# Zero-Shot Robustification of Zero-Shot Models With Auxiliary Foundation Models

Anonymous Author(s) Affiliation Address email

#### Abstract

Zero-shot inference is a powerful paradigm that enables the use of large pretrained 1 models for downstream classification tasks without further training. However, 2 these models are vulnerable to inherited biases that can impact their performance. 3 The traditional solution is fine-tuning, but this undermines the key advantage of 4 pretrained models, which is their ability to be used out-of-the-box. We propose 5 ROBOSHOT, a method that improves the robustness of pretrained model embed-6 dings in a fully zero-shot fashion. First, we use zero-shot language models (LMs) 7 to obtain useful insights from task descriptions. These insights are embedded and 8 used to remove harmful and boost useful components in embeddings-without any 9 supervision. Theoretically, we provide a simple and tractable model for biases in 10 zero-shot embeddings and give a result characterizing under what conditions our 11 approach can boost performance. Empirically, we evaluate ROBOSHOT on nine 12 image and NLP classification tasks and show an average improvement of 15.98% 13 over several zero-shot baselines. Additionally, we demonstrate that ROBOSHOT is 14 15 compatible with a variety of pretrained and language models.

### 16 **1 Introduction**

I7 Zero-shot models are among the most exciting paradigms in machine learning. These models obviate the need for data collection and model training loops by simply asking the model for a prediction on any set of classes. Unfortunately, such models inherit biases or undesirable correlations from their large-scale training data [DLS<sup>+</sup>18, TE11]. In a now-canonical example [KSM<sup>+</sup>21], they often associate waterbirds with water background. This behavior leads to decreased performance, often exacerbated on rare data slices that break in-distribution correlations.

A growing body of literature [YNPM23, GKG<sup>+</sup>22, ZR22] seeks to improve robustness in zero-shot models. While promising, these works require labeled data to train or fine-tune models, and so **do not tackle the zero-shot setting.** A parallel line of research seeking to debias word embeddings [AZS<sup>+</sup>, BCZ<sup>+</sup>16, DP19, LGPV20] often sidesteps the need for labeled data. Unfortunately, these works often require domain expertise and painstaking manual specification in order to identify particular concepts that embeddings must be invariant to. As a result, out-of-the-box word embedding debiasing methods also cannot be applied to zero-shot robustification.

Can we robustify zero-shot models without (i) labeled data, (ii) training or fine-tuning, or (iii) manual identification? Surprisingly, despite this seemingly impoverished setting, it is often possible to do so. Our key observation is that zero-shot models **contain actionable insights** that can be exploited to improve themselves or other zero-shot models. These insights are noisy but cheaply available at scale—and can be easily translated into means of refinement for zero-shot representations. These refinements improve performance, particularly on underperforming slices—at nearly no cost.



Figure 1: ROBOSHOT pipeline (right) vs. vanilla zero-shot classification (left).

We propose ROBOSHOT, a system that robustifies zero-shot models via auxiliary language models *without labels, training, or manual specification.* Using just the task description, ROBOSHOT obtains *positive and negative insights* from a language model (potentially the model to be robustified itself).
It uses embeddings of these noisy insights to recover *harmful, beneficial*, and *benign* subspaces of
zero-shot latent representation spaces. Representations are then modified to neutralize and emphasize
their harmful and beneficial components, respectively.

Theoretically, we introduce a simple and tractable model to capture and quantify failures in zero-shot models. We provide a result that characterizes the *quantity and quality* of insights that must be obtained as a function of the severity of harmful correlations. Empirically, ROBOSHOT achieves 15.98% improvement across nine image and NLP datasets while offering sufficient versatility to apply to a diverse variety of base models. Most excitingly, in certain cases, it reaches comparable or greater improvements even when compared to fine-tuned models that rely on labeled data.

- 48 Our contributions include,
- A simple theoretical model describing zero-shot model failures along with a theoretical
   analysis of our approach that characterizes the amount of information required for obtaining
   improvements as a function of the most harmful unwanted correlation,
- ROBOSHOT, an algorithm that implements our core idea. It extracts insights from foundation models and uses them to improve zero-shot representations,
- Extensive experimental evidence on zero-shot language and multimodal models, showing
   improved worst-group accuracy of 15.98% across nine image and NLP datasets.

## 56 2 Related Work

We describe related work in zero-shot model robustness, debiasing embeddings, guiding multi-modal
 models using language, and using LMs as prior information.

Zero-Shot inference robustness. Improving model robustness to unwanted correlations is heav-59 ily studied [SKHL19, ABGLP19, KCJ<sup>+</sup>21, KIW22, LHC<sup>+</sup>21, LCT<sup>+</sup>22]. Some methods require 60 training from scratch and are less practical when applied to large pretrained architectures. Existing 61 approaches to improve robustness post-pretraining predominantly focus on fine-tuning. [YNPM23] 62 detects spurious attribute descriptions and fine-tunes using these descriptions. Specialized contrastive 63 loss is used to fine-tune a pretrained architecture in [GKG<sup>+</sup>22] and to train an adapter on the frozen 64 embeddings in [ZR22]. While promising, fine-tuning recreates traditional machine learning pipelines 65 (e.g., labeling, training, etc.), which contradicts the promise of zero-shot models. In contrast, our 66 goal is to avoid any training and any use of labeled data. 67

Debiasing embeddings. A parallel line of work seeks to de-bias text embeddings [AZS<sup>+</sup>]
 [BCZ<sup>+</sup>16] [DP19] [LGPV20] and multimodal embeddings [WZS22, BHB<sup>+</sup>22, WLW21] by re-



