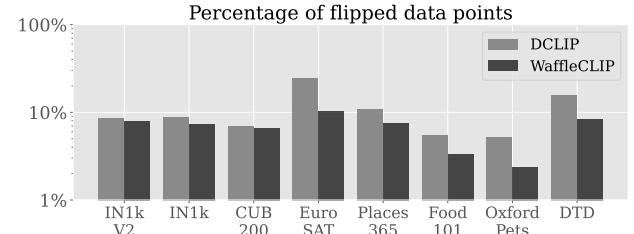


Figure 3: Label flipping experiment from CLIP to DCLIP or WaffleCLIP. Each bar indicates the percentage of data points getting either positively or negatively flipped (i.e. labelled correctly or incorrectly adjusted) when switching from CLIP to either DCLIP or WaffleCLIP. The consistently higher flip percentage indicates structural differences between natural language descriptors and randomized ones.

descriptor list (*shuffled*), and descriptor lists subsampled from all available ones (*random*). This gives a progression from systematic to more independent descriptor randomization. And indeed, our results in Tab. 4 reveal that directly interchanging full *class-dependent* descriptor lists (*interchanged*) drops performance significantly (e.g. from 58.56% to 55.03 on ViT-B/32). In cases where no such shift is happening, we find performances to match that of DCLIP (e.g. 86.54% → 86.28% on Oxford Pets). Similarly, when moving from a systematic shift closer to fully randomized descriptors (*scrambled* with 58.56% → 57.55% to *ran-*

ViT-B/32	ImageNetV2	ImageNet	CUB200	EuroSAT	Places365	Food101	Oxford Pets	DTD	Avg
DCLIP [36]	55.82	63.12	52.47	43.29	40.47	82.79	86.54	43.99	58.56
DCLIP (interchanged)	52.51 ±0.42	59.62 ±0.13	52.52 ±0.41	33.63 ±4.16	35.52 ±0.32	81.71 ±0.35	86.28 ±0.50	38.42 ±1.14	55.03 ±1.56
DCLIP (scrambled)	55.12 ±0.12	62.57 ±0.12	52.18 ±0.28	40.48 ±2.52	39.91 ±0.08	82.46 ±0.13	86.10 ±0.40	41.58 ±0.31	57.55 ±0.92
DCLIP (random, 1x)	54.11 ±0.28	61.37 ±0.18	52.42 ±0.19	36.83 ±4.27	38.80 ±0.26	82.86 ±0.23	85.99 ±0.62	42.20 ±0.85	56.82 ±1.57
DCLIP (random, 5x)	55.43 ±0.12	62.81 ±0.05	52.66 ±0.17	38.57 ±1.52	40.54 ±0.05	84.03 ±0.11	86.75 ±0.21	43.41 ±0.74	58.02 ±0.61

Table 4: **Progression from systematic to fully randomized descriptor scrambling.** To model a systematic semantic shift, we randomly swap descriptor lists between classes (*interchanged*) which often contain class-specific keywords, before progressing to more independently randomized descriptors by shuffling descriptor words within the classes (*scrambled*) and randomly sampling LLM-generated descriptors for each class (*random*) from the complete set of descriptors, with counts as in DCLIP (*1x*), or five times that (*5x*). As can be seen, a systematic shift results in a notable performance drop, while a move towards more independently randomized descriptors can recover the DCLIP performance, aligning with the observation for WaffleCLIP that fully randomized prompt averaging is the main performance driver.

dom with 58.56% → 58.02%, see Supp. for more results), we move closer to DCLIP performance. While this offers further evidence for WaffleCLIP and the fact that class-dependent ensembling drives gains, it does not yet allow us to directly compare the impact on the prediction behavior between LLM-generated descriptors and randomized ones.

Structural differences. We consider the percentages of samples that get positively or negatively flipped - i.e. ones that are classified correctly while previously being classified incorrectly (and vice versa) - when moving from CLIP to either DCLIP or WaffleCLIP in Fig. 3. We find that using LLM-generated fine-grained descriptors flips significantly more predictions than randomized words and characters, even in cases where WaffleCLIP outperforms DCLIP. For example, DCLIP achieves 43.29% compared to WaffleCLIP with 44.31% on EuroSAT or 82.79% to 83.35% on Food101 in Tab. 2, but DCLIP still flips a significantly larger portion of samples than WaffleCLIP on those datasets in Fig. 3. This reveals that full sentence, LLM-generated descriptors have a *structurally different* impact on the zero-shot classification process, which we find to operate *complementary* to randomized ones (see Tab. 2, + *GPT descr.*), where the use of both descriptor types leads to additional performance improvements over WaffleCLIP (e.g. 58.57% → 58.93% for ViT-B/32 or 54.20% → 55.22% for ResNet50). This means that even if additional semantics are not the guiding factor, LLMs for structured descriptor generation can still facilitate the extraction of a more robust class embedding. Furthermore, even with access to an external model, WaffleCLIP can provide additional benefits.

