
DENOISER: Rethinking the Robustness for Open-Vocabulary Action Recognition

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 As one of the fundamental video tasks in computer vision, Open-Vocabulary Action
2 Recognition (OVAR) has recently gained increasing attention, with the develop-
3 ment of vision-language pre-trainings. To enable open-vocabulary generalization,
4 existing methods formulate vanilla OVAR to evaluate the embedding similarity
5 between visual samples and text descriptions. However, one crucial issue is com-
6 pletely ignored: the text descriptions given by users may be noisy, *e.g.*, misspellings
7 and typos, limiting the real-world practicality. To fill the research gap, this paper
8 analyzes the noise rate/type in text descriptions by full statistics of manual spelling;
9 then reveals the poor robustness of existing methods; and finally rethinks to study
10 a practical task: noisy OVAR. One novel *DENOISER* framework, covering two
11 parts: generation and discrimination, is further proposed for solution. Concretely,
12 the generative part denoises noisy text descriptions via a decoding process, *i.e.*,
13 proposes text candidates, then utilizes inter-modal and intra-modal information to
14 vote for the best. At the discriminative part, we use vanilla OVAR models to assign
15 visual samples to text descriptions, injecting more semantics. For optimization, we
16 alternately iterate between generative-discriminative parts for progressive refine-
17 ments. The denoised text descriptions help OVAR models classify visual samples
18 more accurately; in return, assigned visual samples help better denoising. We carry
19 out extensive experiments to show our superior robustness, and thorough ablations
20 to dissect the effectiveness of each component.

21 1 Introduction

22 Action recognition is one of the fundamental tasks in computer vision that involves classifying videos
23 into meaningful semantics. Despite huge progress that has been made, existing researches focus more
24 on closed-set scenarios, where action classes remain constant during training and inference. Such
25 scenarios are an oversimplification of real life, and thus limiting their practical application. Recently,
26 another line of research considers one more challenging scenario, namely open-vocabulary action
27 recognition (OVAR), and receives increasing attention.

28 OVAR allows users to give free texts to describe action classes, and the model needs to match novel
29 (unseen) text descriptions to videos with similar semantics. To tackle OVAR task, Vision-Language
30 Alignment (VLA) paradigm [41, 14, 57] provides one preliminary but popular idea, *i.e.*, measuring
31 the embedding similarity between text descriptions and video embeddings. Following this paradigm,
32 recent works focus on minor improvements, *e.g.*, better align vision-language modalities [16, 49, 62].
33 Although promising, these works all maintain one unrealistic assumption in real-world scenarios, *i.e.*,
34 the given text descriptions are absolutely clean/accurate. The concrete form is that they evaluate open-
35 vocabulary performance by re-partitioning closed-set datasets in which text descriptions of classes are
36 fully human-checked. But in fact, under real-world OVAR, novel text descriptions provided by users
37 are sometimes noisy. Character misspellings (typos, missing, tense error) are inevitable [43, 25] in

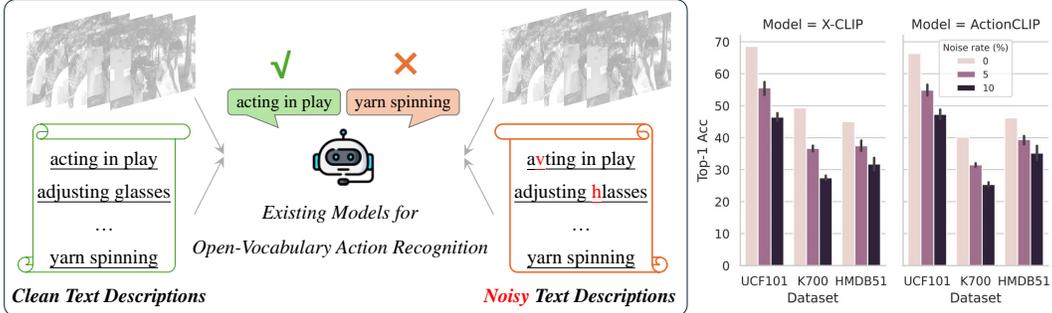


Figure 1: **Left:** For open-vocabulary action recognition (OVAR), existing researches neglect an essential aspect: the text descriptions provided by users may be noisy (*e.g.*, misspelling and typos), resulting in potential classification errors and limiting the real-world practicality. **Right:** Rethinking the robustness for popular OVAR methods [49, 62]. On various datasets, they exhibit high sensitivity to text noises. Besides, as the noise level increases, the performance degrades significantly.

38 thousands of descriptions, since users often don’t double-check, as well as differences in user habits
 39 and diversity of scenarios (Fig. 1 Left).

40 We are hence motivated to fill the research gap of noisy text descriptions in OVAR. We analyze the
 41 noise rate/type in real-world corpora [26, 45, 3]. We also make comprehensive simulations of text
 42 noises, following NLP literature [42, 47]. Fig. 1 Right empirically evaluates noise hazards for existing
 43 OVAR methods [16, 49, 62]. One can find that just a small amount of noise lowers recognition
 44 accuracy by a large margin, implying quite poor robustness.

45 To spur the community to deal with the noisy OVAR task, being necessary and practical, this paper
 46 bravely faces the challenges. One vanilla idea is using a separate language model (*e.g.*, GPT [1]) to
 47 correct noisy class descriptions, and then adapt the off-the-shelf vision-language paradigm [41, 14, 57].
 48 However, there exist two nettlesome issues. 1) *Textual Ambiguity*. One text description is usually a few
 49 compact words, with vague semantics, *e.g.*, for the noisy text “boird”, there could be multiple cleaned
 50 candidates in terms of spelling, such as “bird” and “board”. This short text lacks context, making
 51 phrase correction difficult for uni-modal language models. 2) *Cascaded Errors*. Text correction and
 52 action recognition are independently completed, without sharing knowledge. The noisy output of
 53 text correction is cascaded to the input of action recognition, resulting in continuous propagation of
 54 errors. To address these issues, we design one multi-modal robust framework: *DENOISER*.

55 Our first insight is to treat denoising of text descriptions as one *generative* task: given noisy text
 56 descriptions, decode the clean ones, by considering text-vision information to help denoising. Specif-
 57 ically, it consists of three components: text proposals, inter-modal weighting, and intra-modal
 58 weighting. We first propose potential text candidates based on spelling similarity to limit the decoding
 59 space. Then, two types of weighting are combined to decide the best candidate, that is, inter-modal
 60 weighting uses assigned visual samples to vote; while intra-modal weighting relies solely on text
 61 information. Our other insight is employing existing OVAR models as off-the-shelf tools to assign
 62 visual samples at *discriminative* step. Such tools have been proven to handle clean OVAR tasks well,
 63 also making our framework easier to adapt to previous models. For full usage of information in
 64 the same semantics, we then assign detail-rich visual samples to clarify the semantic ambiguity of
 65 compact text descriptions. To further avoid cascaded errors, we propose a solution of alternating
 66 iterations, to connect *generative* and *discriminative* steps. By progressive refinement, denoised text
 67 descriptions help OVAR models to match visual samples more accurately; assigned visual samples
 68 help better denoising. Under multiple iterations, denoising results and OVAR are both better.

69 Our main contributions are summarized as follows:

- 70 • We pioneer to explore noisy text descriptions for open-vocabulary action recognition (OVAR): first
 71 fully analyze the noise rate/type in text descriptions by extensive statistics in real-world corpora; then
 72 evaluate the robustness for existing methods; finally rethink to study one practical task: noisy OVAR.
- 73 • We propose a novel *DENOISER* framework to tackle the noisy OVAR task, by alternately optimizing
 74 generative-discriminative steps. The generative step leverages knowledge of vision-text alignment to
 75 denoises noisy text descriptions, in the form of progressive decoding; while the discriminative step
 76 assigns visual samples to text descriptions for open-vocabulary action recognition.

