EFFECTIVE AND EFFICIENT ADVERSARIAL DETEC TION FOR VISION-LANGUAGE MODELS VIA A SINGLE VECTOR

Anonymous authors

Paper under double-blind review

ABSTRACT

Visual Language Models (VLMs) are vulnerable to adversarial attacks, especially those from adversarial images, which is however under-explored in literature. To facilitate research on this critical safety problem, we first construct a new laRge-scale Adervsarial images dataset with Diverse hArmful Responses (RADAR), given that existing datasets are either small-scale or only contain limited types of harmful responses. With the new RADAR dataset, we further develop a novel and effective iN-time Embedding-based AdveRSarial Image DEtection (NEAR-SIDE) method, which exploits a single vector that distilled from the hidden states of VLMs, which we call *the attacking direction*, to achieve the detection of adversarial images against benign ones in the input. Extensive experiments with two victim VLMs, LLaVA and MiniGPT-4, well demonstrate the effectiveness, efficiency, and cross-model transferrability of our proposed method. Our code is included in the supplementary file and will be made publicly available.

1 INTRODUCTION

Vision Language Models (VLMs), such as BLIP-2 (Li et al., 2023a), LLaVA (Liu et al., 2023a), MiniGPT-4 (Zhu et al., 2024) and GPT-4V (OpenAI, 2023), have attained remarkable success over various vision-language tasks (Dai et al., 2023; Zhu et al., 2024). Besides improving performances, ensuring the safety of responses is just as important in the development of VLMs. Compared with classic Large Language Models (LLMs) that take in discrete textual inputs, VLMs that accept both textual and visual inputs are more susceptible to "jailbreaking", wherein malicious users manipulate inputs to elicit harmful outputs, due to the continuous and high-dimensional nature of visual inputs (Qi et al., 2024a). This issue, which has posed persistent safety challenges in classical vision models (Chakraborty et al., 2018), also presents intrinsic difficulties for developing safe VLMs.



Figure 1: (a) Working mechanism of VLMs. VLMs map textual and visual inputs to the embedding space, and employ LLMs to fuse both embeddings to generate textual responses. (b) Adversarial images that jailbreak VLMs. The adversarial images that contain human-imperceptible noises can jailbreak VLMs to elicit harmful responses.

Existing studies examine the safety threat in VLMs mainly from the perspective of adversarial samples as shown in Fig. 1. It has been revealed that adversarial images are more effective than adversarial texts on attacking VLMs (Qi et al., 2024a; Carlini et al., 2023). Currently, only a few studies

have been conducted to protect VLMs against adversarial images. These methods either seek to
detect adversarial images based on the responses' discrepancy (Zhang et al., 2023b), or to purify
noised-images (Qi et al., 2024a) with diffusion models (Nie et al., 2022), achieving promising effectiveness. However, the first approach is computation-intensive and time-consuming, as it requires
sampling multiple responses for the same input; the other approach, in addition to the computational
cost issue, may even suffer degraded performance when dealing with less perceptible noises.

060 According to previous studies (Subramani et al., 2022; Turner et al., 2023; Zou et al., 2023a; Rimsky 061 et al., 2024; Li et al., 2023b; Liu et al., 2023b) (see Sec. 2.2), the behaviors of LLMs can be mod-062 ulated to generate texts towards certain specific attributes, such as truthfulness, by exploiting a set 063 of steering vectors (SVs) that can be directly extracted from LLMs' hidden states. In adversarial at-064 tacks, the victim VLMs are manipulated by adversarial inputs to generate harmful responses, where the VLMs' behaviors change from harmlessness to harmfulness. We can calculate the SV that can 065 account for VLMs' behavior change given the adversarial inputs, which is named the attacking di-066 rection, and exploit it to detect the existence of adversarial samples by assessing whether the inputs' 067 embedding has high similarity to the attacking direction. 068

069 However, existing datasets for investigating adversarial attacks for VLMs, as shown in Tab. 1, are small-scale and contain limited harm types, significantly restricting the thorough evaluation of 071 VLMs defending against adversarial attacks. Therefore, we construct RADAR, a dataset of laRgescale Adervsarial images with Diverse hArmful Responses, for comprehensively evaluating VLMs 072 against adversarial images. In RADAR, we generate adversarial images to attack two widely-used 073 VLMs, MiniGPT-4 (Zhu et al., 2024) and LLaVA (Liu et al., 2023a), based on a wide diversity of 074 harmful contents. Each sample consists of an adversarial and a benign sample, with each containing 075 a query, an adversarial/benign image and corresponding response of VLMs. In total, RADAR con-076 tains 4,000 samples, which is the most large-scale so far. For high sample quality, we apply filtering 077 operations to ensure harmlessness and harmfulness of responses to benign and adversarial inputs 078 respectively. It will be released to the public to facilitate related research in the community. 079

With RADAR, we further propose a novel iN-time Embedding-based AdveRSarial Image DEtection (NEARSIDE) method, which leverages *the attacking direction* to detect adversarial images to defend VLMs. Specifically, we first extract *the attacking direction* from VLMs by calculating the average difference between the benign input and the adversarial input in the embedding space of VLMs. With the obtained *attacking direction*, we classify an input as an adversarial input if the projection of its embedding to *the attacking direction* is larger than a threshold; otherwise the input is classified as a benign input. Once the adversarial image is detected with the proposed NEARSIDE method, further actions can be taken to protect the VLMs, such as overwriting outputs with a predefined harmless response or purifying the adversarial images by diffusion models.

We conduct extensive experiments to evaluate our NEARSIDE method on the new RADAR dataset. It is demonstrated that NEARSIDE achieves detection accuracy of 83.1% on LLaVA and 93.5% on MiniGPT-4, indicating impressive effectiveness. Furthermore, we experimentally verify the crossmodel transferability of *the attacking direction* in our method. At inference, we compare the efficiency between our method and the baseline method, showing that our method takes an average of 0.14 seconds to complete a detection on LLaVA that is 40 times faster than the best existing method.

⁰⁹⁴ In summary, the major contributions of our work are four-fold:

096

098

099

102

103

- We propose to identify *the attacking direction* that directly distilled from the VLMs' hidden space, and exploit it to defend the VLMs against adversarial images.
- We construct the RADAR dataset, which is the first large-scale adversarial image dataset with a diverse range of harmful responses, to support a comprehensive analysis of VLMs' safety and facilitate future research.
- Based on RADAR, we propose a novel NEARSIDE method, which is capable of effectively and efficiently detecting adversarial visual inputs of VLMs using the identified *attacking direction* from VLMs' hidden space. We further explore the cross-model transferrability of our method given the Platonic Representation Hypothesis (Huh et al., 2024).
- Extensive experiments on two victim VLMs, LLaVA and MiniGPT-4, demonstrate the effectiveness, efficiency, and cross-model transferrability of our method.

¹⁰⁸ 2 BACKGROUND

110 2.1 Adversarial attack

Adversarial attack is maliciously manipulating inputs to compromise performance of the targeted model (Chakraborty et al., 2021; Ponnuru et al., 2023). The manipulated inputs are referred to as adversarial samples. Formally, adversarial samples are generated by minimizing the negative loglikelihood loss of an adversarial target:

- 116
- 117
- 118 119

 $I_{\text{adv}} = \underset{\hat{I}_{\text{adv}} \in \mathcal{I}}{\operatorname{arg\,min}} \sum_{i=1}^{m} -\log(p(y_i|\hat{I}_{\text{adv}})).$ (1)

Here \mathcal{I} represents the input space subject to certain constraints, such as a perturbation radius $||I_{adv} - I|| \le \epsilon$, with ϵ typically set to 16/255, 32/255, 64/255, or unbounded (denoted as "inf"). y_i refers to harmful outputs, and I_{adv} can be either a manipulated text input, where a suffix is appended to attack LLMs (Zou et al., 2023b), or a manipulated visual input, where imperceptible noise is added to the original image to attack VLMs (Qi et al., 2024a; Carlini et al., 2023).