4.3. Semantic guidance via high-level concepts

While we verified the relevance of additional semantic context through fine-grained descriptors, methods using additional fine-grained class information [36, 42, 46] suffer from the inherent ambiguities in some class names. As proposed in §3.3, our aim is to understand if high-level semantic context can be used to resolve such ambiguities and provide high-level semantic guidance for the class-retrieval process. Our results with extracted high-level concepts in

CUB200	52.23	52.62	51.47	52.47	51.50
EuroSAT	41.89	48.86	40.61	47.81	44.19
Places365	37.65	38.71	39.31	38.12	37.28
Food101	79.86	81.98	83.35	84.66	79.41
Ox. Pets	83.65	82.42	79.91	83.54	86.73
	"bird"	"land use"	"place"	"food"	"breed"

Figure 4: Study of semantic impact of GPT-3 generated high-level concepts. We find that interchanging the concepts generally reduces performance, indicating that high-level concepts provide complementary semantic context.

Tab. 2, i.e. (+ *Concepts*), demonstrate across most benchmarks and backbones consistent and significant improvements, when used with CLIP, with WaffleCLIP, and even when used alongside WaffleCLIP and DCLIP. These improvements are especially evident on benchmarks with ambiguous (e.g. Food101) or generic labeling (e.g. EuroSAT, with labels such as *Industrial* or *Residential*): For ViT-B/32, classification accuracy increases from 40.78% to 48.86% when applied to CLIP, with similarly high improvements for ViT-L/14 (56.03% → 61.23%) or ResNet50 (28.09 → 34.06). Overall, the average classification accuracy also increases consistently (e.g. from 57.34% to 58.96% with ViT-B/32). This beats even DCLIP, while only being applicable on five out of eight benchmarks (58.96% versus 58.56%). When applied to WaffleCLIP, improvements across most benchmark and backbone settings are also significant, although we find diminishing returns on the largest backbone, ViT-L/14, with average performance increasing only from 68.95% to 69.12%. This might be due to its capabilities of retaining the most common concepts associated with specific classes, resulting in a robust class retrieval setup when averaging over multiple randomized descriptor variants.

We verify the benefits of high-level semantics further by looking at performance changes when concepts are interchanged (Fig. 4). For most benchmarks, the highest improvements are obtained with respective GPT-generated concepts. Some off-diagonal terms with higher scores, e.g. CUB200 where "bird" performs similar to/worse than "land use"/"food", do appear out of distribution, and

ViT-B/32	Flowers102	FGVCAircraft	Stanford Cars	Avg
CLIP [47]	62.89	24.99	58.54	48.81
DCLIP [36]	64.01	26.94	57.08	49.34
WaffleCLIP	66.27 ± 0.26	25.66 ± 0.19	58.91 ± 0.17	50.28 ± 0.21
+ Concepts	67.19 ± 0.19	28.44 ± 0.22	59.70 ± 0.12	51.78 ± 0.18
+ GPT desc. + Conc.	66.71 ± 0.39	28.96 ± 0.37	59.33 ± 0.14	51.67 ± 0.32

Table 5: We find similar performance improvements with WaffleCLIP and high-level concept guidance for three additional benchmarks, which in parts do not benefit from LLM-generated descriptors (see *Stanford Cars*).

Benchmarks	ImageNet-R [24]	ImageNet-S [58]	ImageNet-A [25]
CLIP [47]	65.97	40.73	29.63
DCLIP [36]	65.12	41.09	29.19
WaffleCLIP	67.31	42.00	31.52

Table 6: Performance gains of WaffleCLIP on distribution-shifted datasets further highlight general applicability through simple averaging over randomized descriptors, even if natural language ones fail.

Avg.	ViT-B/32	ViT-L/14	RN50
Joint	58.57 ± 0.41	68.95 ± 0.18	54.20 ± 0.23
Random Words	58.18 ± 0.44	68.73 ± 0.58	55.24 ± 0.41
Random Characters	58.59 ± 0.27	68.02 ± 0.14	53.79 ± 0.16

Table 7: **Randomized descriptor modes.** Considering both the joint usage of randomized word and character sequences compared to only randomized word or character sequences, joint usage provides the most consistent performance improvements across benchmarks and backbones.

warrant future research to improve our understanding of how semantics concepts are truly encoded in large VLMs.

However, seeing maximum performances primarily on the diagonal heuristically supports that additional semantics introduced as high-level concepts and commonalities, can offer reliable guidance. Indeed, considering a selection of ambiguous samples such as "Boxer" or "Sphynx" in the Oxford Pets dataset, "Mussels", "Oysters" or "Grilled Salmon" in the Food101 dataset, or highly generic labels such as "Industrial" or "Residential" in the EuroSAT satellite image dataset, we find a consistent increase in average similarity to all associated test images by up to 13%. This confirms that concept guidance can re-align and refine class embeddings based on the respective context.

4.4. Ablation studies

Evaluation on additional (OOD) benchmarks. For further evidence on the generality of WaffleCLIP and concept guidance, we study three additional benchmarks beyond those in Tab. 2 and [36]: Flowers102 [40] (extracted concept: "flower"), FGVCAircraft [34] ("aircraft"), and StanfordCars [30] ("car"). Our results in Tab. 5 (and in the suppl. material for other backbones) again show consistent gains when going from CLIP to WaffleCLIP

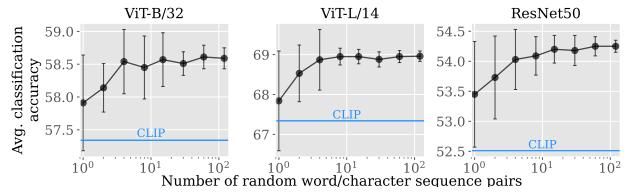


Figure 5: Ablation study on the number of randomized word and character descriptors used in WaffleCLIP. We find consistent competitive performance with just four randomized descriptor pairs (c.f. DCLIP Tab. 2). CLIP (blue line) is outperformed with just a single descriptor pair.