Figure 2: (a) ROBOSHOT debiases original input embedding (left). The projected embedding (right)'s variance in the unwanted direction is reduced, and in the relevant direction increases. (b) Embedding projection. We project embeddings to the space orthogonal to the embeddings of all unwanted insights (e.g., water and land)

moving subspaces that contain harmful or unwanted concepts. We use a similar procedure as a 70 building block. However, these methods either target specific fixed concepts (such as gender) or rely 71 on concept annotations, which limits their applicability across a wide range of tasks. In contrast, our 72 method automates getting both beneficial and unwanted concepts solely from the task descriptions. 73 An additional difference is that our goal is simply to add robustness at low or zero-cost; we not seek 74 to produce fully-invariant representations as is often desired for word embeddings. 75

**Using language to improve visual tasks** A large body of work has shown the efficacy of using 76 language to improve performance on vision tasks [RKH<sup>+</sup>21, FCS<sup>+</sup>13, LCLBC20]. Most relevant 77 are those that focus on robustness, like [PDN<sup>+</sup>22], where attention maps using multimodal models 78 (like CLIP) are used as extra supervision to train a downstream image classifier. [YNPM23] uses 79 text descriptions of spurious attributes in a fine-tuning loss to improve robustness against spurious 80 correlations. In contrast to these works, we focus on using textual concepts to improve zero-shot 81 82 model robustness-without fine-tuning.

**Language model as prior** The basis of our work comes from the observation that language models 83 contain information that can serve as a prior for other learning tasks. [KNST23] finds that LLMs can 84 perform causal reasoning tasks, substantially outperforming existing methods. [CCSE22] explicitly 85 prompts LLMs for task-specific priors, leading to substantial performance improvements in feature 86 selection, reinforcement learning, and causal discovery. Our work shares the spirit of these approaches 87 in using the insights embedded in language models to enhance zero-shot robustness. 88

#### 3 **RoboShot: Robustifying Zero-shot Models** 89

We are ready to provide our setup and describe the algorithm. 90

#### 3.1 Modeling and setup 91

Suppose that the zero-shot model's latent space contains an (unknown) concept set; similar notions 92

have been studied frequently in the literature [DKA<sup>+</sup>]. For simplicity, we assume that this concept 93

set is given by the orthonormal vectors  $\{z_1, \ldots, z_k\}$ . The model's encoder produces, for a particular input a representation x that is a mixture of concepts  $\sum_i \gamma_i z_i$ , where  $\gamma_i \ge 0$  are weights. 94

95

We shall work with the following theoretical model for zero-shot classification. It closely resembles 96 models like CLIP. For simplicity, we assume that there are two classes. It is straightforward to extend 97

### Algorithm 1: ROBOSHOT

- 1: **Parameters:** Input embedding x, class embeddings  $c^0$ ,  $c^1$ , harmful insight representations  $v^1, \ldots, v^{|S|}$ , helpful insight representations  $u^1, \ldots, u^{|R|}$ 2: for  $j \in \{1, 2, \ldots, |S|\}$  do
- Reject harmful insight  $v_j$ : set  $x \leftarrow x \langle x, v^j \rangle / \langle v^j, v^j \rangle v^j$ 3:
- 4: Renormalize x = x/||x||
- 5: end for
- 6: for  $k \in \{1, 2, \dots, |R|\}$  do
- Increase helpful insight  $u_k$ : set  $x \leftarrow x + \langle x, u^k \rangle / \langle u^k, u^k \rangle u^k$ 7:
- 8: end for
- 9:  $\hat{c} = \mathbb{1}\{x^T c^0 < x^T c^1\}$
- 10: **Returns:** Robustified zero-shot prediction  $\hat{c}$

the analysis below to multiple classes. We take  $\sum_i \alpha_i z_i$  to be the embedding of a datapoint, while  $c^0 = \sum_i \beta_{i,0} z_i$  is the embedding of the first class and  $c^1 = \sum_i \beta_{i,1} z_i$  is that of the second. Finally, we assume that we have access to m answers  $v^1, \ldots, v^m$  from the queries to the language model. These are given by  $v^j = \sum_i \gamma_{i,j} z_i$  for  $j \leq m$ . We call these *insight representations*. Without our approach, the prediction is made by  $\mathbb{1}\{(\sum_i \alpha_i z_i)^T (\sum_i \beta_{i,0} z_i) < (\sum_i \alpha_i z_i)^T (\sum_i \beta_{i,1} z_i)\}$ , so that we predict whichever class has higher inner product with the datapoint's embedding. 98 99 100 101 102 103

Next, we assume that each input representation x can be represented by partitioning the mixture 104 components into three groups, 105

$$x = \sum_{s}^{S} \alpha_{s}^{\text{harmful}} z_{s} + \sum_{r}^{R} \alpha_{r}^{\text{helpful}} z_{r} + \sum_{b}^{B} \alpha_{b}^{\text{benign}} z_{b}.$$

The same holds for class and insight representations. 106

**Example** We illustrate how harmful correlations produce errors on rare slices of data through a 107 standard task setting, Waterbirds [KSM<sup>+</sup>21]. In this dataset, the goal is to classify landbirds versus 108 waterbirds, and the background (land or water) is spurious. Suppose that we have these terms 109 relate to concepts such that  $z_{water} = -z_{land}$  and  $z_{waterbird} = -z_{landbird}$ . 110

Consider a datapoint coming from a rare slice infrequently encountered in the training set. This might 111 be an image of a landbird over water. Its embedding might be  $x = 0.7 z_{water} + 0.3 z_{landbird}$ . We may 112 also have that 113

 $c_{\text{waterbird}} = 0.4 z_{\text{water}} + 0.6 z_{\text{waterbird}}$  and  $c_{\text{landbird}} = 0.4 z_{\text{land}} + 0.6 z_{\text{landbird}}$ .