77 • We carry out extensive experiments to show the superior robustness of *DENOISER* against noisy
 78 text descriptions, under various noises and datasets. Great performance improvements are achieved
 79 over existing competitors. Thorough ablations are studied to show effectiveness of every design.

80 2 Related Work

81 **Vision-Language-Audio Pre-training (VLP)** aims to jointly optimize multi-modal embeddings with
 82 large-scale web data, *e.g.*, CLIP [41], ALIGN [14], Florence [57], FILIP [55], VideoCLIP [52], and
 83 LiT [58]. In architectures, VLP uses independent encoders for vision, text, and audio, followed by
 84 cross-modal fusion. For optimization, contrastive learning [5, 61] and cross-modal matching [7, 29]
 85 are mainstream, covering self supervision [32, 34], weak supervision [28, 8] and partial supervi-
 86 sion [19, 33]. VLP benefits various applications: image-text retrieval [6, 18], video understand-
 87 ing [23, 20, 22, 21], action recognition [16, 60], visual grounding [32, 56, 31], AIGC [4, 36].

88 **Open-Vocabulary Concept Learning** aims to understand vision, where conceptual semantics are
 89 described by free/arbitrary text descriptions. It is characterized by using vision-language pre-trainings
 90 to match text descriptions and visual samples in semantic space. Its typical evaluation metric is
 91 the downstream zero-shot performance, *i.e.*, classify unseen classes [49, 62, 17, 38, 54, 48, 37]. To
 92 achieve the evaluation, most methods re-partition closed-set datasets.[49] Although there is some
 93 plausibility, such re-partition implicitly makes an unrealistic assumption: text descriptions of unseen
 94 classes are human-checked, and thus absolutely clean, limiting real-world application. We pioneer
 95 taking noises from text descriptions (misspellings and typos) into consideration. By adding real-world
 96 noise for the above methods, we reveal their poor robustness, and design *DENOISER* for solution.

97 **Robustness of Language Models** is extensively studied by adversarial attack-defense techniques [50,
 98 59]. When text inputs are facing noises, defense methods correct the outputs, dividing into: detection-
 99 purification [63, 39], as well as adversarial training [53, 9, 35, 30, 51]. The former methods detect
 100 and correct the corrupted part of a text phrase. The latter trains a model on adversarial samples to
 101 increase its direct noise-against ability. Overall, all these methods employ solely textual information
 102 for robustness in pure NLP tasks. We differ from them by considering robustness in the context of
 103 multi-modal scenarios and by employing multi-modal information to better assist text denoising.

104 3 Method

105 We explore noisy text descriptions for open-vocabulary action recognition. In Sec 3.1, we introduce
 106 noisy open-vocabulary setting; in Sec 3.2, we detail our *DENOISER* framework, covering *generative*
 107 - *discriminative* sub-parts; in Sec 3.3, we report the accompanying optimization strategy.

108 3.1 Preliminary & Rethinking

109 **Open-Vocabulary Action Recognition (OVAR).** For a video dataset $\mathcal{V} = (v_j \in \mathbb{R}^{T \times H \times W \times 3})_j^N$,
 110 OVAR aims to train one model Φ_{OVAR} that matches target videos with arbitrary text description \mathcal{T} .

$$\mathcal{Y}^{\text{train}} = \Phi_{\text{OVAR}}(\mathcal{V}^{\text{train}}, \mathcal{T}^{\text{train}}) \in \mathbb{R}^{C_{\text{base}}}, \quad \mathcal{Y}^{\text{test}} = \Phi_{\text{OVAR}}(\mathcal{V}^{\text{test}}, \mathcal{T}^{\text{test}}) \in \mathbb{R}^{C_{\text{novel}}}, \quad (1)$$

111 where \mathcal{Y} refers to the matching label between \mathcal{V} and \mathcal{T} . During training, (video, text, matching label)
 112 triplets from the base semantic-classes are provided; while during testing, the model is evaluated
 113 on the novel semantic-classes. Note that, the semantic-classes between training (C_{base}) and testing
 114 (C_{novel}) are disjoint, *i.e.*, $C_{\text{base}} \cap C_{\text{novel}} = \emptyset$.

115 **Vision-Language Alignment (VLA).** To enable open-vocabulary capability, recent OVAR stud-
 116 ies [16, 49, 62, 40] embrace vision-language pre-trainings (VLPs), for their notable ability in cross-
 117 modal alignment. Specifically, OVAR could be achieved by measuring the embedding similarity
 118 between text descriptions \mathcal{T} and video samples \mathcal{V} , which is formally formulated as:

$$\mathcal{Y} = \sigma(\mathcal{F}_v * \mathcal{F}_t), \quad \mathcal{F}_v = \Phi_{\text{pool}}(\Phi_{\text{vis}}(\mathcal{V})) \in \mathbb{R}^{N \times D}, \quad \mathcal{F}_t = \Phi_{\text{txt}}(\mathcal{T}) \in \mathbb{R}^{C \times D}. \quad (2)$$

119 where σ refers to the softmax activation, Φ_{pool} is the spatio-temporal pooling, Φ_{vis} and Φ_{txt} are
 120 visual and textual encoders of VLPs, D is the embedding dimension.

121 **Noisy Text Descriptions in OVAR.** Although great progress has been made, the VLA paradigm
 122 suffers from an unrealistic assumption, *i.e.*, that text descriptions are absolutely clean/accurate,

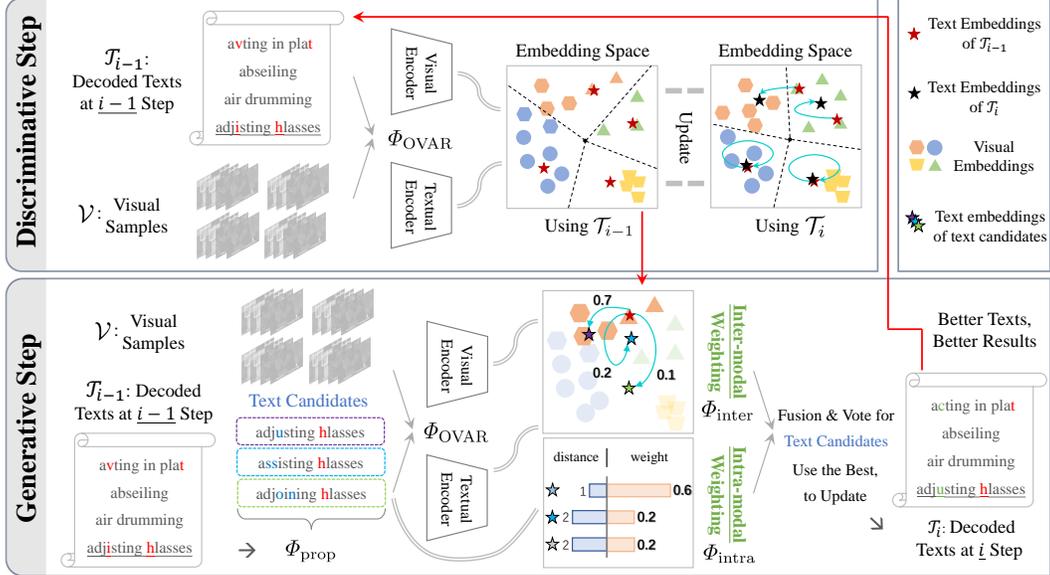


Figure 2: **Framework Overview.** *DENOISER* is composed of one *generative* part Ψ_{gene} and one *discriminative* part Ψ_{disc} . Ψ_{gene} views denoising text descriptions as a decoding process $\mathcal{T}_{i-1} \rightarrow \mathcal{T}_i$. We first propose text candidates Φ_{prop} for \mathcal{T}_{i-1} based on spelling similarity; then choose the best candidate by inter-modal weighting Φ_{inter} and intra-modal weighting Φ_{intra} . Φ_{inter} uses vision-text information, while Φ_{intra} relies solely on texts. Ψ_{disc} assigns text semantics to visual samples, then only visual samples with the same semantics can vote for text candidates. We optimize alternately between *generative* and *discriminative* steps to tackle noisy OVAR.