or WaffleCLIP + *Concepts*. Interestingly, DCLIP is detrimental on very fine-grained benchmarks like Stanford Cars, losing 1.46% against CLIP. We speculate that this is due to semantically similar descriptors for multiple classes that are coarser than the actual class label (e.g. "BMW Active Hybrid" and "BMW 1 Series" being assigned similar generic BMW descriptors). Consequently, embeddings of related classes are systematically moved too close, deteriorating the performance. Meanwhile, WaffleCLIP (+ *Concepts*) can still offer performance boosts ($58.54\% \rightarrow 58.91\% \rightarrow 59.70\%$). Furthermore, we study WaffleCLIP on OOD benchmarks: Adversarial natural images (ImageNet-A, [25]), sketches (ImageNet-S, [58]) and renditions (ImageNet-R, [24]). Results in Tab. 6 show that while DCLIP does not improve consistently, WaffleCLIP operates well even for out-of-distribution data (e.g. $29.63\% \rightarrow 31.52\%$ on ImageNet-A).

Comparison to prompt ensembles. We also compare WaffleCLIP to prompt ensembling (c.f. e.g. [46]) with the same budget of 30 randomly selected prompt options from a list of eighty handcrafted ones (taken from [46], such as "A tattoo of a {class}.", "A {class} in a video game.", ...). Unlike WaffleCLIP, prompt ensembling still requires human input and design. Results on all eleven benchmarks are listed in Tab. 8, which favor WaffleCLIP- outperforming prompt ensembling in eight out of eleven benchmarks and comparable performance on the remaining ones. In particular, improvements over prompt ensembling are higher than the improvement of prompt ensembling over vanilla CLIP ($56.32\% \rightarrow 55.58\% \rightarrow 55.01\%$). This further supports the benefit of extracting more robust semantic representations, for which randomized descriptors provide a cheap and suitable tool. In addition to that, we highlight the complementarity of high-level concept guidance also with prompt ensembling in Tab. 8 (wherever the classname is included, we simply use "a {concept}: a {classname}" instead), raising the average classification accuracy from 55.58% to 56.94% (compared to the base CLIP accuracy of 55.01%).

Comparison to latent space noise. To highlight that a more robust extraction of semantics through class-

ViT-B/32	IN1k-V2	IN1k	CUB	Euro	Places	Food	Pets	DTD	Flowers	FGVC	Cars	Avg
CLIP [47]	54.71	62.01	51.28	40.78	39.12	82.59	85.06	43.18	62.89	24.99	58.54	55.01
DCLIP [36]	55.82	63.12	52.47	43.29	40.47	82.79	86.54	43.99	64.01	26.94	57.08	56.05
P. Ensemble + Concepts	55.49 ± 0.21	62.79 ± 0.29	51.46 ± 0.43	45.76 ± 0.49	40.58 ± 0.06	82.67 ± 0.37	83.20 ± 0.72	42.53 ± 0.54	63.30 ± 0.33	25.14 ± 0.45	58.38 ± 0.29	55.58 ± 0.42
	\downarrow	\downarrow	52.08 ± 0.17	49.80 ± 0.66	40.61 ± 0.14	84.45 ± 0.15	87.42 ± 0.20	\downarrow	65.38 ± 0.27	26.64 ± 0.50	59.12 ± 0.14	56.94 ± 0.34
WaffleCLIP + Concepts	55.92 ± 0.08	63.31 ± 0.09	52.38 ± 0.12	44.31 ± 1.07	40.56 ± 0.07	83.25 ± 0.21	85.70 ± 0.25	43.16 ± 0.25	66.27 ± 0.26	25.66 ± 0.19	58.91 ± 0.17	56.31 ± 0.37
	\downarrow	\downarrow	52.83 ± 0.19	48.51 ± 0.70	40.97 ± 0.08	85.21 ± 0.06	87.52 ± 0.10	\downarrow	67.19 ± 0.19	28.44 ± 0.22	59.70 ± 0.12	57.52 ± 0.26

Table 8: **Prompt ensembling versus WaffleCLIP (+concepts).** Across all visual classification benchmarks, we find matching or improved performance of WaffleCLIP compared to prompt ensembling (improving on eight out of eleven benchmarks), with the increase in average classification performance of WaffleCLIP compared to prompt ensembling higher than the increase of a prompt-ensembled version compared to standard CLIP, **without** requiring a handcrafted list of prompts.

conditioned randomization on the input level is crucial, we also compare to randomization directly in the (hyperpherical) latent space. For that, we choose a von-Mises-Fisher distribution (as commonly utilized to model uni-model distributions on the hypersphere [59, 13, 68, 50]):

$$p(\hat{\phi}^c | \phi_l^c, \kappa) = \mathcal{C}_d(\kappa) \exp\left(\kappa \phi_l^c T \hat{\phi}^c\right), \quad (3)$$

centered around default class embeddings ϕ_l^c with constant normalization function $\mathcal{C}_d(\kappa)$ only dependent on the input dimensionality and concentration κ . To sample from a vMF distribution around each class embedding, we leverage the sampler utilized in [13, 27] with the same budget of 30 noise embeddings. Average classification performance as a function of the (inverse) concentration κ is visualized in Fig. 6. As can be seen, for high concentrations (i.e. random embedding samples are placed close to the mean direction), one can replicate the default performance of a single CLIP embedding. For higher variances, performance continuously drops, with a hard inflection at around $\kappa \approx 500$. This serves to show that class-conditioned randomized descriptors as used in WaffleCLIP are crucial to providing a more robust estimate of semantic concepts, and cannot be simulated through simple embedding space noise.