Then,  $x^T c_{waterbird} = 0.1 > x^T c_{landbird} = -0.1$ , so that the prediction is waterbird, and thus 114 incorrect. This is caused by the presence of harmful components in both the class embedding (caused 115 by seeing too many images with water described as waterbirds) and the datapoint embedding (where 116 the water background appears). Thus our goal is to *remove* harmful components (the  $z_s$ 's) and *boost* 117 helpful components (the  $z_r$ 's). We explain our approach towards doing so next. 118

#### 3.2 ROBOSHOT: Zeroshot robustification with LLM 119

We describe ROBOSHOT in Algorithm 1. It uses representations of insights from language models to 120 121 shape input and class embeddings to remove harmful components and boost helpful ones. Figure 2 is helpful in understanding the intuition behind these procedures. The left part (a) illustrates the 122 effect of ROBOSHOT on a true dataset. Note how unhelpful directions are neutralized while others 123 are boosted. The illustration on the right (b) shows this effect on the waterbirds running example. 124

**Obtaining insight representations from LMs** The first question is how to obtain insight repre-125 sentations without training. To do so in a zero-shot way, we use *textual* descriptions of harmful and 126 helpful concepts by querying language models using only the task description. For example, in the 127 Waterbirds dataset, we use the prompt "What are the biased/spurious differences between waterbirds 128 and landbirds?". We list the details of the prompts used in the Appendix. Let  $s_1, s_2$  be the text 129 insights obtained from the answer (e.g., {'water background,' 'land background'}). We obtain 130 a spurious insight representation by taking the difference of their embedding  $v = \frac{g(s_1) - g(s_2)}{\|g(s_1) - g(s_2)\|}$ 131 where q is the text encoder of our model. 132

In addition to attempting to discover harmful correlations, we seek to discover helpful components 133 in order to boost their magnitudes past remaining harmful ones (or noise). The procedure is similar. 134 We obtain insight representations using language models. For example, we ask "What are the true 135 characteristics of waterbirds and landbirds?' and obtain e.g., {'short beak,' 'long beak'}. The 136 remainder of the procedure is identical to the case of harmful components. Note that since we 137 are seeking to boost (rather than remove) components, it is also possible to fix a multiplicative 138 constant (to be treated as a hyperparameter) for the boosting procedure. That is, we could take 139  $x \leftarrow x + \nu \times \langle x, u^k \rangle / \langle u^k, u^k \rangle u^k$  for some  $\nu > 0$ . While this is possible if we have access to a 140 labeled set that we can tune  $\nu$  over, we intentionally avoid doing so to ensure our procedure is truly 141 zero-shot. 142

Prompting a language model is typically inexpensive, which will enable obtaining multiple insight vectors  $\tilde{v}^1, \ldots, \tilde{v}^m$ . From these, we obtain an orthogonal basis  $v^1, \ldots, v^m$  separately for harmful and helpful components. Thus we have access to recovered subspaces spanned by such components.

Removing and Boosting Components ROBOSHOT applies simple vector rejection to mitigate or
 remove harmful components, which is described in lines 2-5 of Algorithm 1. Similarly, it boosts
 helpful components as described in lines 6-9.

To see the impact of doing so, consider our earlier example. Suppose that  $v^{\text{harmful}} = 0.9z_{\text{water}} + 0.1z_{\text{landbird}}$ , and that this is our only harmful insight. Similarly, suppose that we obtain a single helpful insight given by  $v^{\text{helpful}} = 0.1z_{\text{water}} + 0.9z_{\text{landbird}}$ . Note that even these insights can be imperfect: they do not uniquely identify what are harmful or helpful concepts, as they have non-zero weights on other components.

We first obtain from removing the harmful component (ignoring normalization for ease of calculation) that

$$\hat{x} \leftarrow x - \frac{\langle x, v^{\text{harmful}} \rangle}{\langle v^{\text{harmful}}, v^{\text{harmful}} \rangle} v^{\text{harmful}} = -0.0244 z_{\text{water}} + 0.2195 z_{\text{landbird}}.$$

Then, we already we have that  $x^T c_{waterbird} = -0.1415 < x^T c_{landbird} = 0.1415$ , so that the correct class is obtained. In other words we have already, from having access to a single insight, neutralized a harmful correlation and corrected what had been an error. Adding in the helpful component further helps. We obtain

$$\hat{x} \leftarrow \hat{x} + \frac{\langle \hat{x}, v^{\text{neipful}} \rangle}{\langle v^{\text{helpful}}, v^{\text{helpful}} \rangle} v^{\text{helpful}} = -0.0006 z_{\text{water}} + 0.4337 z_{\text{landbird}}.$$

This further increases our margin. Note that it is not necessary to fully neutralize (i.e., to be fully invariant to) spurious or harmful components in our embeddings. The only goal is to ensure, as much as possible, that their magnitudes are reduced when compared to helpful components (and to benign components). In the following section, we provide a theoretical model for the magnitudes of such components and characterize the conditions under which it will be possible to correct zero-shot errors. We note that there is a variant of our approach that can also update class embeddings as well.