123 limiting the practicality in reality. Actually, the diversity of users and scenarios can easily cause
 124 text descriptions given to be somewhat noisy, especially for unseen semantic-classes, due to their
 125 enormous degree of freedom. Formally, for one text description with n words, the clean/noisy
 126 versions \mathcal{T}/\mathcal{T}' are:

$$\mathcal{T}' = (t'_1, \dots, t'_n) = \Psi_{\text{noise}}(\mathcal{T}; p), \quad \mathcal{T} = (t_1, \dots, t_n). \quad (3)$$

127 where t_i is the i -th word of \mathcal{T} . Ψ_{noise} refers to noise contamination in reality, *e.g.*, *inserting*, *substitut-*
 128 *ing* and *deleting* characters with probability p , following [42, 47]. Since these three atomic operations
 129 defined in Levenshtein edit distance \mathcal{D} are of distance 1, noise rate p can also be deduced by:

$$p = \frac{\mathcal{D}(\mathcal{T}, \mathcal{T}')}{\max(\text{length of } \mathcal{T}, \text{length of } \mathcal{T}')} \quad (4)$$

130 As a result, the noisy OVAR task can be formulated as: given \mathcal{V} and \mathcal{T}' , the model is expected to
 131 maximize the accuracy of action recognition, and even recovering \mathcal{T}' to \mathcal{T} .

132 **Robustness of Existing Methods.** Fig. 1 evaluates for typical OVAR studies [49, 62], across three
 133 public datasets. In terms of Top-1 classification accuracy, existing methods are rather sensitive to
 134 noise and show one trend: the larger the noise, the more significant the performance degradation
 135 (please see quantitative experiments in Tab. 2). Such poor robustness to the noisy OVAR task, proves
 136 excessive idealization of existing studies and also motivates us to fill the research gap.

137 3.2 *DENOISER*: One Robust OVAR Framework

138 **Motivation.** Given the complexity of noisy OVAR, we here divide it into two sub-steps: denoising of
 139 text descriptions, and then vanilla OVAR. The former is viewed as one *generative* decoding form, by
 140 considering both vision-text information for progressive denoising. While the latter is in one natural
 141 *discriminative* form, by assigning text descriptions to video samples. For the joint optimization of
 142 these two sub-steps, we iterate alternately between *generative* and *discriminative* forms. As a result,
 143 our *DENOISER* framework progressively tackles the noisy OVAR task.

144 **Framework.** As shown in Fig. 2, our *DENOISER* framework covers two components: *generative*
 145 sub-step Ψ_{gene} and *discriminative* sub-step Ψ_{disc} . For Ψ_{gene} , we iteratively refine text descriptions
 146 by one decoding process, that is, $(\mathcal{T}_0, \mathcal{T}_1, \dots, \mathcal{T}_n)$, where n is the index of decoding steps. Upon
 147 finishing step i , we will have $\mathcal{T}_i = (t_1, \dots, \bar{t}_i, t'_{i+1}, \dots, t'_n)$, where \bar{t} refers to the decoded version
 148 of t , meaning that the i -th word of text descriptions is decoded at step i . We start with $\mathcal{T}_0 = \mathcal{T}'$, and
 149 finish at \mathcal{T}_n to ensure that all words are denoised. While for Ψ_{disc} , we find it identical to vanilla OVAR
 150 task and thus leveraging the VLA pipeline [16, 49] for help, which is off-the-shelf and well-studied.
 151 Formally, our *DENOISER* framework tackles noisy OVAR as follows:

$$\mathcal{T}_i = \Psi_{\text{gene}}(\mathcal{T}_{i-1}, \mathcal{Y}_{i-1}, \mathcal{V}), \quad \mathcal{Y}_{i-1} = \Psi_{\text{disc}}(\mathcal{T}_{i-1}, \mathcal{V}) = \Phi_{\text{OVAR}}(\mathcal{T}_{i-1}, \mathcal{V}). \quad (5)$$

152 At the *discriminative* step, we calculate the matching label \mathcal{Y}_{i-1} to make coarse semantic classification
 153 of visual samples, *i.e.*, assign \mathcal{T}_{i-1} to \mathcal{V} . At the *generative* step, we first propose K text candidates
 154 $\Phi_{\text{prop}}(\mathcal{T}_{i-1})$ for \mathcal{T}_i base on \mathcal{T}_{i-1} to limit the decoding space. Then, to vote for the best candidate, we
 155 design two novel modules, namely inter-modal weighting Φ_{inter} and intra-modal weighting Φ_{intra} .
 156 Here, Φ_{inter} uses vision information \mathcal{V} , while Φ_{intra} relies on text information \mathcal{T}_{i-1} .

157 We alternate between the *generative* and *discriminative* steps to optimize the decoding result step by
 158 step. Please find in Algorithm 1 for comprehensive details.

159 3.3 Optimization for the *DENOISER* Framework

160 **Discriminative Step** consists in calculating cross-modal matching labels \mathcal{Y} using Ψ_{disc} . Intuitively,
 161 visual samples \mathcal{V}_c whose labels \mathcal{Y} are assigned to semantic-class c , *i.e.* $\text{argmax } \mathcal{Y} = c$, are those who
 162 could help decode $\mathcal{T}_{c,i}$ most efficiently. On the contrary, visual samples from other semantic-classes
 163 may have few connections with the current class and thus provide no meaningful aid. Here, we find
 164 this process is identical to vanilla OVAR, and hence employs Φ_{OVAR} as Ψ_{disc} . We theoretically
 165 prove in the Appendix that, \mathcal{V}_c is the best set of visual samples to choose from. With \mathcal{V}_c defined and
 166 $\text{argmax } \mathcal{Y} = c$, Ψ_{gene} decodes text descriptions $\mathcal{T}_{c,i}$ for each semantic-class c :

$$\Psi_{\text{gene}}(\mathcal{T}_{c,i-1}, \mathcal{Y}, \mathcal{V}) = \Psi_{\text{gene}}(\mathcal{T}_{c,i-1}, \mathcal{V}_c) = \underset{\mathcal{T}_{c,i}}{\text{argmax}} p(\mathcal{T}_{c,i} | \mathcal{T}_{c,i-1}, \mathcal{V}_c). \quad (6)$$

167 Recall $t_{c,i}$ is the i -th word to be decoded, and $\mathcal{T}_{c,i-1}$ is from last decoding, with the first $i-1$
 168 words decoded. As we decode word-by-word, choosing the best $\mathcal{T}_{c,i}$ is exactly choosing the best $t_{c,i}$,
 169 *i.e.* $\text{argmax}_{\mathcal{T}_{c,i}} p(\mathcal{T}_{c,i} | \mathcal{T}_{c,i-1}, \mathcal{V}_c) = \text{argmax}_{t_{c,i}} p(t_{c,i} | \mathcal{T}_{c,i-1}, \mathcal{V}_c)$, as we do in *generative* step.