Dependence on descriptor counts. We study the impact of the randomized word and character sequence pair count for WaffleCLIP in Fig. 5. A value of one indicates a single pair comprising a random words and characters descriptor, respectively. We achieve competitive performance already with 4 to 15 descriptor pairs (c.f. DCLIP in Tab. 2), while consistently outperforming CLIP (blue line) even with a single randomized descriptor pair. As class embeddings can be computed *a priori*, the impact on overall inference time is low, making WaffleCLIP and its extensions very attractive for enhancing image classification performance of VLMs.

Impact of randomization types. Finally, we analyze how performance changes when either only using random character sequences or only random word sequences, instead of a combination of both as in WaffleCLIP. Across benchmarks and architectures (see Tab. 7), we observe dichotomies in performance between either random word or random character sequences, often performing either best or worst on a specific benchmark and backbone, while

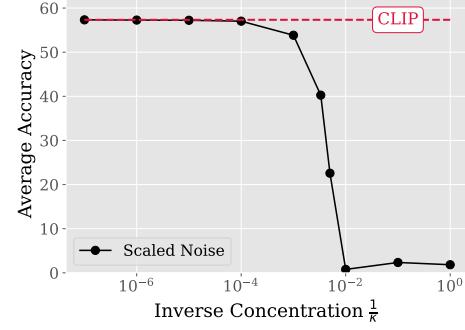


Figure 6: We compare latent space noise (vMF distribution around default CLIP language embeddings) against default CLIP. Reducing latent noise (i.e. increasing concentration κ) converges to initial performance, highlighting no notable benefit of deploying noise in the latent space.

the joint usage of both random words and character sequences strikes a consistent and best transferable average improvement across benchmarks and backbone architectures. Therefore, we chose the joint usage of both random words and characters as our default setup.

5. Conclusion

In this work, we systematically examined the benefits of using LLM-generated additional class descriptors for improved training-free image classification with vision-language models (VLMs). In-depth studies reveal how similar performance gains can be achieved by replacing these LLM-generated descriptors with randomized ones, giving rise to WaffleCLIP. Even though WaffleCLIP is entirely zero-shot as it does not require access to external LLMs, across eleven visual classification benchmarks, we find comparable or better results than those obtained when using fine-grained GPT-3 generated descriptors, making WaffleCLIP very attractive for practical use in true zero-shot scenarios. We also show that VLM struggles to leverage the actual semantics introduced through LLM-generated descriptors, and instead show that if given access to external LLMs, semantics are better exploited through coarse-grained, high-level concepts. Using specific queries, we show how these can be automatically extracted, while jointly helping to address issues of class ambiguity.

Acknowledgements

This work was supported by DFG project number 276693517, by BMBF FKZ: 01IS18039A, by the ERC (853489 - DEXIM), and by EXC number 2064/1 – project number 390727645. Karsten Roth and Jae Myung Kim thank the European Laboratory for Learning and Intelligent Systems (ELLIS) PhD program and the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for support.