#### 166 4 Analysis

Next, we provide an analysis that characterizes under what conditions ROBOSHOT is capable of correcting zero-shot errors. First, we consider the following error model on the weights of the various representations. For all benign representations, we assume that  $\alpha_b, \beta_b, \gamma_b \sim \mathcal{N}(0, \sigma_{\text{benign}}^2)$ . That is, the magnitudes of benign components are drawn from a Gaussian distribution. The value of  $\sigma_{\text{benign}}$  is a function of the amount of data and the training procedure for the zero-shot model.

Next, we assume that the embedding insight  $v_s = \sum_{i=1}^k \gamma_{i,s} z_i$  (where  $1 \le s \le S$ ) satisfies the property that for  $i \ne s$ ,  $\gamma_{i,s} \sim \mathcal{N}(0, \sigma_{\text{insight}}^2)$ , while  $\gamma_{s,s}$  is a constant. In other words, the vectors  $v_1, \ldots, v_S$  spanning the harmful component subspace are well-aligned with genuinely harmful concepts, but are also affected by noise. We seek to understand the interplay between this noise, benign noise, and the coefficients of the other vectors (i.e., helpful components). Let the result of rejecting embedding insights  $v_1, \ldots, v_S$  be

$$\hat{x} = x - \sum_{s=1}^{S} \frac{x^T v_s}{||v_s||^2} v_s = \sum_i A_i z_i.$$

- <sup>178</sup> We provide a bound on  $A_s$ , the coefficient of a targeted harmful concept post-removal.
- **Theorem 4.1.** Under the noise model described above, the post-removal coefficient for harmful concept s satisfies

$$\left|\mathbb{E}\left[A_{s}\right]\right| \leq \left|\frac{(k-1)\alpha_{s}\sigma_{insight}^{2}}{\gamma_{s,s}^{2}}\right| + \left|\sum_{t\neq s}^{S}\frac{\alpha_{s}\sigma_{insight}^{2}}{\gamma_{t,t}^{2}}\right|$$

181 where k is the number of concepts.

The theorem illustrates how and when the rejection component of ROBOSHOT works-it scales 182 down harmful coefficients at a rate inversely proportional to the harmful coefficients of the insight 183 embeddings. As we would hope, when insight embeddings have larger coefficients for harmful vectors 184 (i.e., are more precise in specifying terms that are not useful), ROBOSHOT yields better outcomes. 185 In addition, we observe that the harmful coefficients decrease when the insight embeddings have 186 less noise. In fact, we have that  $\lim_{\sigma_{insight} \to 0} A_s = 0$  — the case of perfectly identifying harmful 187 concepts. In the Appendix, we present additional theoretical results for control of helpful coefficients 188 along with a combined result. 189

#### **190 5 Experimental Results**

- 191 This section evaluates the following claims about ROBOSHOT:
- Improving multi-modal models (Section 5.1): ROBOSHOT improves zero-shot classification robustness of various multi-modal models, even outperforming prompting techniques that include spurious insight descriptions (which we do not have access to) in the label prompts.
- Improving language models (Section 5.2): ROBOSHOT improves zero-shot robustness when
   using language model embeddings for text zero-shot classification.
- Extracting concepts from LM with varying capacities (Section 5.3): ROBOSHOT can extract
   insights from language models with varying capacities. Improvements persist with weaker LMs.
- **Ablations (Section 5.4)** ROBOSHOT benefits from both removing harmful and boosting helpful representations (line 3 and line 7 in ROBOSHOT Algorithm 1).

Metrics and how to interpret the results. We use three metrics: average accuracy % (AVG), worst-group accuracy % (WG), and the gap between the two (Gap). While a model that relies on harmful correlations may achieve high AVG when such correlations are present in the majority of the test data, it may fail in settings where the correlation is absent. A robust model should have high AVG and WG, with a small gap between them.

**Baselines** We compare against the following sets of baselines:

1. Multimodal baselines: We compare against: (i) vanilla zero-shot classification (ZS) and (ii) 207 zero-shot classification with group information (Group Prompt ZS). We do so across a variety of 208 models: CLIP (ViT-B-32 and ViT-L-14) [RKH<sup>+</sup>21], ALIGN [JYX<sup>+</sup>21], and AltCLIP [CLZ<sup>+</sup>22]. 209 Group Prompt ZS assumes access to spurious or harmful insight annotations and includes them 210 in the label prompt. For instance, the label prompts for waterbirds dataset become [waterbird 211 with water background, waterbird with land background, landbird with water 212 background, landbird with land background]. We only report Group Prompt ZS results 213 on datasets where spurious insight annotations are available. 214

 Language model baselines: We compare against zero-shot classification using multiple language model embeddings, including BERT [RG19] and Ada [NXP<sup>+</sup>22] (ZS).

### 217 5.1 Improving multi-modal models

218

Setup. We experimented on five binary and multi-class datasets with spurious correlations and distribution shifts, coming from a variety of domains: Waterbirds [SKHL19], CelebA [LLWT15],
 CXR14 [WPL<sup>+</sup>17], PACS [LYSH17], and VLCS [FXR13]. We use the default test splits of all datasets. Dataset details are provided in the appendix. For CXR14, we use BiomedCLIP [ZXU<sup>+</sup>23],