170 **Generative Step** here consists in, for each semantic-class c , choosing the best $t_{c,i}$ that maximizes
 171 $p(t_{c,i} | \mathcal{T}_{c,i-1}, \mathcal{V}_c)$. With $p(\mathcal{T}_{c,i-1}, \mathcal{V}_c)$ and $p(\mathcal{V}_c)$ same for all possible $t_{c,i}$, we make detailed deriva-
 172 tions in the Appendix to show that:

$$p(t_{c,i} | \mathcal{T}_{c,i-1}, \mathcal{V}_c) \propto p(t_{c,i}, \mathcal{T}_{c,i-1}, \mathcal{V}_c) \propto \prod_{v_j \in \mathcal{V}_c} p(t_{c,i} | v_j) p(\mathcal{T}_{c,i-1} | t_{c,i}, v_j). \quad (7)$$

173 Here, the error model $p(\mathcal{T}_{c,i-1} | t_{c,i}, v_j)$ evaluates how $t_{c,i}$ may be misspelled as $t'_{c,i}$, since the i -th
 174 word in $\mathcal{T}_{c,i-1}$ is still noisy and not decoded. Knowing that errors in text descriptions are independent
 175 of visual samples, it reduces to uni-modal $p(\mathcal{T}_{c,i-1} | t_{c,i})$. As the error that one may make given the
 176 correct text is harder to model while the reverse is much easier, we let $p(\mathcal{T}_{c,i-1} | t_{c,i}) \propto p(t_{c,i} | \mathcal{T}_{c,i-1})$.
 177 Please refer to detailed derivations in the Appendix. As a result, our final objective is:

$$p(t_{c,i} | \mathcal{T}_{c,i-1}) \prod_{v_j \in \mathcal{V}_c} p(t_{c,i} | v_j) = \Phi_{\text{intra}} \prod_{v_j \in \mathcal{V}_c} \Phi_{\text{inter}}. \quad (8)$$

178 **Text Proposals** consists in proposing K candidates $\{t_i^k\}_k$ for t_i with the lowest Levenshtein Edit
 179 Distance $\mathcal{D}(\cdot, t'_i)$ (a metric of spelling similarity). By replacing original noisy word t'_i in \mathcal{T}_{i-1}^k with
 180 $\{t_i^k\}_k$, they form $\Phi_{\text{prop}}(\mathcal{T}_{i-1}) = \mathcal{T}_i^k = (\bar{t}_1, \dots, \bar{t}_{i-1}, t_i^k, t'_{i+1}, \dots, t'_n)$, the K candidates for \mathcal{T}_i .
 181 The benefit of text proposals is to reduce computing complexity. Since text embeddings are quantized
 182 in the semantic space, the search is limited to proposed candidates, rather than in the entire space.

183 **Inter-modal Weighting** $\Phi_{\text{inter}} = p(t_{c,i} | v_j)$, $v_j \in \mathcal{V}_c$ relies on vision samples from semantic-class c
 184 to determine the best $t_{c,i}$ for the next iteration. Concretely, we model the probability of being chosen

Algorithm 1 *DENOISER*: Robust Open-Vocabulary Action Recognition

Require: noisy text descriptions \mathcal{T}' , visual samples \mathcal{V} , iteration number n , temperature λ , candidate number K , edit distance \mathcal{D} , open-vocabulary model Φ_{OVAR}

```
 $\mathcal{T}_0 \leftarrow \mathcal{T}'$   
for  $i = 1, 2, \dots, n$  do  
  for  $c = 1, 2, \dots, C$  do ▷ Text Proposals  
     $t'_{c,i}$  is the  $i$ -th word of  $\mathcal{T}_{c,i-1}$ , which is noisy and not yet decoded  
    Select from corpus,  $K$  candidates  $\{t'_{c,i}\}_k$  with the smallest  $\mathcal{D}$  with  $t'_{c,i}$   
    Replace  $t'_{c,i}$  with  $\{t'_{c,i}\}_k$ , forming  $\{\mathcal{T}_{c,i}\}_k$   
  end for  
  for  $j = 1, 2, \dots, |\mathcal{V}|$  do ▷ Discriminative Step  
     $c \leftarrow \underset{c}{\operatorname{argmax}} \max_k \frac{\exp(\mathcal{S}(v_j, \mathcal{T}_{c,i}^k))}{\sum_{k'} \exp(\mathcal{S}(v_j, \mathcal{T}_{c,i}^{k'}))}$   
    Assign  $v_j$  to class  $c$ ,  $v_j \in \mathcal{V}_c$   
  end for  
  for  $c = 1, 2, \dots, C$  do ▷ Generative Step  
     $\Phi_{\text{intra}}^k \leftarrow \frac{\exp(-\mathcal{D}(t_{c,i}^k, t'_{c,i})/\lambda)}{\sum_{k'} \exp(-\mathcal{D}(t_{c,i}^{k'}, t'_{c,i})/\lambda)}$  ▷ Intra-Modal Weighting  
     $\Phi_{\text{inter}}^k \leftarrow \prod_{v_j \in \mathcal{V}_c} \frac{\exp(\mathcal{S}(v_j, \mathcal{T}_{c,i}^k))}{\sum_{k'} \exp(\mathcal{S}(v_j, \mathcal{T}_{c,i}^{k'}))}$  ▷ Inter-Modal Weighting  
     $\mathcal{T}_{c,i} \leftarrow \mathcal{T}_{c,i}^k, k = \operatorname{argmax}_k \Phi_{\text{intra}}^k \times \Phi_{\text{inter}}^k$   
  end for  
end for
```

185 for each proposed candidate to be:

$$\mathbb{P}(t_{c,i} = t'_{c,i} | v_j) = \mathbb{P}(\mathcal{T}_{c,i} = \mathcal{T}_{c,i}^k | v_j) = \frac{\exp(\mathcal{S}(v_j, \mathcal{T}_{c,i}^k))}{\sum_{k'} \exp(\mathcal{S}(v_j, \mathcal{T}_{c,i}^{k'}))}, v_j \in \mathcal{V}_c. \quad (9)$$

186 where $\mathcal{S}(\cdot, \cdot)$ is the cosine similarity between video-text embeddings, both encoded by Φ_{OVAR} . The
187 intuition is that the more unanimously visual samples agree on candidate $\mathcal{T}_{c,i}^k$, the more likely it is the
188 text descriptions corresponding to semantic-class c . Besides, by letting visual samples vote on $\mathcal{T}_{c,i}^k$
189 instead of $t'_{c,i}$, we take into consideration not only the current word $t_{c,i}$ but also context implicitly.

190 *Intra-modal Weighting* $\Phi_{\text{intra}} = p(t_{c,i} | \mathcal{T}_{c,i-1})$ relies solely on text information to decide the best $t_{c,i}$
191 for next iteration. Although Φ_{intra} may be solved by uni-modal spell-checkers [15] or large language
192 models [1], we here design a simple model by considering only spelling similarity (ignore contexts),
193 to save computing costs. That is, choose $t_{c,i}$ depending solely on $t'_{c,i}$ instead of on entire $\mathcal{T}_{c,i-1}$:

$$\mathbb{P}(t_{c,i} = t'_{c,i} | \mathcal{T}_{c,i-1}) = \mathbb{P}(t_{c,i} = t'_{c,i} | t'_{c,i}) = \frac{\exp(-\mathcal{D}(t_{c,i}^k, t'_{c,i})/\lambda)}{\sum_{k'} \exp(-\mathcal{D}(t_{c,i}^{k'}, t'_{c,i})/\lambda)}. \quad (10)$$

194 The intuition is that, the more similar a word candidate $t_{c,i}^k$ is, compared to the noisy word $t'_{c,i}$, the
195 more likely it is the corresponding denoised word. Here, we introduce one temperature parameter λ to
196 balance Φ_{intra} and Φ_{inter} . A larger λ indicates that different edit distance gives similar probabilities,
197 meaning that we rely more on visual samples for decision, and vice versa.