References

- [1] Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. Exploring visual prompts for adapting large-scale models. *arXiv:2203.17274*, 2022.
- [2] Hritik Bansal and Aditya Grover. Leaving reality to imagination: Robust classification via generated datasets. *arXiv:2302.02503*, 2023.
- [3] Yonatan Belinkov and Yonatan Bisk. Synthetic and natural noise both break neural machine translation. In *ICLR*, 2018.
- [4] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In *ECCV*, 2014.
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *NeurIPS*, 2020.
- [6] Sebastian Bujwid and Josephine Sullivan. Large-scale zero-shot image classification from rich and diverse textual descriptions. In *Workshop on Beyond Vision and LANguage: inTEGRating Real-world kNowledge*, 2021.
- [7] Aochuan Chen, Yuguang Yao, Pin-Yu Chen, Yihua Zhang, and Sijia Liu. Understanding and improving visual prompting: A label-mapping perspective. *arXiv:2211.11635*, 2022.
- [8] Guangyi Chen, Weiran Yao, Xiangchen Song, Xinyue Li, Yongming Rao, and Kun Zhang. PLOT: Prompt learning with optimal transport for vision-language models. In *ICLR*, 2023.
- [9] Jiaao Chen, Zichao Yang, and Diyi Yang. Mixtext: Linguistically-informed interpolation of hidden space for semi-supervised text classification. In *ACL*, 2020.
- [10] Yong Cheng, Zhaopeng Tu, Fandong Meng, Junjie Zhai, and Yang Liu. Towards robust neural machine translation. In *ACL*, 2018.
- [11] Subhabrata Choudhury, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. The curious layperson: Fine-grained image recognition without expert labels. In *BMVC*, 2021.
- [12] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *CVPR*, 2014.
- [13] Tim R. Davidson, Luca Falorsi, Nicola De Cao, Thomas Kipf, and Jakub M. Tomczak. Hyperspherical variational auto-encoders. *Conference on Uncertainty in Artificial Intelligence*, 2018.
- [14] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [15] Lisa Dunlap, Clara Mohri, Devin Guillory, Han Zhang, Trevor Darrell, Joseph E. Gonzalez, Aditi Raghunathan, and Anja Rohrbach. Using language to extend to unseen domains, 2023.
- [16] Mohamed Elhoseiny, Yizhe Zhu, Han Zhang, and Ahmed Elgammal. Link the head to the “beak”: Zero shot learning from noisy text description at part precision. In *CVPR*, 2017.
- [17] Steven Y Feng, Varun Gangal, Jason Wei, Sarah Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. A survey of data augmentation approaches for nlp. In *ACL-IJCNLP*, 2021.
- [18] Yanming Guo, Yu Liu, Erwin M Bakker, Yuanhao Guo, and Michael S Lew. Cnn-rnn: a large-scale hierarchical image classification framework. *Multimedia tools and applications*, 2018.
- [19] Xiaoshuai Hao, Yi Zhu, Srikanth Appalaraju, Aston Zhang, Wanqian Zhang, Bo Li, and Mu Li. Mixgen: A new multi-modal data augmentation. In *WACV*, 2023.
- [20] Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing Zhang, Philip Torr, Song Bai, and XIAOJUAN QI. Is synthetic data from generative models ready for image recognition? In *ICLR*, 2023.
- [21] Xiangteng He and Yuxin Peng. Fine-grained image classification via combining vision and language. In *CVPR*, 2017.
- [22] Georg Heigold, Stalin Varanasi, Günter Neumann, and Josef van Genabith. How robust are character-based word embeddings in tagging and mt against wrod scramlbng or randdm noise? In *Association for Machine Translation in the Americas*, 2018.
- [23] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2017.
- [24] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. *ICCV*, 2021.
- [25] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. *CVPR*, 2021.
- [26] Tony Huang, Jack Chu, and Fangyun Wei. Unsupervised prompt learning for vision-language models. *arXiv:2204.03649*, 2022.
- [27] Michael Kirchhof, Karsten Roth, Zeynep Akata, and Enkelejda Kasneci. A non-isotropic probabilistic take on proxy-based deep metric learning. In *ECCV*, 2022.
- [28] Sosuke Kobayashi. Contextual augmentation: Data augmentation by words with paradigmatic relations. In *NAACL-HLT*, 2018.
- [29] Simon Kornblith, Jonathon Shlens, and Quoc V. Le. Do better imagenet models transfer better? In *CVPR*, 2019.
- [30] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *IEEE Workshop on 3D Representation and Recognition*, 2013.

- [31] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hi-roaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 2023.
- [32] Jochem Loedeman, Maarten C Stol, Tengda Han, and Yuki M Asano. Prompt generation networks for efficient adaptation of frozen vision transformers. *arXiv:2210.06466*, 2022.
- [33] Yuning Lu, Jianzhuang Liu, Yonggang Zhang, Yajing Liu, and Xinmei Tian. Prompt distribution learning. In *CVPR*, 2022.
- [34] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv:1306.5151*, 2013.
- [35] Chengzhi Mao, Revant Teotia, Amrutha Sundar, Sachit Menon, Junfeng Yang, Xin Wang, and Carl Vondrick. Doubly right object recognition: A why prompt for visual rationales, 2023.
- [36] Sachit Menon and Carl Vondrick. Visual classification via description from large language models. In *ICLR*, 2023.
- [37] George A Miller. *WordNet: An electronic lexical database*. MIT press, 1998.
- [38] Muhammad Ferjad Naeem, Muhammad Gul Zain Ali Khan, Yongqin Xian, Muhammad Zeshan Afzal, Didier Stricker, Luc Van Gool, and Federico Tombari. I2mvformer: Large language model generated multi-view document supervision for zero-shot image classification. In *CVPR*, 2023.
- [39] Muhammad Ferjad Naeem, Yongqin Xian, Luc Van Gool, and Federico Tombari. I2dformer: Learning image to document attention for zero-shot image classification. In *NeurIPS*, 2022.
- [40] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *Indian Conference on Computer Vision, Graphics and Image Processing*, 2008.
- [41] Xing Niu, Prashant Mathur, Georgiana Dinu, and Yaser Al-Onaizan. Evaluating robustness to input perturbations for neural machine translation. In *ACL*, 2020.
- [42] Zachary Novack, Saurabh Garg, Julian McAuley, and Zachary C Lipton. Chils: Zero-shot image classification with hierarchical label sets. *arXiv:2302.02551*, 2023.
- [43] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *CVPR*, 2012.
- [44] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*. 2019.
- [45] Tzuf Paz-Argaman, Reut Tsarfaty, Gal Chechik, and Yuval Atzmon. Zest: Zero-shot learning from text descriptions using textual similarity and visual summarization. In *EMNLP*, 2020.
- [46] Sarah Pratt, Rosanne Liu, and Ali Farhadi. What does a platypus look like? generating customized prompts for zero-shot image classification. *arXiv:2209.03320*, 2022.
- [47] 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 *ICML*, 2021.
- [48] Scott Reed, Zeynep Akata, Honglak Lee, and Bernt Schiele. Learning deep representations of fine-grained visual descriptions. In *CVPR*, 2016.
- [49] Karsten Roth, Oriol Vinyals, and Zeynep Akata. Integrating language guidance into vision-based deep metric learning. In *CVPR*, 2022.
- [50] Karsten Roth, Oriol Vinyals, and Zeynep Akata. Non-isotropy regularization for proxy-based deep metric learning. In *CVPR*, 2022.
- [51] Gözde Güл Şahin. To augment or not to augment? a comparative study on text augmentation techniques for low-resource nlp. *Computational Linguistics*, 2022.
- [52] Sheng Shen, Chunyuan Li, Xiaowei Hu, Yujia Xie, Jianwei Yang, Pengchuan Zhang, Anna Rohrbach, Zhe Gan, Lijuan Wang, Lu Yuan, et al. K-lite: Learning transferable visual models with external knowledge. *arXiv:2204.09222*, 2022.
- [53] Connor Shorten and Taghi M Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of big data*, 2019.
- [54] Manli Shu, Weili Nie, De-An Huang, Zhiding Yu, Tom Goldstein, Anima Anandkumar, and Chaowei Xiao. Test-time prompt tuning for zero-shot generalization in vision-language models. *arXiv:2209.07511*, 2022.
- [55] Lichao Sun, Congying Xia, Wengpeng Yin, Tingting Liang, S Yu Philip, and Lifang He. Mixup-transformer: Dynamic data augmentation for nlp tasks. In *Computational Linguistics*, 2020.
- [56] Vishaal Udandarao, Ankush Gupta, and Samuel Albanie. Sus-x: Training-free name-only transfer of vision-language models. *arXiv:2211.16198*, 2022.
- [57] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- [58] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. In *Advances in Neural Information Processing Systems*, pages 10506–10518, 2019.
- [59] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *ICML*, 2020.
- [60] Jason Wei and Kai Zou. Eda: Easy data augmentation techniques for boosting performance on text classification tasks. In *EMNLP-IJCNLP*, 2019.
- [61] Junyang Wu, Xianhang Li, Chen Wei, Huiyu Wang, Alan Yuille, Yuyin Zhou, and Cihang Xie. Unleashing the power of visual prompting at the pixel level. *arXiv:2212.10556*, 2022.
- [62] Yinghui Xing, Qirui Wu, De Cheng, Shizhou Zhang, Guoqiang Liang, and Yanning Zhang. Dual modality prompt tuning for vision-language pre-trained model. *arXiv:2208.08340*, 2022.