Dataset	Model		ZS		Gr	oupProm	pt ZS	]	<b>R</b> OBOSHOT			
		AVG	WG(↑)	Gap(↓)	AVG	WG(↑)	Gap(↓)	AVG	WG(↑)	Gap(↓)		
	CLIP (ViT-B-32)	80.7	27.9	52.8	81.6	<u>43.5</u>	<u>38.1</u>	82.0	54.4	28.6		
Waterbirds	CLIP (ViT-L-14)	88.7	<u>27.3</u>	61.4	70.7	10.4	<u>60.3</u>	79.9	45.2	34.7		
	ALIGN	72.0	50.3	21.7	72.5	5.8	66.7	50.9	41.0	9.9		
	AltCLIP	90.1	<u>35.8</u>	54.3	82.4	29.4	<u>53.0</u>	78.5	54.8	23.7		
	CLIP (ViT-B-32)	80.1	72.7	7.4	80.4	<u>74.9</u>	<u>5.5</u>	84.8	80.5	4.3		
CelebA	CLIP (ViT-L-14)	80.6	<u>74.3</u>	<u>6.3</u>	77.9	68.9	9.0	85.5	82.6	2.9		
	ALIGN	81.8	77.2	4.6	78.3	67.4	10.9	86.3	83.4	2.9		
	AltCLIP	82.3	79.7	2.6	82.3	79.0	3.3	86.0	77.2	8.8		
	CLIP (ViT-B-32)	96.7	82.1	<u>14.6</u>	97.9	82.7	15.2	97.0	86.3	10.7		
PACS	CLIP (ViT-L-14)	98.1	79.8	18.3	98.2	86.6	11.6	98.1	83.9	14.2		
	ALIGN	95.8	77.1	18.7	96.5	65.0	31.5	95.0	<u>73.8</u>	21.2		
	AltCLIP	98.5	82.6	15.9	98.6	<u>85.4</u>	13.2	98.7	89.5	9.2		
	CLIP (ViT-B-32)	75.6	20.5	55.1		-		76.5	33.0	43.5		
VLCS	CLIP (ViT-L-14)	72.6	4.20	68.4		-		71.1	12.6	58.5		
	ALIGN	78.8	33.0	45.8		-		77.6	39.8	37.8		
	AltCLIP	78.3	24.7	53.6		-		78.9	25.0	53.9		
CXR14	BiomedCLIP	55.3	28.9	26.4		-		56.2	41.6	14.6		

Table 1: Main results. Best WG and Gap performance **bolded**, second best underlined.



Figure 3: (a) Original (green) and projected (red) input embeddings x, and label embeddings  $c^0$  and  $c^1$ . (b) label embeddings  $c^0$  and  $c^1$ , harmful insight embeddings  $v^k$  (black star) and helpful insight embeddings  $u^j$  (blue star)

which is a variant of CLIP finetuned on biomedical images and articles. All experiments are conducted using frozen pretrained models.

**Results.** Table 1 shows that **ROBOSHOT significantly improves the worst group performance** (WG) and maintains (and sometimes also improves) the overall average (AVG) without any auxiliary

227 information (in contrast to Group Prompt, which requires access to spurious insight annotation).

Improved robustness nearly across-the-board suggests that both the insights extracted from LMs and 228 the representation modifications are useful. We also provide insights insights into the case where 229 our method does not improve the baseline (ALIGN model on Waterbirds) in Fig. 3. In Fig. 3a, we 230 visualize the original and projected input embeddings (x in green and red points, respectively), and 231 the label embeddings ( $c^0$  and  $c^1$ ). Fig. 3a (left) shows the embeddings from the ALIGN model. We 232 observe that the projected embeddings (red) still lie within the original embedding space, even with 233 reduced variance. In contrast, when examining the CLIP model embeddings (Figure 3a (right)), we 234 observe that the projected embeddings are significantly distant from the original ones. Unsurprisingly, 235 Figure 3b (left) reveals that  $v^{j}$  and  $u^{k}$  (harmful and helpful insight embeddings in black and blue 236 stars, respectively) are not distinguishable in the text embedding space of ALIGN, collapsing the 237 input embeddings after ROBOSHOT is applied. 238

Dataset	Model		ZS			RoboShot			
		AVG	$WG(\uparrow)$	$Gap(\downarrow)$	AVG	$WG(\uparrow)$	$Gap(\downarrow)$		
CivilComments	BERT	48.1	33.3	14.8	49.7	42.3	7.4		
	Ada	56.2	43.2	13.0	56.6	44.9	11.7		
HateXplain	BERT	60.4	0.0	60.4	57.3	14.0	43.3		
	Ada	62.8	14.3	48.5	63.6	21.1	42.5		
Amazon	BERT	81.1	64.2	16.8	81.0	64.4	<b>16.6</b>		
	Ada	81.2	63.4	<b>17.8</b>	82.9	63.8	19.1		
Gender Bias	BERT	84.8	83.7	1.1	85.1	84.9	<b>0.2</b>		
	Ada	77.9	60.0	17.9	78.0	60.1	17.9		

Table 2: ROBOSHOT text zero-shot classification. Best WG in **bold**.

Table 3: ROBOSHOT with LMs of varying capacity. Best WG bolded, second best underlined

Dataset	ZS		Ours (ChatGPT)		Ours (Flan-T5)		Ours (GPT2)		Ours (LLaMA)	
	AVG	WG	AVG	WG	AVG	WG	AVG	WG	AVG	WG
Waterbirds	80.7	27.9	82.0	54.4	72.1	32.4	88.0	<u>39.9</u>	84.8	36.5
CelebA	80.1	72.7	84.8	<u>80.5</u>	77.5	68.2	80.3	74.1	84.2	82.0
PACS	96.7	82.1	97.0	86.3	96.2	80.3	97.2	74.0	94.8	71.9
VLCS	75.6	20.5	76.5	33.0	69.6	20.5	75.5	<u>26.1</u>	72.0	18.2