198 4 Experiments

199 **Typical Models for Vanilla OVAR.** To illustrate the generalizability of our framework, we leverage
200 two typical models from the VLA pipeline as Φ_{OVAR} , that is, ActionCLIP [49] and XCLIP [62].
201 These two models adopt hand-crafted prompts and visual-conditioned prompt tuning, respectively.
202 Under both models, we choose ViT-B/16-32F as the network backbones, for simplicity.

203 **Datasets.** HMDB51 [26] contains 7k videos covering 51 action categories. UCF101 [45] contains
204 13k videos spanning 101 action categories. Kinetics700 [3] (K700) is simply an extension of K400,
205 with around 650k video clips sourced from YouTube. To partition these datasets for open-vocabulary
206 action recognition, this paper follows the standard consensus [49, 62], for the sake of fairness.

Figure 3: **Statistics for Noises in Reality.** Text noises may be classified into 4 types: inserting, substituting, swapping, and deleting characters.[2] In terms of edit distance, based on TOEFL-Spell dataset[10], most of the text noises have an edit distance = 1 compared to the clean version. Nevertheless, the distribution tends to be positively skewed towards larger noise.

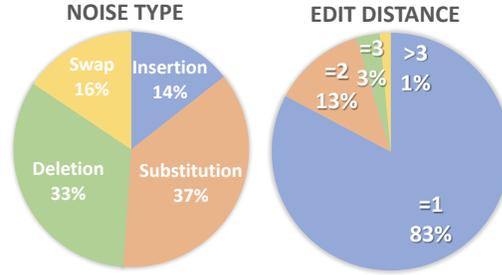


Table 1: **Comparisons between Various Competitors.** Using ActionCLIP [49] as Φ_{OVAR} while evaluating on UCF101, we compare with statistical text spell-checkers (PySpellChecker [15]), neural based ones (Bert from NeuSpell) [13], and GPT 3.5 [1]. Our method remarkably outperforms others in terms of Top-1 classification accuracy, and semantic similarity of recovered text descriptions.

Noise Type	Noise Rate	Competitors	Top-1 Acc	Label Acc	Semantic Similarity
–	0%	Upper Bound	66.3	100	100
Real	~5.52%	GPT 3.5 [1]	61.2 \pm 1.4	74.7 \pm 1.9	97.1 \pm 0.4
		Bert (NeuSpell) [13]	56.0 \pm 1.1	64.7 \pm 2.0	94.5 \pm 0.4
		PySpellChecker [15]	59.9 \pm 1.2	79.6 \pm 1.6	96.7 \pm 0.3
		Ours	61.5\pm0.7	82.3\pm1.6	97.2\pm0.3
Simulated	5%	GPT 3.5 [1]	59.7 \pm 1.2	47.6 \pm 3.1	95.9 \pm 0.4
		Bert (NeuSpell) [13]	56.6 \pm 0.5	66.2 \pm 2.3	94.6 \pm 0.4
		PySpellChecker [15]	60.9 \pm 1.1	82.5 \pm 2.9	97.1 \pm 0.4
	Ours	63.8\pm0.7	86.4\pm2.3	97.7\pm0.2	
	10%	GPT 3.5 [1]	58.5 \pm 1.3	51.6 \pm 2.3	95.8 \pm 0.3
		Bert (NeuSpell) [13]	51.0 \pm 0.5	50.4 \pm 3.6	91.6 \pm 0.6
PySpellChecker [15]		55.7 \pm 1.1	69.3 \pm 1.5	94.8 \pm 0.3	
Ours	61.2\pm0.8	75.9\pm1.9	96.4\pm0.3		

207 **Metric.** We use three metrics for full evaluations from multiple perspectives. Top-1 Acc refers to
 208 the top-1 classification accuracy of noisy open-vocabulary action recognition. Label Acc counts the
 209 percentage of denoised text descriptions that match exactly with ground truth. Semantic Similarity
 210 calculates the cosine similarity of embeddings, between denoised and clean text descriptions. Label
 211 Acc and Semantic Similarity measure how well noisy text descriptions are recovered.

212 **Implementations.** We set the proposal number $K = 10$. Intra-modal weighting and inter-modal
 213 weighting are both used to determine the best candidate. Temperature λ follows a linear schedule
 214 from 0.01 to 1. We use the same corpus as in PySpellChecker, which contains 70317 English words,
 215 for text proposals. For typical OVAR methods [49, 62], we choose the ViT-B/16-32F checkpoint
 216 pretrained on K400 [24] to evaluate their zero-shot robustness on HMDB51 [27], UCF101 [46] and
 217 K700 [44]. Since K700 and K400 have overlapped categories, we exclude them when evaluating on
 218 K700. For UCF101, we use the separated lowercase text label. All ablation studies are conducted on
 219 UCF101 under 20% noise. For statistical significance, We do each simulation 10 times and report the
 220 mean and confidence interval of 95%. All experiments are done using a single RTX 3090.

221 4.1 Statistics on Noise Type/Rate for Text Descriptions

222 **Real Noise.** We adopt two large-scale corpora [11, 10] of misspellings to analyze noise type in text
 223 descriptions. As shown in Fig. 3, the conclusion is similar to the NLP community [42, 47], *i.e.*, three
 224 atomic types of noise are inserting, substituting, and deleting text characters. More complicated noise
 225 patterns, *e.g.* swapping, can be constructed by mixing atomic noise types. Then, following previous
 226 literature, we quantify noise rate through Levenshtein Edit Distance, a generally accepted metric,
 227 to calculate the occurrence number of atomic noise types. Specifically, GitHub Typo Corpus [11]
 228 contains over 350k edits of typos from GitHub. The average noise rate (per sentence) is 3.3%.
 229 Nevertheless, the distribution is highly positively skewed (skewness = 2.9). For the worst 5% cases,
 230 the noise rate (per sentence) is larger than 9.4%. TOEFL-Spell Corpus [10] samples essays written
 231 by candidates from various language backgrounds in TOEFL[®] iBT test. There are, on average, 6.9
 232 spelling mistakes per essay. For misspelled words, the noise rate (per word) is on average 16.0%.

Table 2: **Comparison Across Datasets and Models.** On three standard datasets, facing multiple noise types (real or simulated), and under various noise rates, our *DENOISER* consistently improves the performance for noisy OVAR, regardless of underlying OVAR methods Φ_{OVAR} .

Dataset	Noise Type	Noise Rate	Φ_{OVAR} : Typical Models for Vanilla OVAR task			
			ActionCLIP [49]		XCLIP [62]	
			w/o Ours	w Ours	w/o Ours	w Ours
UCF101	Upper Bound		66.3		68.6	
	Real	$\sim 5.52\%$	$54.0_{\pm 2.3}$	$61.5_{\pm 0.7}$	$53.8_{\pm 2.7}$	$63.4_{\pm 0.9}$
	Simulated	5%	$54.9_{\pm 1.8}$	$63.2_{\pm 0.7}$	$55.6_{\pm 2.2}$	$64.2_{\pm 1.4}$
10%		$47.3_{\pm 1.4}$	$61.2_{\pm 1.2}$	$46.4_{\pm 1.3}$	$62.9_{\pm 2.3}$	
HMDB51	Upper Bound		46.2		45.0	
	Real	$\sim 6.71\%$	$37.6_{\pm 1.6}$	$40.0_{\pm 1.4}$	$35.3_{\pm 1.5}$	$38.4_{\pm 1.4}$
	Simulated	5%	$39.4_{\pm 1.4}$	$41.3_{\pm 1.4}$	$37.5_{\pm 1.8}$	$39.7_{\pm 1.0}$
10%		$35.2_{\pm 2.3}$	$39.6_{\pm 1.4}$	$31.8_{\pm 2.2}$	$37.3_{\pm 1.5}$	
K700	Upper Bound		40.2		49.3	
	Real	$\sim 5.47\%$	$30.8_{\pm 0.51}$	$35.9_{\pm 0.4}$	$35.6_{\pm 0.6}$	$43.5_{\pm 0.7}$
	Simulated	5%	$31.5_{\pm 0.5}$	$36.8_{\pm 0.3}$	$36.7_{\pm 0.9}$	$44.1_{\pm 0.6}$
10%		$25.4_{\pm 0.8}$	$35.3_{\pm 0.5}$	$27.5_{\pm 0.7}$	$41.8_{\pm 0.9}$	