- [63] Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it? In *The Eleventh International Conference on Learning Representations*, 2023.
- [64] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018.
- [65] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE TPAMI*, 2017.
- [66] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *CVPR*, 2022.
- [67] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *IJCV*, 2022.
- [68] Roland S. Zimmermann, Yash Sharma, Steffen Schneider, Matthias Bethge, and Wieland Brendel. Contrastive learning inverts the data generating process. In *ICML*, 2021.

SUPPLEMENTARY MATERIAL:

Waffling around for Performance: Visual Classification with Random Words and Broad Concepts

In this supplementary material, we first provide a collection of additional results in §A which extend those presented in the main paper to more backbone models. Finally, we showcase the GPT-generated descriptors for our additionally used benchmarks beyond [36] (§B), and some exemplary images from the eleven benchmarks used in this work in Fig. 7.

A. Additional results

Motivational experiments for random class descriptors. In Tab. 9, we extend our motivational experiments on random class descriptor assignment to motivate WaffleCLIP from Tab. 1, highlighting similar behaviour on both a larger ViT-L/14 and a ResNet50 backbone network, where descriptor randomization does not result in a significant drop in performance, but rather matching initial DCLIP performance.

Comparison of WaffleCLIP and DCLIP. Tab. 10 extends results from Tab. 2 on the ViT-L/14 and ResNet50 backbones, in which WaffleCLIP as standalone, as well as equipped with high-level concepts and/or joint usage of LLM-generated descriptors, is compared to DCLIP. Results reinforce our conclusions drawn in §4.1, wherein WaffleCLIP, without access to any external LLM, can match the performance of LLM-descriptor-based approaches like DCLIP. In addition to that, we again find complementarity between randomized descriptors and LLM-generated ones, as well as performance gains through the usage of automatic high-level concept selection.

Progression from systematic to fully randomized descriptor scrambling. Finally, Tab. 11 extends the descriptor scrambling progression studies from Tab. 4 to two additional backbones, namely, ViT-L/14 and ResNet50. Similar to the ViT-B/32 backbone, a move from systematic semantic shifts to independently subsamples descriptors can recover and even beat DCLIP performance.

B. Exemplary GPT-3 generated descriptors for additional benchmarks

As we introduce descriptions for three additional datasets compared to [36], we provide four example descriptors for three random classes for each dataset.

Flowers102

Pink Primrose

- "delicate flower"
- "five petals in a star shape"
- "pink in color"
- "often has yellow center"

Balloon Flower

- "a delicate flower with five petals"
- "a unique balloon-like shape"
- "a star-shaped center in the middle of the flower"
- "vibrant colors such as pink, purple, blue, white, and yellow"

Sunflower

- "large, bright yellow petals"
- "a dark center surrounded by disk florets"
- "long stem"
- "a single, long, narrow leaves tapered to a point"

FGVCAircraft

A300

- "black or silver color"
- "a rectangular body with rounded edges"
- "two lens ports"
- "a mode dial"

EMB-120

- "a cabin with 30–33 seats"
- "a distinctive high-wing design"
- "two Pratt and Whitney PW118 turboprop engines"
- "a T-tail configuration"

Tornado

- "dark, rotating funnel-shaped cloud"
- "strong winds"
- "dark clouds"
- "heavy precipitation"