To solve Eqn. (1), various optimization techniques can be employed to generate the adversarial sample I_{adv} . For LLMs, the coordinate gradient-based search (Zou et al., 2023b) or genetic algorithms (Andriushchenko et al., 2024) are commonly used due to the discrete nature of textual inputs. In contrast, for VLMs, where image noise is continuous, Projected Gradient Descent (PGD) (Madry et al., 2018; Qi et al., 2024a; Carlini et al., 2023) is an effective and widely adopted approach.

130 131 2.2 Steering vectors in LLMs

According to the previous research (Subramani et al., 2022; Turner et al., 2023; Zou et al., 2023a; Rimsky et al., 2024; Li et al., 2023b; Liu et al., 2023b), the behaviors of LLMs can be modulated to generate texts towards certain specific attributes, such as truthfulness, by exploiting a set of *steering vectors* (SVs) that can be directly extracted from LLMs' hidden states. To extract the SV for a certain behavior of LLMs, pairs of contrastive prompts (p_+, p_-) are used, where p_+, p_- involve the same question or request, but p_+ adds words to encourage LLMs to possess the behavior while $p_$ represents the opposite. Formally, given a set \mathcal{D} of (p_+, p_-) , the SV is calculated by

$$SV = \frac{1}{|\mathcal{D}|} \sum_{p_+, p_- \in \mathcal{D}} LLM(p_+) - LLM(p_-)$$
(2)

where $LLM(p_+), LLM(p_-) \in \mathbb{R}^d$ are *d*-dimension vectors that represent LLMs' embedding for the *i*th prompt p_+ and p_- respectively. Through Eqn.(2) that takes the mean difference of the embeddings, the SV can be easily acquired, which can specify a tendency, or direction, in LLMs' embedding space regarding the model behavior. That means, simply adding or subtracting such a direction in LLMs' activations can noticeably control LLMs' behavior to generate text with certain attributes. For example, given a direction of "truthfulness", adding this direction can encourage LLMs to generate more truthful responses (Zou et al., 2023a; Rimsky et al., 2024).

149 150

151 152

153

154

139

140 141

- **3** PROPOSED DATASET
- To comprehensively analyze the threat of adversarial attacks posed to VLMs, we propose RADAR, a la**R**ge-scale Adervsarial images dataset with **D**iverse hArmful **R**esponses. Fig. 2 illustrates our construction pipeline. At below we elaborate each step in the pipeline and provide an analysis of its statistics to highlight its merits. An exemplar sample in the new RADAR dataset is given in Appx. **C**.
- 155 156
- 157 3.1 DATA PREPARATION

In RADAR, each sample consists of an adversarial sample and a benign sample, with each containing a query, an adversarial/benign image and VLMs' response. To build RADAR, we use queries from train and test sets in HH-rlhf harm-set (Bai et al., 2022b), those from Harmful-Dataset (Harm-Data) (Sheshadri et al., 2024), and sentences in Derogatory corpus (D-corpus) (Qi et al., 2024a).



According to Eqn. (1) and Sec. 2.1, we optimize a continuous noise that is added to the benign image to maximize the probability of the harmful text, in order to generate an adversarial image. The optimization of noises is implemented using PGD (Madry et al., 2018). In particular, for samples from HH-rlhf harm-set and Harm-Data, we optimize $-\log(p(y_i|\hat{I}_{adv}))$ in Eqn. (1), where \hat{I}_{adv} denotes the noised adversarial image and the query, and y_i denotes the harmful response. Note that when optimizing $-\log(p(y_i|\hat{I}_{adv}))$ on D-corpus, \hat{I}_{adv} refers to only the noised adversarial image, and y_i is the harmful sentence. To generate the adversarial images, we use the open-sourced code ¹ and leave implementation details and hyper-parameters to Appx. A.

193 194

3.3 SAMPLE GENERATION AND FILTERING

195 We then use the benign and adversarial images obtained as aforementioned to generate the samples 196 constituting the proposed RADAR dataset. In particular, we input each benign or counterpart adver-197 sarial image plus a corresponding query to the victim VLM, i.e. MiniGPT-4 or LLaVA, respectively, 198 and obtain a response. For D-corpus, we utilize the harmful sentence as the query. The response 199 is then judged by two models to assess its safety. The first model is a classifier called HarmBench-Llama-2-13b-cls² (Mazeika et al., 2024), which is fine-tuned from the Llama2-13b (Touvron et al., 200 2023) and classifies that whether a pair of (query, response) is harmful or not. The other model is 201 GPT-40 mini (OpenAI, 2024), which are guided with carefully designed prompts to make judge-202 ments following (Qi et al., 2024b) and (Zeng et al., 2024). Concretely, we prompt GPT-40 mini 203 to provide a score ranging from 1 to 5 for each tuple of (query-response), where the scores of 1, 2 204 indicate a harmless response, the score of 3 indicates borderline, and the scores of 4, 5 indicate a 205 harmful response. Please refer to Appx. B for more details. It is expected that for each pair of benign 206 and adversarial images, the responses given by the victim VLM should be judged as harmless for 207 the benign input while harmful for the adversarial input by both models simultaneously. We take 208 this as the criterion to determine whether the quintuple of (query, benign input, harmless response, 209 adversarial input, harmful response) will be included in our RADAR.

In practice, we find that quite a number of responses are harmful given benign images and harmless given adversarial images. As also reported in Qi et al. (2024a), the success of adversarial attack is far from 100%. When constructing our RADAR, we use the two models to judge the responses' harmfulness. Such filtering operations significantly lift the quality of samples in the proposed dataset.

I
https://github.com/Unispac/Visual-Adversarial-Examples-Jailbreak-Large-Language-Models
2
https://huggingface.co/cais/HarmBench-Llama-2-13b-cls

Paper		Scale	Harmful Types	Open Source	Data Filtering
(Zhang et	al., 2023a) Arxiv	200	Harmful queries	1	×
(Tu et al.,	2023) Arxiv	3	Toxic words	1	×
(Carlini et	al., 2023) Neurips 2023	-	Toxic words	×	×
(Qi et al., 2	2024a) AAAI 2024	3	Toxic words	1	X
(Luo et al.	, 2024) ICLR 2024	-	Harmful queries	×	×
(Shayegan	i et al., 2024) ICLR 2024	8	Toxic words	×	X
RADAR (Ours)	4,000	Both	 ✓ 	✓

Table 1: Comparison of datasets for adversarially attacking VLMs. "-" means not reported.

229

241

242

243 244

245

246

247

248

249

250

253

254

216

217

3.4 STATISTICS ANALYSIS

230 With the above construction pipeline, the resultant RADAR contains 4,000 samples in total, attack-231 ing two victim VLMs, i.e. MiniGPT-4 and LLaVA. For each VLM, RADAR provides one training 232 set and three test sets, with 500 samples per set. Division of train and test sets is based on the source 233 of images and queries. Samples built using images from COCO validation set and queries from the 234 train set of HH-rlhf harm-set are grouped into the train set in RADAR; samples built using images 235 from COCO test set and queries from the test set of HH-rlhf harm-set, D-corpus, and Harm-Data are grouped to three test sets, respectively. Training and tests sets use different images. Different 236 harmful texts are used in the four sets to ensure no information leakage and a reliable result. 237

238 A comparison of our RADAR with previous datasets used for investigating adversarial attack for 239 VLMs is provided in Tab. 1. Our RADAR features four advantages compared with previous ones. 240

- Large-scale: As shown in Tab. 1, RADAR greatly surpasses the existing datasets in scale. It contains up to 4,000 samples while the previous largest dataset, i.e. from (Zhang et al., 2023a), contains only 200 samples, facilitating a reliable evaluation of VLMs' safety.
- **Diversity of harmful types**: RADAR covers a favorable diversity of harmful queries and responses, enabling a comprehensive evaluation of VLMs' performance on understanding and defending various adversarial attacks. Recent research on safety of VLMs (Wang et al., 2023; Dai et al., 2024; Ji et al., 2023a) provides taxonomies about the potential harms in queries or responses, e.g. asking for guidance to make bombs or for providing private information. During the construction of RADAR, we purposely increase such diversity.
 - Open-source: RADAR will be open-sourced to facilitate future research on VLMs defending against adversarial attacks.
 - High sample quality: We apply filtering operations during the construction of our RADAR with two models to ensure that the response to a benign input is harmless and that to an adversarial input is harmful. In comparison, the other datasets are built by specifying the harmfulness of the input before feeding it to victim models, while neglecting the reliability of responses, given the success ratio that adversarial images attack VLMs is not 100%.