233 **Noise Scenarios.** In the "Simulated" noise type, we mix three atomic noises: insertion, substitution,
 234 and deletion. Concretely, for each character, we perturb it with probability p . For each perturbation,
 235 it will be insertion, substitution, and deletion with equal probability. To further ensure real-world
 236 generalizability, we ask GPT3.5 to give examples of perturbation according to real-world scenarios.
 237 We mix them into simulated noises. Noise rate p of the "Real" noise type is estimated with Eq. (3).

238 4.2 Comparison with State-of-the-art Methods

239 **Comparison to Competitors.** Tab. 1 compares from three axes: Top-1 Acc of Φ_{OVAR} after correction,
 240 Label Acc and Semantic Similarity. PySpellChecker is a uni-modal statistical model that corrects
 241 each word by edit distance and appearance frequency. Bert (NeuSpell) [13] employs a uni-modal
 242 Bert-based model to translate noisy text descriptions into clean ones. We also ask GPT 3.5 to denoise
 243 text descriptions using the prompt "The following words may contain spelling errors by deleting,
 244 inserting, and substituting letters. You are a corrector of spelling errors. Give only the answer
 245 without explication. What is the correct spelling of the action of <noisy text description>?". Our
 246 method outperforms all competitors by large margins, which is impressive because our method is
 247 unsupervised without prior knowledge other than those contained in the OVAR model. Note that the
 248 output of GPT 3.5 tends to be unstable depending on prompts, which requires manual cleaning to
 249 remove irrelevant parts contained in the output, thus impeding real-world usage.

250 **Comparisons Across Datasets/Models.** Tab. 2 compares Top-1 Acc to further reveal our solution is
 251 scalable/generalizable. Under various noise rates, our model is robust to achieve huge improvements.
 252 In terms of scalability across models, our method is not only applicable to hand-crafted prompts as in
 253 ActionCLIP but also to learnable visual-conditioned prompts as in XCLIP. Furthermore, we notice
 254 that, whenever XCLIP outperforms ActionCLIP, our method also yields a better result. A better
 255 visual encoder and well-tuned prompt may significantly increase our performance, showing that our
 256 method's upper limit could become higher, as the community continues to train better OVAR models.

257 4.3 Ablation Study

258 **Inter-modal Weighting Φ_{inter} & Intra-modal Weighting Φ_{intra} .** Tab. 3 shows that, both Φ_{inter}
 259 and Φ_{intra} contribute to denoising text descriptions and to improving the robustness of underlying
 260 Φ_{OVAR} . In terms of Top-1 Acc and Semantic Similarity, Φ_{inter} performs better than Φ_{intra} , since
 261 Φ_{inter} uses visual information as one direct optimization guideline to improve video recognition.
 262 While Φ_{intra} performs better in terms of Label Acc, which focuses more on spelling correctness.
 263 Besides, Φ_{inter} and Φ_{intra} turn out to be complementary: visual information helps to understand
 264 noisy text descriptions; while textual information prevents the model from being misled by visual
 265 samples. We achieve the best performance when combining these two weightings.

Table 3: **Ablations for Inter-modal Weighting Φ_{Inter} , Intra-modal Weighting Φ_{Intra} , Schedule of Temperature λ .** Φ_{Inter} alone outperforms Φ_{Intra} . Both contribute to correcting class texts, and give the best results when combined. Linear schedule of balancing factor λ outperforms the constant one, meaning that it helps to rely more on Φ_{Intra} at first, and then gradually switch to Φ_{Inter} .

	Φ_{Inter}	Φ_{Intra}	Schedule λ	Top-1 Acc	Label Acc	Semantic Similarity
A1		✓	/	48.1 \pm 2.2	38.2 \pm 2.5	88.9 \pm 0.4
A2	✓		/	52.9 \pm 1.4	34.1 \pm 2.4	89.1 \pm 0.6
A3	✓	✓	Constant	54.5 \pm 2.5	54.9 \pm 4.5	92.4 \pm 0.8
A4	✓	✓	Linear	55.2\pm1.5	55.1\pm3.0	92.9\pm0.6

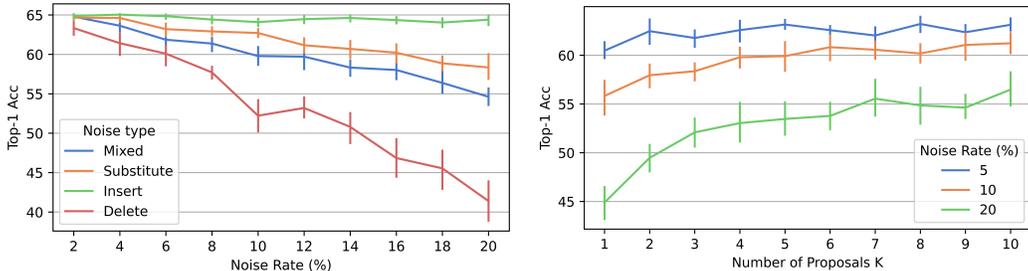


Figure 4: We evaluate on UCF101 by using ActionCLIP as Φ_{OVAR} . **Left: Ablation Study on Noise Type.** “Mixed” means that all types of text noises: “Substitute”, “Insert”, “Delete” take place with equal probability. Our *DENOISER* shows good resilience, especially against noises of inserting or substituting. **Right: Ablation Study on Proposal Number K .** As K increases, Top-1 Acc increases and converges gradually towards the upper bound, but it also brings heavier computing costs.

266 **Temperature Schedule λ** balances intra-modal weighting and inter-modal weighting. One larger λ
 267 indicates more reliance on inter-modal weighting. “Linear” means that λ augments from 0.01 to 1
 268 linearly. Tab. 3 reports that it is beneficial to rely more on intra-modal at the beginning of decoding,
 269 and then gradually turn to inter-modal for more help. This indicates that, when text noises are high,
 270 Φ_{intra} offers more help; when text noises are slight, Φ_{inter} could help more.

271 **Noise Type.** Fig. 4 Left reports our robustness under various noise types/rates. “Mixed” means that
 272 three noise types: “Substitute”, “Insert”, “Delete” are equally possible to appear. Our method shows
 273 remarkable resilience when texts are perturbed by inserting or substituting characters. Performance
 274 degradation is observed when texts are perturbed by deleting characters. It is reasonable, as deleting
 275 characters causes huge information loss, making the model difficult to recover clean text descriptions.

276 **Number of Candidates K .** Fig. 4 Right shows as K increases, inter-modal weighting can reveal
 277 its full power, hence improving performance. Otherwise, if a good candidate is excluded from the
 278 proposal stage due to a small K , it can be selected by neither of the inter- or intra-modal weighting,
 279 thus decreasing performance. Moreover, the performance tends towards one plateau, showing a
 280 decreasing marginal contribution of more proposals to performance. Since a larger K means more
 281 computing costs for text encoding, we select $K = 10$ by default to make reasonable trade-offs.