ViT-L/14	ImageNetV2	ImageNet	CUB200	EuroSAT	Places365	Food101	Oxford Pets	DTD	Avg
CLIP [47]	67.90	73.37	62.24	56.03	40.46	92.55	93.30	52.87	67.34
DCLIP [36]	69.72	75.26	63.53	58.72	42.60	92.81	93.89	56.60	69.14
DCLIP (same, 1x)	69.27 ± 0.23	75.05 ± 0.15	64.21 ± 0.36	57.59 ± 1.72	42.01 ± 0.23	93.15 ± 0.13	93.97 ± 0.22	55.16 ± 0.47	68.80 ± 0.66
DCLIP (same, 2x)	69.58 ± 0.21	75.30 ± 0.16	64.30 ± 0.26	59.32 ± 1.63	42.28 ± 0.17	93.31 ± 0.05	94.04 ± 0.11	55.31 ± 0.50	69.18 ± 0.62
ResNet50	ImageNetV2	ImageNet	CUB200	EuroSAT	Places365	Food101	Oxford Pets	DTD	Avg
CLIP [47]	51.34	58.16	45.20	28.09	36.63	78.37	83.76	38.51	52.51
DCLIP [36]	52.70	59.66	47.76	34.27	38.39	78.59	85.77	41.01	54.77
DCLIP (same, 1x)	52.63 ± 0.28	59.69 ± 0.30	47.76 ± 0.39	32.74 ± 1.49	38.63 ± 0.22	80.08 ± 0.58	85.36 ± 0.52	40.77 ± 0.63	54.71 ± 0.67
DCLIP (same, 1x)	52.89 ± 0.23	59.90 ± 0.26	47.70 ± 0.29	34.37 ± 1.27	38.93 ± 0.21	80.11 ± 0.30	85.34 ± 0.29	40.91 ± 0.79	55.02 ± 0.58

Table 9: **Motivating random class descriptors - additional backbones.** Extension of our motivational experiments from Tab. 1 with ViT-L/14 and ResNet50 backbones.

ViT-L/14	ImageNetV2	ImageNet	CUB200	EuroSAT	Places365	Food101	Oxford Pets	DTD	Avg
CLIP [47]	67.90	73.37	62.24	56.03	40.46	92.55	93.30	52.87	67.34
+ Concepts	↓	↓	63.01	61.23	41.07	93.52	93.65	↓	68.32
DCLIP [36]	69.72	75.26	63.53	58.72	42.60	92.81	93.89	56.60	69.14
WaffleCLIP (ours)	69.48 ± 0.08	75.30 ± 0.04	64.18 ± 0.13	61.17 ± 0.35	42.26 ± 0.10	93.31 ± 0.09	91.98 ± 0.11	53.94 ± 0.29	68.95 ± 0.18
+ Concepts	↓	↓	63.40 ± 0.17	60.20 ± 0.87	42.57 ± 0.09	93.65 ± 0.05	94.38 ± 0.08	↓	69.12 ± 0.33
+ GPT descr.	69.80 ± 0.13	75.57 ± 0.06	64.32 ± 0.21	60.63 ± 1.23	42.96 ± 0.12	93.28 ± 0.08	93.35 ± 0.22	56.33 ± 0.42	69.53 ± 0.48
+ GPT descr. + Concepts	↓	↓	63.14 ± 0.16	61.82 ± 1.07	42.95 ± 0.09	93.49 ± 0.04	94.12 ± 0.09	↓	69.65 ± 0.42
ResNet50	ImageNetV2	ImageNet	CUB200	EuroSAT	Places365	Food101	Oxford Pets	DTD	Avg
CLIP [47]	51.34	58.16	45.20	28.09	36.63	78.37	83.76	38.51	52.51
+ Concepts	↓	↓	46.60	34.06	37.43	80.89	83.43	↓	53.80
DCLIP [36]	52.70	59.66	47.76	34.27	38.39	78.59	85.77	41.01	54.77
WaffleCLIP (ours)	52.89 ± 0.15	60.12 ± 0.12	47.68 ± 0.15	31.34 ± 0.47	38.32 ± 0.10	79.68 ± 0.17	84.32 ± 0.20	39.25 ± 0.27	54.20 ± 0.23
+ Concepts	↓	↓	48.34 ± 0.13	35.08 ± 0.42	39.03 ± 0.08	81.38 ± 0.08	85.80 ± 0.12	↓	55.24 ± 0.21
+ GPT descr. + Concepts	↓	↓	48.41 ± 0.21	37.36 ± 0.62	39.43 ± 0.07	81.17 ± 0.09	85.82 ± 0.16	↓	55.75 ± 0.26

Table 10: **Performance of WaffleCLIP with additional backbones.** Extending the comparison of WaffleCLIP (Tab. 2) against GPT-generated fine-grained class-descriptors in DCLIP [36] over both ViT-L/14 and ResNet50 backbones, we find similarly consistent insights, where LLM-free WaffleCLIP can match the performance of DCLIP. Joint usage of both randomized and LLM-generated descriptors again reveal complementarity (*WaffleCLIP + GPT descr.*). In addition to that, the usage of automatically extracted high-level semantic concepts can provide consistent additional performance gains (+ Concepts). We use (↓) to denote same results as previous lines where high-level concept guidance is not applicable.