#### 239 5.2 Improving language models

Setup. We experimented on four text classification datasets: CivilComments-WILDS [BDS+19, 240 KSM<sup>+</sup>21], HateXplain [MSY<sup>+</sup>21], Amazon-WILDS [NLM19, KSM<sup>+</sup>21] and Gender Bias clas-241 sification dataset  $DFW^+20$ , MFB<sup>+</sup>17]. We use the default test splits of all datasets. In text 242 experiments, the distinctions between harmful and helpful insights are less clear than for images. 243 For this reason, we only use harmful vector rejection (line 3 in ROBOSHOT) in text experiments. 244 CivilComments and HateXplain are toxic classification datasets with unwanted correlation between 245 toxicity labels and mentions of demographics (e.g., male, female, mentions of religions). The datasets 246 are annotated with demographic mentions of each text, and we directly use them to construct  $v^{j}$ . 247 For Amazon and Gender Bias datasets, we query LMs with task descriptions. All experiments are 248 conducted using frozen pretrained models. 249

**Results.** Table 2 shows that ROBOSHOT also improves zero-shot text classification in text datasets, as shown by our consistent boost over the baselines across all datasets.

#### 252 5.3 Extracting concepts from LMs with varying capacities

Setup. We use LMs with different capacities: ChatGPT [OWJ<sup>+</sup>22], Flan-T5 [CHL<sup>+</sup>22], GPT2 [RWC<sup>+</sup>19], and LLaMA [TLI<sup>+</sup>23], to get harmful and helpful features insights ( $v^j$  and  $u^k$ ).

**Results.** Table 3 shows that ROBOSHOT can get insights on  $v^j$  and  $u^k$  from LMs of various capacities

and improves zero-shot performance. Even though the the LM capacity correlates with the zero-shot performance, ROBOSHOT with weaker LMs still outperforms zero-shot (ZS) baseline.

#### 258 5.4 Ablations

Setup. We run ROBOSHOT with only harmful component mitigation (reject  $v^j$ : ROBOSHOT line 3), only boosting helpful vectors (increase  $u^k$ : ROBOSHOT line 7), and both.

**Results.** The combination of both projections often achieves the best performance, as shown in Table 4. Figure 4 provides insights into the impact of each projection. Rejecting  $v^j$  reduces variance in one direction, while increasing  $u^k$  amplifies variance in the orthogonal direction. When both projections are applied, they create a balanced mixture. We note that when doing both projections does not

Dataset	Model		ZS		Ou	rs ( $v^j$	only)	Ours $(u^k \text{ only})$			Ours (both)		
		AVG	WG(†)	Gap(↓)	AVG	WG(†)	Gap(↓)	AVG	WG(†)	Gap(↓)	AVG	WG(†)	$\overline{\text{Gap}(\downarrow)}$
	CLIP (ViT-B-32)	80.7	27.9	52.8	82.0	<u>50.4</u>	<u>31.6</u>	82.6	30.2	52.4	83.0	54.4	28.6
Waterbirds	sCLIP (ViT-L-14)	88.7	27.3	61.4	82.7	<u>35.8</u>	<u>46.9</u>	88.3	29.8	58.5	79.9	45.2	<b>34</b> .7
	ALIGN	72.0	<u>50.3</u>	21.7	56.4	41.6	14.8	62.8	56.4	6.4	50.9	41.0	<u>9.9</u>
	AltCLIP	90.1	35.8	54.3	81.4	59.0	22.4	89.1	35.2	53.9	78.5	<u>54.8</u>	23.7
	CLIP (ViT-B-32)	80.1	72.7	7.4	85.2	81.5	3.7	79.6	71.3	8.3	84.8	80.5	4.3
CelebA	CLIP (ViT-L-14)	80.6	74.3	6.3	85.9	82.8	<u>3.1</u>	80.0	73.1	6.9	85.5	<u>82.6</u>	2.9
	ALIGN	81.8	77.2	4.6	83.9	78.0	5.7	83.9	<u>81.4</u>	2.5	86.3	83.4	2.9
	AltCLIP	82.3	<b>79.</b> 7	2.6	86.1	75.6	10.5	81.9	<u>79.0</u>	<u>2.9</u>	86.0	77.2	8.8
	CLIP (ViT-B-32)	96.7	82.1	14.6	97.0	83.7	13.3	96.6	84.2	12.4	97.0	86.3	10.7
PACS	CLIP (ViT-L-14)	98.1	79.8	18.3	98.0	79.8	18.2	98.1	<u>83.8</u>	<u>14.3</u>	98.1	83.9	14.2
	ALIGN	95.8	77.1	18.7	95.8	78.0	17.8	95.1	71.1	24.0	95.0	73.8	21.2
	AltCLIP	98.5	82.6	15.9	98.4	83.0	15.4	98.6	<u>88.8</u>	<u>9.8</u>	98.7	89.5	9.2
	CLIP (ViT-B-32)	75.6	20.5	55.1	75.6	22.7	52.9	76.4	29.5	46.9	76.5	33.0	43.5
VLCS	CLIP (ViT-L-14)	72.6	4.2	68.4	70.9	6.8	<u>64.1</u>	73.4	8.9	64.5	71.1	12.6	58.5
	ALIGN	78.8	33.0	45.8	78.2	30.7	47.5	78.0	43.2	34.8	77.6	<u>39.8</u>	37.8
	AltCLIP	78.3	<u>24.7</u>	53.6	77.5	24.4	<u>53.1</u>	79.0	20.5	58.5	78.9	25.0	53.9
CXR14	BiomedCLIP	55.3	28.9	26.4	55.7	41.8	13.9	54.8	21.8	33.0	56.2	41.6	14.6

Table 4: Main results. Best WG and Gap performance **bolded**, second best <u>underlined</u>.