282 5 Conclusion

283 This paper investigates how noises in class-text descriptions negatively interference OVAR; and
 284 one novel framework *DENOISER* is proposed for solutions. By incorporating visual information
 285 during denoising, we clarify the ambiguity induced by short and context-lacking text descriptions; by
 286 iteratively refining the denoised output through one generative-discriminative process, we mitigate
 287 cascaded errors which may propagate from spell-checking models to outputs of OVAR model. We
 288 conduct extensive experiments to demonstrate the generalizability of *DENOISER* across multiple
 289 models and datasets, and also show our superiority over uni-modal spell-checking solutions.

290 **Limitations.** 1) We focus more on spelling noises; while in the real world, text noises can be more
 291 complex, involving semantic ambiguity. Equipping *DENOISER* with large language models may
 292 be a feasible solution. 2) Using more text candidates or visual samples brings better results for
 293 *DENOISER*, but also costs more. There is a trade-off between performance and computational cost.

294 **References**

- 295 [1] Gpt-3.5 turbo, <https://platform.openai.com/docs/models/gpt-3-5-turbo/>
- 296 [2] Al-Oudat, A.: Spelling errors in english writing committed by english-major students at bau.
297 *Journal of Literature, Languages and Linguistics* **32**(2) (2017)
- 298 [3] Carreira, J., Noland, E., Hillier, C., Zisserman, A.: A short note on the kinetics-700 human
299 action dataset. arXiv preprint arXiv:1907.06987 (2019)
- 300 [4] Chen, M., Chen, X., Zhai, Z., Ju, C., Hong, X., Lan, J., Xiao, S.: Wear-any-way: Manipulable
301 virtual try-on via sparse correspondence alignment. arXiv preprint arXiv:2403.12965 (2024)
- 302 [5] Chen, X., Chen, S., Yao, J., Zheng, H., Zhang, Y., Tsang, I.W.: Learning on attribute-missing
303 graphs. *IEEE transactions on pattern analysis and machine intelligence* (2020)
- 304 [6] Chen, X., Cheng, Z., Yao, J., Ju, C., Huang, W., Lan, J., Zeng, X., Xiao, S.: Enhancing
305 cross-domain click-through rate prediction via explicit feature augmentation. arXiv preprint
306 arXiv:2312.00078 (2023)
- 307 [7] Cheng, F., Wang, X., Lei, J., Crandall, D., Bansal, M., Bertasius, G.: Vindlu: A recipe for
308 effective video-and-language pretraining. In: *Proceedings of the IEEE Conference on Computer
309 Vision and Pattern Recognition* (2023)
- 310 [8] Cheng, Z., Xiao, S., Zhai, Z., Zeng, X., Huang, W.: Mixer: Image to multi-modal retrieval
311 learning for industrial application. arXiv preprint arXiv:2305.03972 (2023)
- 312 [9] Dinan, E., Humeau, S., Chintagunta, B., Weston, J.: Build it break it fix it for dialogue safety:
313 Robustness from adversarial human attack. arXiv preprint arXiv:1908.06083 (2019)
- 314 [10] Flor, M., Fried, M., Rozovskaya, A.: A benchmark corpus of english misspellings and a
315 minimally-supervised model for spelling correction. In: *Proceedings of the Fourteenth Workshop
316 on Innovative Use of NLP for Building Educational Applications*. pp. 76–86 (2019)
- 317 [11] Hagiwara, M., Mita, M.: Github typo corpus: A large-scale multilingual dataset of misspellings
318 and grammatical errors. arXiv preprint arXiv:1911.12893 (2019)
- 319 [12] Hu, X., Zhang, K., Xia, L., Chen, A., Luo, J., Sun, Y., Wang, K., Qiao, N., Zeng, X., Sun,
320 M., et al.: Reclip: Refine contrastive language image pre-training with source free domain
321 adaptation. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer
322 Vision*. pp. 2994–3003 (2024)
- 323 [13] Jayanthi, S.M., Pruthi, D., Neubig, G.: Neuspell: A neural spelling correction toolkit. arXiv
324 preprint arXiv:2010.11085 (2020)
- 325 [14] Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q.V., Sung, Y., Li, Z., Duerig,
326 T.: Scaling up visual and vision-language representation learning with noisy text supervision.
327 In: *Proceedings of the International Conference on Machine Learning* (2021)
- 328 [15] Jiang, Y.G., Liu, J., Zamir, A.R., Toderici, G., Laptev, I., Shah, M., Sukthankar, R.:
329 pypellchecker: Action recognition with a large number of classes, [https://github.com/
330 barrust/pypellchecker/](https://github.com/barrust/pypellchecker/)
- 331 [16] Ju, C., Han, T., Zheng, K., Zhang, Y., Xie, W.: Prompting visual-language models for efficient
332 video understanding. In: *Proceedings of the European Conference on Computer Vision*. Springer
333 (2022)
- 334 [17] Ju, C., Li, Z., Zhao, P., Zhang, Y., Zhang, X., Tian, Q., Wang, Y., Xie, W.: Multi-modal
335 prompting for low-shot temporal action localization. arXiv preprint arXiv:2303.11732 (2023)
- 336 [18] Ju, C., Wang, H., Li, Z., Chen, X., Zhai, Z., Huang, W., Xiao, S.: Turbo: Informativity-driven
337 acceleration plug-in for vision-language models. arXiv preprint arXiv:2312.07408 (2023)
- 338 [19] Ju, C., Wang, H., Liu, J., Ma, C., Zhang, Y., Zhao, P., Chang, J., Tian, Q.: Constraint and union
339 for partially-supervised temporal sentence grounding. arXiv preprint arXiv:2302.09850 (2023)
- 340 [20] Ju, C., Zhao, P., Chen, S., Zhang, Y., Wang, Y., Tian, Q.: Divide and conquer for single-frame
341 temporal action localization. In: *Proceedings of the International Conference on Computer
342 Vision* (2021)
- 343 [21] Ju, C., Zhao, P., Chen, S., Zhang, Y., Zhang, X., Wang, Y., Tian, Q.: Adaptive mutual supervision
344 for weakly-supervised temporal action localization. *IEEE Transactions on Multimedia* (2022)