Stanford Cars

Acura TL Sedan 2012

- "silver, grey, or black exterior"
- "Acura logo on the front grille"
- "distinctive headlights"
- "chrome accents on the exterior"
- "a curved hood"
- "wide, round headlights"
- "a Honda logo"

BMW X6 SUV 2012

- "four-door SUV"
- "sloping roof-line"
- "signature BMW kidney grille"
- "round headlights and taillights"

Honda Odyssey Minivan 2012

- "four doors and a hatchback"

ViT-L/14	ImageNetV2	ImageNet	CUB200	EuroSAT	Places365	Food101	Oxford Pets	DTD	Avg
DCLIP [36]	69.72	75.26	63.53	58.72	42.60	92.81	93.89	56.60	69.14
DCLIP (interchanged)	66.44 ±0.12	72.07 ±0.15	63.62 ±0.44	51.49 ±4.89	37.06 ±0.41	91.30 ±0.30	93.74 ±0.28	49.84 ±0.78	65.69 ±1.77
DCLIP (scrambled)	68.68 ±0.21	74.47 ±0.11	63.78 ±0.13	55.98 ±2.01	41.29 ±0.23	92.29 ±0.20	93.52 ±0.18	53.28 ±1.12	67.91 ±0.83
DCLIP (random, 1x)	68.01 ±0.22	73.89 ±0.08	63.81 ±0.22	55.72 ±2.01	40.32 ±0.29	92.37 ±0.31	93.60 ±0.19	52.83 ±0.46	67.57 ±0.76
DCLIP (random, 5x)	69.27 ±0.17	75.11 ±0.08	64.25 ±0.16	58.34 ±1.55	42.11 ±0.14	93.22 ±0.12	93.88 ±0.09	55.28 ±0.23	68.93 ±0.57

ResNet50	ImageNetV2	ImageNet	CUB200	EuroSAT	Places365	Food101	Oxford Pets	DTD	Avg
DCLIP [36]	52.70	59.66	47.76	34.27	38.39	78.59	85.77	41.01	54.77
DCLIP (interchanged)	49.80 ±0.22	56.35 ±0.06	47.68 ±0.32	28.17 ±4.43	33.77 ±0.34	77.59 ±0.29	84.60 ±0.63	35.81 ±1.12	51.72 ±1.64
DCLIP (scrambled)	52.20 ±0.20	59.21 ±0.06	47.60 ±0.39	34.98 ±2.00	37.90 ±0.18	78.33 ±0.14	85.07 ±0.34	39.19 ±0.95	54.31 ±0.81
DCLIP (random, 1x)	51.60 ±0.29	58.29 ±0.15	47.37 ±0.23	30.18 ±4.18	36.82 ±0.26	78.87 ±0.24	84.52 ±0.17	38.89 ±0.85	53.32 ±1.52
DCLIP (random, 5x)	52.81 ±0.09	59.73 ±0.05	47.74 ±0.10	34.53 ±0.74	38.62 ±0.15	80.20 ±0.13	85.30 ±0.15	40.29 ±0.46	54.90 ±0.32

Table 11: **Progression from systematic to fully randomized descriptor scrambling - additional benchmarks.** We extend our descriptor scrambling progression studies from Tab. 4 to two additional backbones: ViT-L/14 and ResNet50. In both cases, the same trend can be seen, in which a move from systematic semantic shift to independently subsampled descriptors can recover DCLIP performance after an initial performance drop.

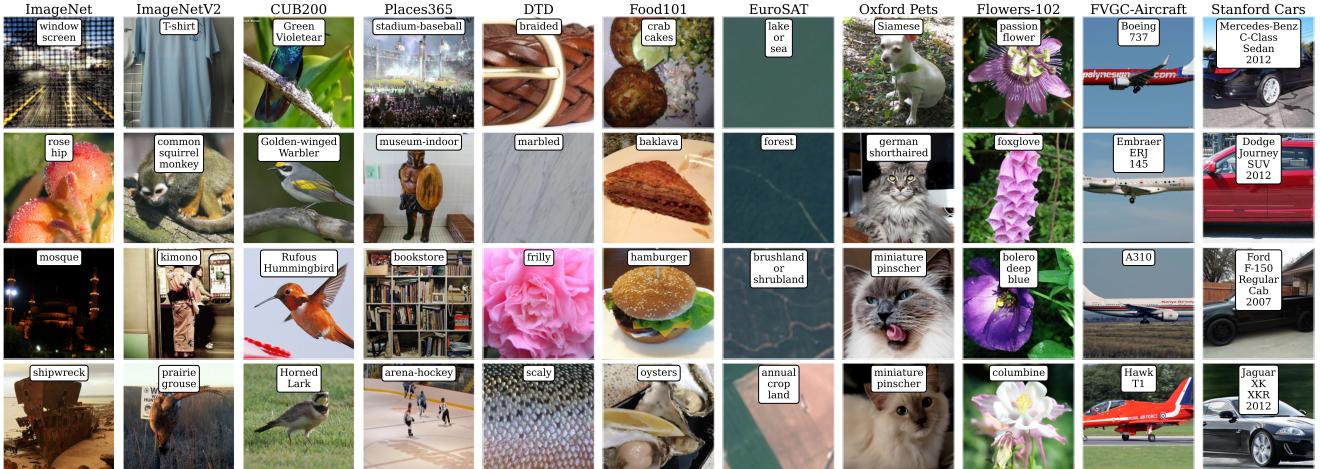


Figure 7: To get an intuition of the different visual classification tasks, we showcase samples of four randomly selected classes for each of the eleven utilized visual classification benchmarks.