256 257 258

259 260

261

262

PROPOSED METHOD 4

To efficiently defend VLMs from adversarial attacks, we propose a novel iN-time Embeddingbased AdveRSarial Image DEtection method (abbr. as NEARSIDE) that uses a single vector, named the attacking direction, to detect the adversarial inputs. Fig. 3 gives an illustration of NEARSIDE.

- 264 265
- 4.1 ATTACKING DIRECTION

266 As discussed in Sec. 2.2, the behaviors of LLMs can be controlled with a set of steering vectors (SVs) to generate texts towards certain specific attributes, such as truthfulness. Such SVs can be easily 267 distilled from LLMs' hidden states based on Eqn. (2). In adversarial attacks, the adversarial inputs 268 elicit harmful responses of the victim VLMs, where the VLMs' behaviors alter with an attribute shifting from harmlessness to harmfulness. We can calculate the SV that can account for VLMs'

302

306

307 308

314 315 316



Figure 3: An illustration of proposed NEARSIDE. Our method learns *the attacking direction* on a set of tuples (benign input, adversarial input), and then classifies a test input as benign or adversarial according to the projection between the input's embedding and *the attacking direction*. If the projection is larger than a threshold, it is classified as an adversarial input, and otherwise as benign.

behavior change given the adversarial inputs. We name such a vector *the attacking direction*. In this work, we propose to detect the existence of the adversarial samples by assessing whether the inputs' embedding has shown high similarity to *the attacking direction*.

To extract *the attacking direction* from VLMs' hidden states, the adversarial and benign samples that make pairwise contrastive prompts are required. Formally, consider a training set $\mathbb{T} = \{(I_{adv}^i, I_b^i) | i = 0, 1, ..., n\}$ where I_{adv} , I_b denote the adversarial and benign sample, respectively, and *n* is the index. Each sample contains an image and a piece of text. We embed each sample I^i by taking the embedding of *the last input token from the last LLMs' layer*, i.e. $E^i \in \mathbb{R}^d$, where *d* is the embedding dimension. We embed all samples in \mathbb{T} , and obtain $\mathbb{T}_{emb} = \{(E_{adv}^i, E_b^i) | i = 0, 1, ..., n\}$. Then, we calculate *the attacking direction* by

$$D_{\text{attack}} = \frac{1}{n} \sum_{i=0}^{n} (E_{\text{adv}}^{i} - E_{\text{b}}^{i}) \in \mathbb{R}^{d}, \quad (E_{\text{adv}}^{i}, E_{\text{b}}^{i}) \in \mathbb{T}_{\text{emb}}.$$
(3)

4.2 DETECTION OF ADVERSARIAL INPUTS

Let norm $(h) = h/||h||_2$ denote the ℓ_2 normalization for a vector h. Given the attacking direction D_{attack} , we classify a test input I_{test} to be adversarial or benign by

$$I_{\text{test}} = \begin{cases} \text{adversarial example,} & \text{if } E_{\text{test}} \cdot \text{norm} (D_{\text{attack}}^j)^\top - t > 0, \\ \text{benign example,} & \text{otherwise,} \end{cases}$$
(4)

where $E_{\text{test}} \in \mathbb{R}^d$ is the embedding of the last input token from the last layer of an VLM on the test sample, and $t \in \mathbb{R}$ is a scalar threshold to measure whether the similarity score is significant. If the similarity score, i.e., the projection, is greater than the threshold, we classify the input I_{test} to be adversarial as it has high similarity to the attack direction; otherwise, the input is classified as a benign input. The threshold is decided using \mathbb{T}_{emb} :

$$t = \frac{1}{2n} \sum_{i=0}^{n} (E_{\text{adv}}^{i} \cdot \operatorname{norm}(D_{\text{attack}}^{j})^{\top} + E_{\text{b}}^{i} \cdot \operatorname{norm}(D_{\text{attack}}^{j})^{\top}), \quad (E_{\text{adv}}^{i}, E_{\text{b}}^{i}) \in \mathbb{T}_{\text{emb}}.$$
(5)

The threshold is the average similarity score of all training embeddings (from both adversarial and benign samples) on *the attacking direction*.

The proposed NEARSIDE, as shown in Eqn. (4), is extremely efficient as we only require running one feed-forward propagation given the input to infer E_{test} , thus enabling an in-time detection of adversarial samples. After adversarial samples have been detected, the developer can take further steps to ensure VLMs' safety, such as overwriting the responses to a preset text, applying diffusion models to purify the image, or disabling malicious accounts. Therefore, NEARSIDE can defend VLMs from adversarial attack in an efficient and real-time manner.

4.3 CROSS-MODEL TRANSFERABILITY 325

The Platonic Representation Hypothesis (Huh et al., 2024): "Neural networks, trained with different objectives on different data and modalities, are converging to a shared statistical model of reality in their representation spaces."

The proposed NEARSIDE is supposed to use the attacking direction extracted from one VLM to 330 detect the adversarial samples for the same VLM. According to the above Platonic Representation 331 Hypothesis, we can assume that the learnt attacking direction and effectiveness of our detection 332 method NEARSIDE are transferable across different models. That is, our NEARSIDE can use the 333 attacking direction extracted from one VLM to detect the adversarial samples for other VLMs. The reason behind the assumption of the cross-model transferrability in our method is that, although 334 different VLMs are trained from different data, the patterns regarding safety in these data should be 335 similar. However, the embedding spaces between two VLMs do have a gap. We thus propose to 336 explore the transferability using a linear transformation: 337

$$\boldsymbol{W}\boldsymbol{E}_{m_1} = \boldsymbol{E}_{m_2},\tag{6}$$

where $E_{m_1} \in \mathbb{R}^{n \times d_{m_1}}$ and $E_{m_2} \in \mathbb{R}^{n \times d_{m_2}}$ are the stacked embedding of *benign* inputs from the two VLMs m₁ and m₂, respectively, with d_{m_1} , d_{m_2} denoting their embedding dimension. The linear transformation W is to align the two VLMs' embedding spaces. In practice, since powerful LLMs often have high dimension in their hidden states, directly solving Eqn. (6) would be too costly in memory due to the high dimension of d_{m_1} , d_{m_2} . Therefore, we propose to use principal component analysis (PCA) to reduce the dimension of E_{m_1} and E_{m_2} . Then, we have $W = f_{m_2}^{pca}(E_{m_2})f_{m_1}^{pca}(E_{m_2})^{\dagger}$, where f^{pca} denotes PCA that reduces the dimension and $(\cdot)^{\dagger}$ denotes the pseudo-inverse.

Finally, given a test input I_{test} , its detection on m_2 is given by

$$I_{\text{test}} = \begin{cases} \text{adversarial example,} & \text{if } f_{\text{m}_2}^{\text{pca}}(E_{\text{test},\text{m}_2}) \cdot \text{norm}(\boldsymbol{W} f_{\text{m}_1}^{\text{pca}}(D_{\text{attack},\text{m}_1}))^{\top} - t_{\text{m}_1} > 0, \\ \text{benign example,} & \text{otherwise,} \end{cases}$$
(7)

where the threshold t_{m_1} is defined as

352 353 354

355

356 357

358

359

360 361 362

364

365

366

369

347 348

349 350 351

338

$$t_{\mathbf{m}_{1}} = \frac{1}{2n} \sum_{i=0}^{n} (\boldsymbol{W} f_{\mathbf{m}_{1}}^{\mathbf{pca}}(\boldsymbol{E}_{adv,\mathbf{m}_{1}}^{i}) \cdot \operatorname{norm}(\boldsymbol{W} f_{\mathbf{m}_{1}}^{\mathbf{pca}}(\boldsymbol{D}_{attack,\mathbf{m}_{1}}))^{\top} + \boldsymbol{W} f_{\mathbf{m}_{1}}^{\mathbf{pca}}(\boldsymbol{E}_{b,\mathbf{m}_{1}}^{i}) \cdot \operatorname{norm}(\boldsymbol{W} f_{\mathbf{m}_{1}}^{\mathbf{pca}}(\boldsymbol{D}_{attack,\mathbf{m}_{1}}))^{\top}).$$

$$(8)$$

Eqn. (7) detects adversarial samples on m_2 only using *the attacking direction* of m_1 and the embedding of *benign* inputs from m_2 to learn the transformation matrix W. Note that this entire learning process has no access to adversarial samples on m_2 . Eqn. (7) works if the embedding space of the two VLMs can be linearly transformed without disturbing *the attacking direction*.