Figure 4: The effect of  $v^{j}$  (reject),  $u^{j}$  (increase), and both projections

improve the baseline, using only  $u^k$  or  $v^j$  still outperforms the baseline. For instance, the ALIGN model in the Waterbirds dataset achieves the best performance with only  $u^k$  projection. This suggests that in certain cases, harmful and helpful concepts are intertwined in the embedding space, and using just one projection can be beneficial. We leave further investigation to future work.

### 269 6 Conclusion

We introduced ROBOSHOT, a fine-tuning-free system that robustifies zero-shot pretrained models in a truly zero-shot way. Theoretically, we characterized the quantities required to obtain improvements over vanilla zero-shot classification. Empirically, we found that ROBOSHOT improves both multimodal and language model zero-shot performance, has sufficient versatility to apply to various base models, and can use insights from less powerful language models.

### 275 **References**

[ABGLP19] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk
 minimization. arXiv preprint arXiv:1907.02893, 2019.

[AZS<sup>+</sup>] Prince Osei Aboagye, Yan Zheng, Jack Shunn, Chin-Chia Michael Yeh, Junpeng
 Wang, Zhongfang Zhuang, Huiyuan Chen, Liang Wang, Wei Zhang, and Jeff Phillips.

280 281		Interpretable debiasing of vectorized language representations with iterative orthogo- nalization. In <i>The Eleventh International Conference on Learning Representations</i> .
282 283 284 285 286	[BCZ <sup>+</sup> 16]	Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, <i>Advances in Neural Information Processing Systems</i> , volume 29. Curran Associates, Inc., 2016.
287 288 289	[BDS <sup>+</sup> 19]	Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Nuanced metrics for measuring unintended bias with real data for text classification. In <i>Companion proceedings of the 2019 world wide web conference</i> , pages 491–500, 2019.
290 291 292 293	[BHB <sup>+</sup> 22]	Hugo Berg, Siobhan Mackenzie Hall, Yash Bhalgat, Wonsuk Yang, Hannah Rose Kirk, Aleksandar Shtedritski, and Max Bain. A prompt array keeps the bias away: Debiasing vision-language models with adversarial learning. <i>arXiv preprint arXiv:2203.11933</i> , 2022.
294 295	[CCSE22]	Kristy Choi, Chris Cundy, Sanjari Srivastava, and Stefano Ermon. Lmpriors: Pre-trained language models as task-specific priors. <i>arXiv preprint arXiv:2210.12530</i> , 2022.
296 297 298	[CHL <sup>+</sup> 22]	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. <i>arXiv preprint arXiv:2210.11416</i> , 2022.
299 300 301	[CLZ <sup>+</sup> 22]	Zhongzhi Chen, Guang Liu, Bo-Wen Zhang, Fulong Ye, Qinghong Yang, and Ledell Wu. Altclip: Altering the language encoder in clip for extended language capabilities. <i>arXiv preprint arXiv:2211.06679</i> , 2022.
302 303 304 305	[DFW <sup>+</sup> 20]	Emily Dinan, Angela Fan, Ledell Wu, Jason Weston, Douwe Kiela, and Adina Williams. Multi-dimensional gender bias classification. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 314–331, Online, November 2020. Association for Computational Linguistics.
306 307 308	[DKA <sup>+</sup> ]	Fahim Dalvi, Abdul Rafae Khan, Firoj Alam, Nadir Durrani, Jia Xu, and Hassan Sajjad. Discovering latent concepts learned in bert. In <i>International Conference on Learning Representations</i> .
309 310	[DLS <sup>+</sup> 18]	Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Measuring and mitigating unintended bias in text classification. 2018.
311 312 313	[DP19]	Sunipa Dev and Jeff Phillips. Attenuating bias in word vectors. In <i>The 22nd Inter-</i> <i>national Conference on Artificial Intelligence and Statistics</i> , pages 879–887. PMLR, 2019.
314 315 316	[FCS <sup>+</sup> 13]	Andrea Frome, Greg S Corrado, Jon Shlens, Samy Bengio, Jeff Dean, Marc'Aurelio Ranzato, and Tomas Mikolov. Devise: A deep visual-semantic embedding model. <i>Advances in neural information processing systems</i> , 26, 2013.
317 318 319	[FXR13]	Chen Fang, Ye Xu, and Daniel N Rockmore. Unbiased metric learning: On the utilization of multiple datasets and web images for softening bias. In <i>Proceedings of the IEEE International Conference on Computer Vision</i> , pages 1657–1664, 2013.
320 321 322	[GKG <sup>+</sup> 22]	Sachin Goyal, Ananya Kumar, Sankalp Garg, Zico Kolter, and Aditi Raghunathan. Finetune like you pretrain: Improved finetuning of zero-shot vision models. <i>arXiv</i> preprint arXiv:2212.00638, 2022.
323 324 325 326	[JYX <sup>+</sup> 21]	Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In <i>International Conference on Machine Learning</i> , pages 4904–4916. PMLR, 2021.