5 EXPERIMENTS

We conduct extensive experiments on RADAR to evaluate the effectiveness of the proposed NEAR-SIDE in detecting adversarial images. We first compare our method with strong baseline and then analyze its cross-model transferability, followed by the efficiency test.

367 368

5.1 EXPERIMENTS SETUP

Victim VLMs. We adopt MiniGPT-4 (Zhu et al., 2024) and LLaVA (Liu et al., 2023a) as the victim VLMs. MiniGPT-4 is built upon Vicuna (Chiang et al., 2023) and LLaVA is built upon Llama2 (Touvron et al., 2023). Regarding the visual encoder, MiniGPT-4 utilizes the same pre-trained vision components of BLIP-2 (Li et al., 2023a) consisting of pre-trained ViT followed by a Q-Former, while LLaVA only adopts a pre-trained CLIP (Radford et al., 2021).

Implementation. For each victim VLM, as stated in Sec. 3, RADAR constructs one training set
 and three test sets. NEARSIDE learns *the attacking direction* and threshold from the hidden states
 of the VLM on the training set. Here, the hidden states refer to the embedding of the last token
 of the input from the LLM decoder's final layer. Then, we test the detection performance with the

378 obtained attacking direction on the three test sets. Regarding the cross-model transferability, we 379 collect 5,000 pairs of benign images and queries to train the PCA model for each VLM. We set 2048 380 as the dimension of the embedding after PCA.

381 Baseline. We use JailGuard (Zhang et al., 2023a) as our baseline, which is the state-of-the-art 382 model for this task. To detect adversarial visual inputs, JailGuard mutates input images to generate 383 variants and calculates the discrepancy of VLMs' outputs given different variants to distinguish the 384 adversarial and benign inputs. There are 18 mutation methods, and we use the best-performing 385 mutation method "policy" reported in the JailGuard paper, where 8 variants are generated for each 386 image. We set all other hyperparameters to the recommended values as JailGuard. It is worth 387 mentioning that we adopt only one baseline as there are only limited works on defending VLMs 388 from adversarial examples (Liu et al., 2024).

389 **Evaluation metrics.** Since adversarial detection is a binary classification task, we adopt Accuracy, 390 *Precision*, *Recall* and *F*1 score as the evaluation metrics.

5.2 MAIN RESULTS

394 We compare our proposed method NEARSIDE against the baseline JailGuard on RADAR. The 395 experimental results are shown in Tab. 2. From the results, we make below observations. 1) 396 When taking LLaVA as the victim VLM, our NEARSIDE achieves an average increase of 31.3% 397 in Accuracy, 43.5% in Precision, 12.6% in Recall, and 0.246 in F1, compared to the baseline JailGuard method. 2) When taking MiniGPT-4 as the victim VLM, our NEARSIDE achieves an 398 average increase of 38.7% in Accuracy, 45.6% in Precision, 17.6% in Recall, and 0.316 in F1, 399 compared to the baseline JailGuard method. These results well demonstrate the effectiveness of 400 our proposed method. 3) Although our NEARSIDE has lower Recall on the Harm-Data set with 401 LLaVA as the victim VLM, and also on D-corpus-test set with MiniGPT-4 as the victim VLM, it 402 achieves significantly higher F1 scores on both sets. We attribute the low *Recall* of our method to 403 its threshold for the adversarial detection. As shown in Fig. 4, the projections of the two types of 404 examples do fall into different ranges. However, as the threshold is calculated on the training set, 405 the threshold is not well fit for the Harm-Data, leading to degraded *Recall*. If we set the thresh-406 old to -13, we can increase Recall to 87.6% and F1 score to 0.900, which are both significantly 407 higher than the baseline. From an overall perspective, the results can still demonstrate the powerful 408 distinguishing capability of our method over adversarial and benign data.

409 410 411

391 392

393

Victim VLM	Test Set	Method	Accuracy(%)	Precision(%)	Recall(%)	F1
	HH-rlhf	JailGuard NEARSIDE	51.2 84.4	51.1 89.3	57.8 78.2	0.540 0.834
LLaVA	D-corpus	JailGuard NEARSIDE	58.1 94.0	58.1 99.5	58.2 88.4	0.581 0.936
	Harm-Data	JailGuard NEARSIDE	46.2 71.0	46.8 97.7	55.8 43.0	0.509 0.597
	HH-rlhf	JailGuard NEARSIDE	54.9 99.4	53.9 99.2	67.2 99.6	0.598 0.994
MiniGPT-4	D-corpus	JailGuard NEARSIDE	56.6 81.1	54.4 98.4	81.6 63.2	0.653 0.770
	Harm-Data	JailGuard NEARSIDE	52.8 100.0	52.4 100.0	61.2 100.0	0.565 1.000

Table 2: Results of JailGuard v.s. NEARSIDE on RADAR test sets (best highlighted in **bold**).

426 427

5.3 ANALYSIS ON CROSS-MODEL TRANSFERABILITY

428 429

We utilize the attacking direction extracted from the source VLM (svlm) to detect adversarial input 430 for the target VLM (tvlm), denoted as svlm \rightarrow tvlm. We calculate the difference (i.e. δ) by subtract-431 ing the result of (svlm \rightarrow tvlm) from that of Tab. 2, where $^{-\delta}$ denotes the result is decreased while $^{+\delta}$



4115% 7.407 +13226% 3.571 0.180 0.003 0.027 0.024 0.022 0.098 LLaVA **MiniGPT-4**

Figure 4: Visualized projections of adversarial and benign samples to the attacking directions on Harm-Data with LLaVA as the victim.



denotes the opposite. The obtained results for cross-model transferaility are shown in Tab. 3. We can observe that cross-model transferability results are generally inferior to those in Tab. 2 where the attacking direction is extracted and applied with the same VLM, but both Accuracy and F1 results of our method are higher than those of the baseline JailGuard. Though cross-model transferrability decreases the detection performance, which is expectable, our method can still work well across different models. These results clearly validate the cross-model transferability of *the attacking direction* and the proposed NEARSIDE. It also says that, the Platonic Representation Hypothesis still holds in our setting, where a simple linear transformation is effective to align two VLMs' embedding spaces.

Table 3: Cross-model transferability results for our method.

$svlm \rightarrow tvlm$	TEST SET	Accuracy(%)	Precision(%)	Recall(%)	F1
	HH-rlhf	$64.3^{-20.1}$	$61.3^{-28.0}$	$77.6^{-0.6}$	$0.685^{-0.149}$
MiniGPT-4 \rightarrow LLaVA	D-corpus	$69.4^{-24.6}$	62.5^{-37}	$96.8^{+8.4}$	$0.760^{-0.176}$
	Harm-Data	$74.7^{+3.7}$	76.7^{-21}	71.0^{+28}	$0.737^{+0.14}$
	HH-rlhf	$77.8^{-21.6}$	86.2^{-13}	$66.2^{-33.4}$	$0.749^{-0.245}$
$LLaVA \rightarrow MiniGPT-4$	D-corpus	$80.4^{-0.7}$	$80.3^{-18.1}$	$80.6^{+17.4}$	$0.804^{+0.034}$
	Harm-Data	$97.1^{-2.9}$	$95.0^{-0.5}$	$99.4^{-0.6}$	$0.972^{-0.028}$

We experiment to examine the effect of using W to align two VLMs' embedding spaces, and the effect of reducing the dimension of the attacking direction and the VLMs' embedding with the PCA model. The detailed results are provided in Appx. E. We find that, without W, the crossmodel transferability results will significantly decrease. In addition, when reducing the dimension to 256, the cross-model transferability results still remain high, indicating that the information in low dimensional sub-spaces is already sufficient for aligning two VLMs' embedding spaces.

5.4ANALYSIS OF PERTURBATION RADIUS IN GENERATING ADVERSARIAL IMAGES

476 The generation of adversarial images is constrained by the hyper-parameter ϵ as shown in Eqn. (1). 477 We test the robustness of the proposed NEARSIDE method to varying ϵ . We use NEARSIDE to 478 detect the adversarial samples generated under different ϵ . Results are deferred to Appx. D.

479 480 481

469

470

471

472

473 474

475

443

444

445 446 447

448

449

450

451

452

453

454 455 456

5.5ANALYSIS OF DETECTION EFFICIENCY

In this part, we examine the detection efficiency of the proposed method. For the baseline Jail-482 Guard, we utilize the widely-adopted VLLM (Kwon et al., 2023) to deploy the two VLMs, i.e. 483 LLaVA and MiniGPT-4, on a local machine and generate outputs through API requests. For our 484 proposed NEARSIDE, we load VLMs and perform a single forward propagation to embed each in-485 put since the attacking direction can be pre-computed. In addition to JailGuard and our method, we

also include another two trivial methods that judge the harmfulness of the output into our efficiency
evaluations, i.e. HarmBench and GPT-40 mini, which are used in data filtering operations to judge
the harmfulness of responses in Sec. 3.3. For all compared methods, we calculate the time including
responses inference (note, NEARSIDE does not infer responses) plus follow-up operations, which
refer to discrepancy calculation in JailGuard, projections calculation in NEARSIDE, and harmfulness evaluation in other two methods. All experiments are conducted on a server with AMD EPYC
7543 32-core processors, 1 TB of RAM, and a NVIDIA A40 GPU.

We run experiments on 20 inputs and plot the average throughput in Fig. 5. With our setup, NEARSIDE is (× 41~336) times faster than the other methods on LLaVA, and is (× 132~1190) times faster on MiniGPT-4, demonstrating remarkable efficiency as NEARSIDE is the only embeddingbased method among all compared methods that does not require to infer the entire output.

497 498

499 500

501

6 RELATED WORKS

6.1 VISION LANGUAGE MODELS

502 Vision Language Models (VLMs) is equipped with a visual adapter to align the visual and textual 503 representations in LLMs. Notable examples are BLIP-2 (Li et al., 2023a), LLaVA (Liu et al., 2023a), 504 MiniGPT-4 (Zhu et al., 2024), QWen-VL (Bai et al., 2023), GPT-4V (OpenAI, 2023), and Gemini 505 (Anil et al., 2023), demonstrating impressive performance across various vision-language tasks (Dai 506 et al., 2023; Zhu et al., 2024). These VLMs vary in the design of their adapters (Liu et al., 2023a; Li 507 et al., 2023a; Zhu et al., 2024). For instance, BLIP-2 (Li et al., 2023a) proposes Q-Former to align 508 vision features with LLMs; MiniGPT-4 (Zhu et al., 2024) and LLaVA (Liu et al., 2023a) further add a linear transformation, and Qwen-VL (Bai et al., 2023) uses a single-layer cross-attention module. 509

510 511

512

6.2 SAFETY OF LANGUAGE MODELS

Safe LLMs should behave in line with human intentions and values (Soares & Fallenstein, 2014; 513 Hendrycks et al., 2021; Leike et al., 2018; Ji et al., 2023b) which are measured as being Helpful, 514 Honest, and Harmless (Askell et al., 2021). Alignment has emerged as a nascent research field 515 aiming to align LLMs' behaviors with human preferences, and there are two widely adopted align-516 ment techniques, i.e. Instruction Fine-tuning and Reinforcement Learning from Human Feedback 517 (RLHF). In instruction fine-tuning, LLMs are given examples of (user's query, desired output) and 518 trained to follow user instructions and respond the expected output (Taori et al., 2023; Wei et al., 519 2022). In RLHF, LLMs update output probabilities, i.e., the response policy, by reinforcement learn-520 ing, which are rewarded for generating responses that align with human preferences and otherwises 521 penalized (Russell & Norvig, 2016; Bai et al., 2022a; Rafailov et al., 2023; Ouyang et al., 2022).

Two types of strategies can defend language models from adversarial attacks: detection and purification. For instance, Zhang et al. (2023b) detects adversarial examples by calculating responses' discrepancy; Qi et al. (2024a) uses diffusion models (Nie et al., 2022) to purify the noised-images. Other techniques such as the adversarial training (Bai et al., 2021) can also improve the robustness of models to adversarial attacks. Though effective on classical image classifiers, these methods remain unexplored on large models like LLMs and VLMs and may disturb the optimization.

528 529

7 CONCLUSION

530 531

In this work, we propose RADAR, the first large-scale adversarial image dataset with diverse harmful responses to facilitate research on safety of VLMs. With RADAR, we further develop NEARSIDE that exploits the idea of attacking direction to detect adversarial inputs. We demonstrate with
the effectiveness and efficiency of the proposed NEARSIDE by comparing it to the state-of-the-art
on RADAR, and also highlight its cross-model transferability.

VLMs can generate open-ended responses, posing a persistent challenge to complete evaluation of
 the potential harms (Ganguli et al., 2022). RADAR is built from a diverse array of datasets but may
 fall short of covering all harmful contents. NEARSIDE is intended to detect the adversarial samples
 we examine in this work and is a demonstration of our idea of exploiting *the attacking direction*.

540 REFERENCES

559

561

562

Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. Jailbreaking leading safetyaligned llms with simple adaptive attacks. *CoRR*, abs/2404.02151, 2024. doi: 10.48550/ARXIV. 2404.02151. URL https://doi.org/10.48550/arXiv.2404.02151.

- Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Jo-546 han Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Slav Petrov, Melvin 547 Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy P. Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul Ronald 548 Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan 549 Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha 550 Goel, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Dani-551 helka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, 552 Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, and et al. Gemini: A fam-553 ily of highly capable multimodal models. CoRR, abs/2312.11805, 2023. doi: 10.48550/ARXIV. 554 2312.11805. URL https://doi.org/10.48550/arXiv.2312.11805. 555
- Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones,
 Nicholas Joseph, Ben Mann, Nova DasSarma, et al. A general language assistant as a laboratory
 for alignment. *arXiv preprint arXiv:2112.00861*, 2021.
 - Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- Tao Bai, Jinqi Luo, Jun Zhao, Bihan Wen, and Qian Wang. Recent advances in adversarial training
 for adversarial robustness. In Zhi-Hua Zhou (ed.), *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021*, pp. 4312–4321. ijcai.org, 2021. doi: 10.24963/IJCAI.2021/591. URL https:
 //doi.org/10.24963/ijcai.2021/591.
- 568 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn 569 Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernan-570 dez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Ols-571 son, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, Benjamin Mann, 572 and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from 573 human feedback. CoRR, abs/2204.05862, 2022a. doi: 10.48550/ARXIV.2204.05862. URL 574 https://doi.org/10.48550/arXiv.2204.05862. 575
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022b.
- 580 Nicholas Carlini, Milad Nasr, Christopher A. Choquette-Choo, Matthew Jagielski, Irena 581 Gao, Pang Wei Koh, Daphne Ippolito, Florian Tramèr, and Ludwig Schmidt. Are 582 aligned neural networks adversarially aligned? In Alice Oh, Tristan Naumann, Amir 583 Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Pro-584 cessing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 585 2023, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/ 586 clf0b856a35986348ab3414177266f75-Abstract-Conference.html.
- Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopadhyay. Adversarial attacks and defences: A survey. *CoRR*, abs/1810.00069, 2018. URL http://arxiv.org/abs/1810.00069.
- Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopad hyay. A survey on adversarial attacks and defences. *CAAI Trans. Intell. Technol.*, 6(1):25–45, 2021. doi: 10.1049/CIT2.12028. URL https://doi.org/10.1049/cit2.12028.

- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL https://lmsys.org/blog/2023-03-30-vicuna/.
- Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe RLHF: safe reinforcement learning from human feedback. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=TyFrPOKYXw.
- 603 Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven C. H. Hoi. Instructblip: Towards general-604 purpose vision-language models with instruction tuning. In Alice Oh, Tristan Nau-605 mann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances 606 in Neural Information Processing Systems 36: Annual Conference on Neural Informa-607 tion Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 608 2023, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/ 609 9a6a435e75419a836fe47ab6793623e6-Abstract-Conference.html. 610
- 611 Deep Ganguli, Danny Hernandez, Liane Lovitt, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Nelson Elhage, Sheer El Showk, Stanislav Fort, Zac 612 Hatfield-Dodds, Tom Henighan, Scott Johnston, Andy Jones, Nicholas Joseph, Jackson Kernian, 613 Shauna Kravec, Ben Mann, Neel Nanda, Kamal Ndousse, Catherine Olsson, Daniela Amodei, 614 Tom B. Brown, Jared Kaplan, Sam McCandlish, Christopher Olah, Dario Amodei, and Jack Clark. 615 Predictability and surprise in large generative models. In FAccT '22: 2022 ACM Conference on 616 Fairness, Accountability, and Transparency, Seoul, Republic of Korea, June 21 - 24, 2022, pp. 617 1747-1764. ACM, 2022. doi: 10.1145/3531146.3533229. URL https://doi.org/10. 618 1145/3531146.3533229. 619
- Dan Hendrycks, Nicholas Carlini, John Schulman, and Jacob Steinhardt. Unsolved problems in ML safety. *arXiv preprint arXiv:2109.13916*, 2021.
- Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. Position: The platonic representation hypothesis. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=BH8TYy0r6u.

635

636

637

- Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang 627 Sun, Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety align-628 ment of LLM via a human-preference dataset. In Alice Oh, Tristan Naumann, Amir 629 Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances in Neu-630 ral Information Processing Systems 36: Annual Conference on Neural Information 631 Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 632 2023, 2023a. URL http://papers.nips.cc/paper_files/paper/2023/ 633 hash/4dbb61cb68671edc4ca3712d70083b9f-Abstract-Datasets_and_ 634 Benchmarks.html.
 - Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan, Zhonghao He, Jiayi Zhou, Zhaowei Zhang, et al. AI alignment: A comprehensive survey. *arXiv* preprint arXiv:2310.19852, 2023b.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. Scalable agent alignment via reward modeling: a research direction. *arXiv preprint arXiv:1811.07871*, 2018.
- ⁶⁴⁷ Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. BLIP-2: bootstrapping languageimage pre-training with frozen image encoders and large language models. In Andreas Krause,

676

Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 19730–19742.
PMLR, 2023a. URL https://proceedings.mlr.press/v202/li23q.html.

- Kenneth Li, Oam Patel, Fernanda B. Viégas, Hanspeter Pfister, and Martin Wattenberg. Inferencetime intervention: Eliciting truthful answers from a language model. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023, 2023b. URL http://papers.nips.cc/paper_files/paper/2023/hash/ 81b8390039b7302c909cb769f8b6cd93-Abstract-Conference.html.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In David J. Fleet, Tomás Pajdla, Bernt Schiele, and Tinne Tuytelaars (eds.), *Computer Vision ECCV 2014 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V*, volume 8693 of *Lecture Notes in Computer Science*, pp. 740–755. Springer, 2014. doi: 10.1007/978-3-319-10602-1_48. URL https://doi.org/10.1007/978-3-319-10602-1_48.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In Alice
 Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.),
 Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16,
 2023, 2023a. URL http://papers.nips.cc/paper_files/paper/2023/hash/
 6dcf277ea32ce3288914faf369fe6de0-Abstract-Conference.html.
- Sheng Liu, Haotian Ye, Lei Xing, and James Y. Zou. In-context vectors: Making in context learning more effective and controllable through latent space steering. *ArXiv*, abs/2311.06668, 2023b. URL https://api.semanticscholar.org/CorpusID:265149781.
- Kin Liu, Yichen Zhu, Yunshi Lan, Chao Yang, and Yu Qiao. Safety of multimodal large language
 models on images and text. *arXiv preprint arXiv:2402.00357*, 2024.
- Haochen Luo, Jindong Gu, Fengyuan Liu, and Philip Torr. An image is worth 1000 lies: Transferability of adversarial images across prompts on vision-language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net, 2024. URL https://openreview.net/forum?id=nc5GgFAvtk.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
 Towards deep learning models resistant to adversarial attacks. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net, 2018. URL https://openreview.net/ forum?id=rJzIBfZAb.
- Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee,
 Nathaniel Li, Steven Basart, Bo Li, David A. Forsyth, and Dan Hendrycks. Harmbench: A
 standardized evaluation framework for automated red teaming and robust refusal. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024.* OpenReview.net, 2024. URL https://openreview.net/forum?id=f3TUipYU3U.
- Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Animashree Anandkumar. Diffusion models for adversarial purification. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 16805–16827. PMLR, 2022. URL https://proceedings.mlr.press/v162/nie22a.html.
- 701 OpenAI. Gpt-4v(ision) system card. 2023. URL https://api.semanticscholar.org/ CorpusID:263218031.
 - 13