[KCJ<sup>+</sup>21] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, 327 Dinghuai Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization 328 via risk extrapolation (rex). In International Conference on Machine Learning, pages 329 5815-5826. PMLR, 2021. 330 [KIW22] Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training 331 is sufficient for robustness to spurious correlations. arXiv preprint arXiv:2204.02937, 332 2022. 333 [KNST23] Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. Causal reasoning 334 and large language models: Opening a new frontier for causality. arXiv preprint 335 arXiv:2305.00050, 2023. 336 [KSM<sup>+</sup>21] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, 337 Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena 338 Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In International 339 Conference on Machine Learning, pages 5637–5664. PMLR, 2021. 340 [LCLBC20] Yannick Le Cacheux, Hervé Le Borgne, and Michel Crucianu. Using sentences as 341 semantic representations in large scale zero-shot learning. In Computer Vision-ECCV 342 2020 Workshops: Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16, pages 343 641-645. Springer, 2020. 344 [LCT<sup>+</sup>22] Yoonho Lee, Annie S Chen, Fahim Tajwar, Ananya Kumar, Huaxiu Yao, Percy Liang, 345 and Chelsea Finn. Surgical fine-tuning improves adaptation to distribution shifts. arXiv 346 preprint arXiv:2210.11466, 2022. 347 [LGPV20] Anne Lauscher, Goran Glavaš, Simone Paolo Ponzetto, and Ivan Vulić. A general 348 framework for implicit and explicit debiasing of distributional word vector spaces. 349 In Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, pages 350 8131-8138, 2020. 351 [LHC<sup>+</sup>21] Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, 352 Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group 353 robustness without training group information. In Marina Meila and Tong Zhang, 354 editors, Proceedings of the 38th International Conference on Machine Learning, volume 355 139 of Proceedings of Machine Learning Research, pages 6781–6792. PMLR, 18–24 356 Jul 2021. 357 [LLWT15] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes 358 in the wild. In Proceedings of the IEEE international conference on computer vision, 359 pages 3730-3738, 2015. 360 [LYSH17] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Deeper, broader and 361 artier domain generalization. In Proceedings of the IEEE international conference on 362 computer vision, pages 5542–5550, 2017. 363 [MFB<sup>+</sup>17] Alexander Miller, Will Feng, Dhruv Batra, Antoine Bordes, Adam Fisch, Jiasen Lu, 364 Devi Parikh, and Jason Weston. ParlAI: A dialog research software platform. In 365 Proceedings of the 2017 Conference on Empirical Methods in Natural Language 366 Processing: System Demonstrations, pages 79-84, Copenhagen, Denmark, September 367 2017. Association for Computational Linguistics. 368 [MSY<sup>+</sup>21] Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and 369 Animesh Mukherjee. Hatexplain: A benchmark dataset for explainable hate speech 370 detection. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, 371 pages 14867-14875, 2021. 372 [NLM19] Jianmo Ni, Jiacheng Li, and Julian McAuley. Justifying recommendations using 373 distantly-labeled reviews and fine-grained aspects. In Proceedings of the 2019 confer-374 ence on empirical methods in natural language processing and the 9th international 375 joint conference on natural language processing (EMNLP-IJCNLP), pages 188–197, 376 2019. 377

378 379 380	[NXP <sup>+</sup> 22]	Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al. Text and code embeddings by contrastive pre-training. <i>arXiv preprint arXiv:2201.10005</i> , 2022.
381 382 383 384	[OWJ <sup>+</sup> 22]	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744, 2022.
385 386 387 388	[PDN <sup>+</sup> 22]	Suzanne Petryk, Lisa Dunlap, Keyan Nasseri, Joseph Gonzalez, Trevor Darrell, and Anna Rohrbach. On guiding visual attention with language specification. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 18092–18102, 2022.
389 390	[RG19]	Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> , 2019.
391 392 393 394	[RKH <sup>+</sup> 21]	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learn- ing transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR, 2021.
395 396	[RWC+19]	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9, 2019.
397 398 399	[SKHL19]	Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distribution- ally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. <i>arXiv preprint arXiv:1911.08731</i> , 2019.
400 401	[TE11]	Antonio Torralba and Alexei A. Efros. Unbiased look at dataset bias. In <i>CVPR 2011</i> , pages 1521–1528, 2011.
402 403 404 405	[TLI+23]	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> , 2023.
406 407 408	[WLW21]	Jialu Wang, Yang Liu, and Xin Eric Wang. Are gender-neutral queries really gender- neutral? mitigating gender bias in image search. <i>arXiv preprint arXiv:2109.05433</i> , 2021.
409 410 411 412 413	[WPL+17]	Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers. Chestx-ray8: Hospital-scale chest x-ray database and bench- marks on weakly-supervised classification and localization of common thorax diseases. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 2097–2106, 2017.
414 415 416	[WZS22]	Junyang Wang, Yi Zhang, and Jitao Sang. Fairclip: Social bias elimination based on attribute prototype learning and representation neutralization. <i>arXiv preprint arXiv:2210.14562</i> , 2022.
417 418 419	[YNPM23]	Yu Yang, Besmira Nushi, Hamid Palangi, and Baharan Mirzasoleiman. Mitigat- ing spurious correlations in multi-modal models during fine-tuning. <i>arXiv preprint</i> <i>arXiv:2304.03916</i> , 2023.
420 421	[ZR22]	Michael Zhang and Christopher Ré. Contrastive adapters for foundation model group robustness. <i>arXiv preprint arXiv:2207.07180</i> , 2022.
422 423 424 425	[ZXU <sup>+</sup> 23]	Sheng Zhang, Yanbo Xu, Naoto Usuyama, Jaspreet Bagga, Robert Tinn, Sam Preston, Rajesh Rao, Mu Wei, Naveen Valluri, Cliff Wong, Matthew Lungren, Tristan Naumann, and Hoifung Poon. Large-scale domain-specific pretraining for biomedical vision-language processing, 2023.