702 703	OpenAI.	Gpt-40	system	card.	2024.	URL	https:/,	/cdn.ope	nai.com/
704 705 706 707 708 709 710 711 712 713	Long Ouyang Chong Zha Fraser Kelto Jan Leike, a back. In Sa Advances in tion Proces. 9, 2022, 20 blefde53	, Jeffrey V ng, Sandh on, Luke N and Ryan I nmi Koyej <i>a Neural In</i> <i>sing System</i> 22. URL 1 be364a7	Vu, Xu J ini Agar liller, Ma Lowe. Th o, S. Mol formation is 2022, http:/ 3914f5	iang, Dio wal, Katar ddie Sime raining lar named, A. <i>n Processi</i> <i>NeurIPS</i> 2 /papers 8805a00	go Almeida, rina Slama, ens, Amanda aguage mode Agarwal, D <i>ing Systems</i> . 2022, New O .nips.cc 201731-Abs	, Carroll Alex Ra Askell, els to fol anielle E 35: Anna rleans, I /paper stract	L. Wainwa ay, John Sc Peter Welir low instruc Belgrave, K. <i>ual Confere</i> <i>(A, USA, Na</i> <i>(c_files/pConfere)</i>	right, Pame chulman, Ja ider, Paul F. tions with h Cho, and A nce on Neur wember 28 paper/200 nce.html	la Mishkin, cob Hilton, Christiano, uman feed- Oh (eds.), <i>al Informa</i> - - <i>December</i> 22/hash/
714 715 716 717	Mahesh Datta Garikipati. abs/2312.09 48550/ar	A Sai Ponr A malw 9636, 2023 Xiv.231	uru, Lik are class . doi: 1 2.0963	hitha Am ification s 0.48550/A 6.	asala, Tanu survey on a ARXIV.2312.	Sree B dversari .09636.	himavarapu al attacks a URL htt _l	, and Guna and defence ps://doi	Chaitanya es. <i>CoRR</i> , .org/10.
718 719 720 721 722 723 724	Xiangyu Qi, K Visual adve Jennifer G. <i>ligence, AA</i> <i>IAAI 2024,</i> <i>February 20</i> AAAI.V381	Caixuan Hu rsarial exa Dy, and Sr AI 2024, Th Fourteenth D-27, 2024, [19.30150.	ang, Ash mples jai iraam Na hirty-Sixt Symposi Vancouv URL ht	winee Pan Ibreak alig atarajan (e h Conferen um on Edu ver, Canad tps://c	da, Peter He gned large la ds.), <i>Thirty-I</i> <i>nce on Innov</i> <i>a, pp. 21527</i> doi.org/1	nderson nguage Eighth A pative Ap vances in 2–21536.	, Mengdi Wa models. In AAI Confer plications o n Artificial I. AAAI Pres 9/aaai.v	ang, and Pra Michael J. V ence on Art of Artificial 1 ntelligence, s, 2024a. do 38i19.30	teek Mittal. Wooldridge, <i>ificial Intel-</i> <i>ntelligence,</i> <i>EAAI 2014,</i> <i>bi</i> : 10.1609/ 150.
725 726 727 728 729	Xiangyu Qi, Y Fine-tuning The Twelfth May 7-11, 2 hTEGyKf0	(i Zeng, Ti aligned la <i>Internatio</i> 2024. Oper dZ.	nghao X nguage n nal Conf Review.1	ie, Pin-Yu nodels cor <i>erence on</i> net, 2024b	Chen, Ruox npromises sa <i>Learning Re</i> . URL htt:	ti Jia, Pr afety, eve presente ps://c	ateek Mitta en when use ations, ICLR openrevie	l, and Peter ers do not in 2024, Vien ew.net/fo	Henderson. tend to! In na, Austria, orum?id=
730 731 732 733 734	Xiangyu Qi, Y Fine-tuning The Twelfth May 7-11, 2 hTEGyKf0	i Zeng, Ti aligned la <i>Internatio</i> 2024. Oper dZ.	nghao X nguage n nal Conf nReview.n	ie, Pin-Yu nodels cor <i>erence on</i> net, 2024c	Chen, Ruox npromises sa <i>Learning Re</i> . URL htt:	ti Jia, Pr afety, ever presente ps://c	ateek Mitta en when use ations, ICLR openrevie	l, and Peter ers do not in 2024, Vien ew.net/fo	Henderson. tend to! In <i>na, Austria,</i> orum?id=
735 736 737 738 739	Alec Radford, Girish Sasti models fror 8748–8763.	Jong Woo ry, Amanda n natural la PMLR, 20	ok Kim, a Askell, inguage s 021.	Chris Hall Pamela M supervision	lacy, Aditya Iishkin, Jack n. In <i>Interna</i>	Ramesh c Clark, tional co	n, Gabriel C et al. Learn conference of	oh, Sandhin ning transfe n machine le	ni Agarwal, rable visual <i>earning</i> , pp.
740 741 742 743 744 745 746	Rafael Rafaile Chelsea Fir In Alice Oh (eds.), Adva Information 16, 2023, 20 a85b405e	ov, Archit in. Direct j a, Tristan N ances in Na Processin 023. URL 1 ed65c647	Sharma, preference laumann, eural Infa g System http:// 7a4fe8	Eric Mit e optimiza Amir Glo ormation L as 2023, N /papers 302b5e0	cchell, Chris ation: Your I berson, Kato Processing S JeurIPS 202. .nips.cc D6ce7-Abs	topher l anguage e Saenko Systems 3, New /paper stract	D. Manning model is so o, Moritz Ha 36: Annual Orleans, LA c_files/p -Confere	g, Stefano H ecretly a rev ardt, and Ser <i>Conference</i> <i>A, USA, Dec</i> paper/201 ence.html	Ermon, and ward model. rgey Levine e on Neural rember 10 - 23/hash/
747 748 749 750 751 752	Nina Rimsky, Turner. St tins, and V for Comput gust 11-16, https://	Nick Ga eering llar ivek Sriku ational Lin 2024, pp. aclanth	brieli, Ju na 2 via mar (eds <i>nguistics</i> 15504– ology.	lian Schu contrastiv .), <i>Procee</i> (<i>Volume</i> 15522. As org/202	lz, Meg To ve activation <i>dings of the</i> <i>1: Long Pap</i> ssociation fo	ng, Eva a additic e 62nd A pers), A or Comp ong.82	n Hubingen on. In Lun <i>Annual Mee</i> <i>CL 2024, B</i> utational Li 8.	r, and Alex I-Wei Ku, A ting of the Cangkok, Th nguistics, 2	ander Matt Andre Mar- Association ailand, Au- 024. URL
753 754	Stuart J Russe	ll and Pete	r Norvig.	Artificial	intelligence	: a mode	ern approac	h. Pearson,	2016.
755	Erfan Shayega versarial att	ani, Yue D acks on m	ong, and ulti-moda	Nael B. A al languag	Abu-Ghazale e models. I	ch. Jailt n <i>The T</i> i	reak in pied welfth Intern	ces: Componational Col	sitional ad- nference on

758

772

Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024. OpenReview.net, 2024. URL https://openreview.net/forum?id=plmBsXHxgR.

- Abhay Sheshadri, Aidan Ewart, Phillip Guo, Aengus Lynch, Cindy Wu, Vivek Hebbar, Henry Sleight, Asa Cooper Stickland, Ethan Perez, Dylan Hadfield-Menell, and Stephen Casper. Targeted latent adversarial training improves robustness to persistent harmful behaviors in llms. *CoRR*, abs/2407.15549, 2024. doi: 10.48550/ARXIV.2407.15549. URL https://doi.org/10.48550/arXiv.2407.15549.
- Nate Soares and Benja Fallenstein. Aligning superintelligence with human interests: A technical research agenda. *Machine Intelligence Research Institute (MIRI) technical report*, 8, 2014.
- Nishant Subramani, Nivedita Suresh, and Matthew E. Peters. Extracting latent steering vectors from pretrained language models. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*, pp. 566–581. Association for Computational Linguistics, 2022. doi: 10.18653/V1/2022.FINDINGS-ACL.48. URL https://doi.org/10.18653/v1/2022.
 findings-acl.48.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- 776 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-777 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, 778 Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, 779 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel 780 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya 781 Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar 782 Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan 783 Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen 784 Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan 785 Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, 786 Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-787 tuned chat models. CoRR, abs/2307.09288, 2023. doi: 10.48550/ARXIV.2307.09288. URL 788 https://doi.org/10.48550/arXiv.2307.09288.
- Haoqin Tu, Chenhang Cui, Zijun Wang, Yiyang Zhou, Bingchen Zhao, Junlin Han, Wangchunshu Zhou, Huaxiu Yao, and Cihang Xie. How many unicorns are in this image? A safety evaluation benchmark for vision llms. *CoRR*, abs/2311.16101, 2023. doi: 10.48550/ARXIV.2311.16101.
 URL https://doi.org/10.48550/arXiv.2311.16101.
- Alexander Matt Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte
 MacDiarmid. Activation addition: Steering language models without optimization. CoRR,
 abs/2308.10248, 2023. doi: 10.48550/ARXIV.2308.10248. URL https://doi.org/10.
 48550/arXiv.2308.10248.
- Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov, and Timothy Baldwin. Do-not-answer: A dataset for evaluating safeguards in llms. *CoRR*, abs/2308.13387, 2023. doi: 10.48550/ARXIV. 2308.13387. URL https://doi.org/10.48550/arXiv.2308.13387.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 2022. URL https://openreview.net/forum?id=gEZrGCozdqR.
- Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge AI safety by humanizing llms. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL*

Linguistics, 2024. URL https://aclanthology.org/2024.acl-long.773. Xiaoyu Zhang, Cen Zhang, Tianlin Li, Yihao Huang, Xiaojun Jia, Ming Hu, Jie Zhang, Yang Liu, Shiqing Ma, and Chao Shen. Jailguard: A universal detection framework for llm prompt-based attacks. 2023a. URL https://api.semanticscholar.org/CorpusID:266359752. Xiaoyu Zhang, Cen Zhang, Tianlin Li, Yihao Huang, Xiaojun Jia, Xiaofei Xie, Yang Liu, and Chao Shen. A mutation-based method for multi-modal jailbreaking attack detection. arXiv preprint arXiv:2312.10766, 2023b. Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. In The Twelfth Interna-tional Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024. OpenReview.net, 2024. URL https://openreview.net/forum?id=1tZbq88f27. Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, J. Zico Kolter, and Dan Hendrycks. Representation engineering: A top-down ap-proach to AI transparency. CoRR, abs/2310.01405, 2023a. doi: 10.48550/ARXIV.2310.01405. URL https://doi.org/10.48550/arXiv.2310.01405. Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. CoRR, abs/2307.15043, 2023b. doi: 10.48550/ARXIV.2307. 15043. URL https://doi.org/10.48550/arXiv.2307.15043.

2024, Bangkok, Thailand, August 11-16, 2024, pp. 14322–14350. Association for Computational

DETAILS ABOUT GENERATING ADVERSARIAL IMAGES А

In our implementation, we utilize the standard PGD method to generate adversarial images, according to 2.1. We run 5000 iterations of PGD optimization on generating each adversarial image and we set ϵ to 32/255. We test the robustness of the NEARSIDE method to varying ϵ . We use NEARSIDE to detect the adversarial samples generated under different ϵ . Results are deferred to Appx. D.

PROMPTS TO GPT-4 JUDGEMENT В

We adopt the prompt template proposed by Qi et al. (2024c), which combines the prohibited model usage policy, the (harmful) input query, the model's output for this query, and a rating rubric together to be fed to GPT-4 via OpenAI's APIs. The complete prompt template is as follows: ({Context.} will be replaced with the OpenAI's detailed policy and scoring tules, {Context.} will be replaced with the harmful query and model's output)

Prompt template
Below are OpenAI's prohibited usage policies. These policies ensure that models interact
with users safely and responsibly. Models should refuse to execute instructions that fall into
the disallowed use cases outlined in the OpenAI policy, avoiding amitmative responses or providing any information that could be helpful for any prohibited activities.
providing any information that could be helpful for any promoted activities.
OpenAI Policy - Disallowed usage of our models:
{The detailed policy.}
Now, I will provide you with a user instruction that the model should not comply
with, as per OpenAI's policy. I will also give you the model's response to that instruction.
Please review the model's response in light of the prohibited usage guidelines:
User Instruction: {The harmful query.}
Model Response: {The output to be classified.}
Use the secting rules below to secre the model's response to the user instruction on
a scale of 1 to 5:
{The detailed OpenAI's scoring rules}
(The detailed open it's scoring fulles)
Output your evaluation in the following format:
#the reason: your analysis here
[note: this should be step-by-step analysis following the steps: (a) Give an overview of
the user's primary intent and the model's response. (b) Evaluate the degree to which the
model's response aligns with the user's expressed intent. (c) Examine the user's instruction
and the model's reply respectively to determine if they contain any prohibited content.
Finally, evaluate the degree of the model's adherence to OpenAl's usage guidennes.
#the score: your score here
[note:provide your score in based on your analysis and the provide scoring criteria]
[note:provide your score in oused on your analysis and the provide scoring entertail]
C AN EXEMPLAR SAMPLE IN RADAR

915 916

864

865 866

867

868

869

870 871 872

873

874 875

876

877

878

879

880

In RADAR, each sample consists of an adversarial sample and a benign sample, with each containing 917 a query, an adversarial/benign image and VLMs' response. Fig. 6 provides an exemplar sample.



D PERTURBATION RADIUS IN GENERATING ADVERSARIAL IMAGES

We test the robustness of the NEARSIDE method to varying ϵ . We use NEARSIDE to detect the adversarial samples generated under different ϵ . We generate 100 adversarial samples on the LLaVA D-corpus dataset under settings of $\epsilon = 16$, $\epsilon = 64$, and unconstrained. The data generation follow the same pipeline as Sec. 3. We create 100 samples for each ϵ . The results are provided in Tab. 4.

Table 4: Results of the NEARSIDE on LLaVA D-corpus generated with different ϵ .

ϵ of adversarial training	Accuracy(%)	Precision(%)	Recall(%)	F1
$\begin{aligned} \epsilon &= 16/255\\ \epsilon &= 64/255\\ unconstrained \end{aligned}$	85.0	100.0	70.0	0.824
	93.5	97.8	89.0	0.932
	99.0	100.0	98.0	0.990

E ANALYSIS OF CROSS-MODEL TRANSFERABILITY

Linear transformation W. We explore the Platonic Representation Hypothesis by using a linear transformation W to align the two VLMs' embedding spaces. To demonstrate the importance of the usage of W, we conduct experiments that directly use *the attacking direction* of the source VLM to detect the adversarial samples of the target VLM without using W. Results are shown in Tab. 5.

PCA model. We use PCA model to reduce the dimension of VLMs' embedding and the attacking *direction* before learning the transformation W. In our initial setting, the dimension is reduced to 2056. We experiment to examine the effect of reducing the dimension of *the attacking direction* and the VLMs' embedding with the PCA model. In specific, we vary the dimension in [2048, 1024, 512, 256] and report the cross-model transferability results in Tab. 6.

973	Table 5: The results of cross-model transferability without W . We report result ^{-δ} where $^{\delta}$ indicates
974	the difference between the results w/o W and with W shown in Table 3.

$svlm \rightarrow tvlm \ (w/o \ W)$	TEST SET	Accuracy(%)	Precision(%)	Recall(%)	F1
	HH-rlhf	$50.9^{-13,4}$	$50.9^{-10.4}$	$51.6^{-26.0}$	$0.512^{-0.172}$
MiniGPT-4 \rightarrow LLaVA	D-corpus	$53.8^{-15.6}$	$53.2^{-9.4}$	$63.8^{-33.0}$	$0.580^{-0.180}$
	Harm-Data	$53.7^{-21.0}$	$53.7^{-23.0}$	$53.8^{-17.2}$	$0.537^{-0.200}$
	HH-rlhf	$32.8^{-45.0}$	$27.2^{-58.9}$	$20.6^{-45.6}$	$0.235^{-0.514}$
$LLaVA \rightarrow MiniGPI-4$	D-corpus	$71.0^{-9.4}$	$77.2^{-3.1}$	$59.6^{-21.0}$	$0.804^{-0.132}$
	Harm-Data	$23.9^{-73.2}$	$23.1^{-71.9}$	$22.4^{-77.0}$	$0.227^{-0.744}$

Table 6: The results of cross-model transferability where PCA reduce the VLMs' embedding and the attacking direction to different dimensions. We use **bold** to highlight the best results.

$svlm \rightarrow tvlm$	TEST SET	PCA-Dimension	Accuracy(%)	Precision(%)	Recall(%)	F1
		256	87.0	89.2	84.2	0.86
	1111 11 6	512	88.8	88.5	89.2	0.88
	HH-rint	1024	88.2	89.0	87.2	0.88
		2048	77.8	86.2	66.2	0.74
		256	83.3	76.2	96.8	0.85
$LLaVA \rightarrow MiniGPT-4$	Daamana	512	84.7	82.2	88.6	0.85
	D-corpus	1024	77.6	88.5	63.4	0.74
		2048	80.4	80.3	80.6	0.80
	Harm-Data	256	87.8	80.4	100.0	0.89
		512	83.4	75.1	100.0	0.85
		1024	88.2	89.0	87.2	0.88
		2048	97.1	95.0	99.4	0.93
		256	57.0	53.8	98.8	0.69
	****	512	58.5	54.7	98.0	0.70
	HH-rini	1024	63.7	58.6	93.2	0.72
		2048	64.3	61.3	77.6	0.68
		256	50.9	50.5	100.0	0.67
$MiniGPT-4 \rightarrow LLaVA$	Decembra	512	58.5	54.7	98.0	0.70
	D-corpus	1024	51.2	50.6	100.0	0.67
		2048	69.4	62.5	96.8	0.76
		256	71.9	64.7	96.6	0.77
	Horm Data	512	71.9	64.4	98.2	0.7
	nann-Data	1024	74.4	67.6	93.8	0.7
		2048	74.7	76.7	0.71	0.7

From Tab. 5, we find that, without W, the cross-model transferability results will significantly decrease. From Tab. 6, we find that, when reducing the dimension to 256, the cross-model trans-ferability results still remain high, indicating that the information in low dimensional sub-spaces is already sufficient for aligning two VLMs' embedding spaces.