- 345 [22] Ju, C., Zhao, P., Zhang, Y., Wang, Y., Tian, Q.: Point-level temporal action localization: Bridging
346 fully-supervised proposals to weakly-supervised losses. arXiv preprint arXiv:2012.08236 (2020)
- 347 [23] Ju, C., Zheng, K., Liu, J., Zhao, P., Zhang, Y., Chang, J., Tian, Q., Wang, Y.: Distilling vision-
348 language pre-training to collaborate with weakly-supervised temporal action localization. In:
349 Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2023)
- 350 [24] Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F.,
351 Green, T., Back, T., Natsev, P., et al.: The kinetics human action video dataset. arXiv preprint
352 arXiv:1705.06950 (2017)
- 353 [25] Keller, Y., Mackensen, J., Eger, S.: Bert-defense: A probabilistic model based on bert to combat
354 cognitively inspired orthographic adversarial attacks. arXiv preprint arXiv:2106.01452 (2021)
- 355 [26] Kuehne, H., Jhuang, H., Garrote, E., Poggio, T., Serre, T.: HMDB: A large video database
356 for human motion recognition. In: Proceedings of the International Conference on Computer
357 Vision (2011)
- 358 [27] Kuehne, H., Jhuang, H., Garrote, E., Poggio, T., Serre, T.: HMDB: a large video database for
359 human motion recognition. In: Proceedings of the International Conference on Computer Vision
360 (ICCV) (2011)
- 361 [28] Li, J., Li, D., Savarese, S., Hoi, S.: Blip-2: Bootstrapping language-image pre-training with
362 frozen image encoders and large language models. In: International conference on machine
363 learning. PMLR (2023)
- 364 [29] Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified
365 vision-language understanding and generation. In: International conference on machine learning.
366 pp. 12888–12900. PMLR (2022)
- 367 [30] Liu, H., Zhang, Y., Wang, Y., Lin, Z., Chen, Y.: Joint character-level word embedding and
368 adversarial stability training to defend adversarial text. In: Proceedings of the AAAI Conference
369 on Artificial Intelligence (2020)
- 370 [31] Liu, J., Ju, C., Ma, C., Wang, Y., Wang, Y., Zhang, Y.: Audio-aware query-enhanced transformer
371 for audio-visual segmentation. arXiv preprint arXiv:2307.13236 (2023)
- 372 [32] Liu, J., Ju, C., Xie, W., Zhang, Y.: Exploiting transformation invariance and equivariance
373 for self-supervised sound localisation. In: Proceedings of ACM International Conference on
374 Multimedia (2022)
- 375 [33] Liu, J., Liu, Y., Zhang, F., Ju, C., Zhang, Y., Wang, Y.: Audio-visual segmentation via unlabeled
376 frame exploitation. arXiv preprint arXiv:2403.11074 (2024)
- 377 [34] Liu, J., Wang, Y., Ju, C., Ma, C., Zhang, Y., Xie, W.: Annotation-free audio-visual segmentation.
378 In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision
379 (2024)
- 380 [35] Liu, K., Liu, X., Yang, A., Liu, J., Su, J., Li, S., She, Q.: A robust adversarial training approach
381 to machine reading comprehension. In: Proceedings of the AAAI Conference on Artificial
382 Intelligence (2020)
- 383 [36] Ma, C., Yang, Y., Ju, C., Zhang, F., Liu, J., Wang, Y., Zhang, Y., Wang, Y.: Diffusionseg:
384 Adapting diffusion towards unsupervised object discovery. arXiv preprint arXiv:2303.09813
385 (2023)
- 386 [37] Ma, C., Yang, Y., Ju, C., Zhang, F., Zhang, Y., Wang, Y.: Open-vocabulary semantic segmen-
387 tation via attribute decomposition-aggregation. Advances in Neural Information Processing
388 Systems (2024)
- 389 [38] Nag, S., Zhu, X., Song, Y.Z., Xiang, T.: Zero-shot temporal action detection via vision-language
390 prompting. In: Proceedings of the European Conference on Computer Vision. Springer (2022)
- 391 [39] Pruthi, D., Dhingra, B., Lipton, Z.C.: Combating adversarial misspellings with robust word
392 recognition. arXiv preprint arXiv:1905.11268 (2019)
- 393 [40] Qian, R., Li, Y., Xu, Z., Yang, M.H., Belongie, S., Cui, Y.: Multimodal open-vocabulary video
394 classification via pre-trained vision and language models. arXiv preprint arXiv:2207.07646
395 (2022)

- 396 [41] Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell,
397 A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language
398 supervision. In: Proceedings of the International Conference on Machine Learning. PMLR
399 (2021)
- 400 [42] Rychalska, B., Basaj, D., Gosiewska, A., Biecek, P.: Models in the wild: On corruption robust-
401 ness of neural nlp systems. In: Neural Information Processing: 26th International Conference,
402 ICONIP 2019, Sydney, NSW, Australia, December 12–15, 2019, Proceedings, Part III 26.
403 Springer (2019)
- 404 [43] Sakaguchi, K., Duh, K., Post, M., Van Durme, B.: Robust word recognition via semi-character
405 recurrent neural network. In: Proceedings of the AAAI Conference on Artificial Intelligence
406 (2017)
- 407 [44] Smaira, L., Carreira, J., Noland, E., Clancy, E., Wu, A., Zisserman, A.: A short note on the
408 kinetics-700-2020 human action dataset. arXiv preprint arXiv:2010.10864 (2020)
- 409 [45] Soomro, K., Zamir, A.R., Shah, M.: UCF101: A dataset of 101 human actions classes from
410 videos in the wild. arXiv preprint arXiv:1212.0402 (2012)
- 411 [46] Soomro, K., Zamir, A.R., Shah, M.: Ucf101: A dataset of 101 human actions classes from
412 videos in the wild. arXiv preprint arXiv:1212.0402 (2012)
- 413 [47] Sun, S., Gu, J., Gong, S.: Benchmarking robustness of text-image composed retrieval. arXiv
414 preprint arXiv:2311.14837 (2023)
- 415 [48] Wang, H., Yan, C., Wang, S., Jiang, X., Tang, X., Hu, Y., Xie, W., Gavves, E.: Towards
416 open-vocabulary video instance segmentation. In: Proceedings of the International Conference
417 on Computer Vision (2023)
- 418 [49] Wang, M., Xing, J., Liu, Y.: Actionclip: A new paradigm for video action recognition. arXiv
419 preprint arXiv:2109.08472 (2021)
- 420 [50] Wang, W., Wang, R., Wang, L., Wang, Z., Ye, A.: Towards a robust deep neural network in
421 texts: A survey. arXiv preprint arXiv:1902.07285 (2019)
- 422 [51] Wang, Z., Wang, H.: Defense of word-level adversarial attacks via random substitution encoding.
423 In: Knowledge Science, Engineering and Management: 13th International Conference, KSEM
424 2020, Hangzhou, China, August 28–30, 2020, Proceedings, Part II 13. Springer (2020)
- 425 [52] Xu, H., Ghosh, G., Huang, P.Y., Okhonko, D., Aghajanyan, A., Metze, F., Zettlemoyer, L.,
426 Feichtenhofer, C.: Videoclip: Contrastive pre-training for zero-shot video-text understanding.
427 arXiv preprint arXiv:2109.14084 (2021)
- 428 [53] Xu, J., Zhao, L., Yan, H., Zeng, Q., Liang, Y., Sun, X.: Lexicalat: Lexical-based adversarial re-
429 inforcement training for robust sentiment classification. In: Proceedings of the 2019 conference
430 on empirical methods in natural language processing and the 9th international joint conference
431 on natural language processing (EMNLP-IJCNLP). pp. 5518–5527 (2019)
- 432 [54] Yang, Y., Ma, C., Ju, C., Zhang, Y., Wang, Y.: Multi-modal prototypes for open-set semantic
433 segmentation. arXiv preprint arXiv:2307.02003 (2023)
- 434 [55] Yao, L., Huang, R., Hou, L., Lu, G., Niu, M., Xu, H., Liang, X., Li, Z., Jiang, X., Xu, C.:
435 Filip: Fine-grained interactive language-image pre-training. In: Proceedings of the International
436 Conference on Learning Representations (2022)
- 437 [56] Ye, Z., Ju, C., Ma, C., Zhang, X.: Unsupervised domain adaptation via similarity-based prototypes
438 for cross-modality segmentation. In: Domain Adaptation and Representation Transfer, and
439 Affordable Healthcare and AI for Resource Diverse Global Health: Third MICCAI Workshop,
440 DART 2021, and First MICCAI Workshop, FAIR 2021, Held in Conjunction with MICCAI
441 2021, Strasbourg, France, September 27 and October 1, 2021, Proceedings 3 (2021)
- 442 [57] Yuan, L., Chen, D., Chen, Y.L., Codella, N., Dai, X., Gao, J., Hu, H., Huang, X., Li, B., Li, C.,
443 et al.: Florence: A new foundation model for computer vision. arXiv preprint arXiv:2111.11432
444 (2021)
- 445 [58] Zhai, X., Wang, X., Mustafa, B., Steiner, A., Keysers, D., Kolesnikov, A., Beyer, L.: Lit:
446 Zero-shot transfer with locked-image text tuning. In: Proceedings of the IEEE Conference on
447 Computer Vision and Pattern Recognition (2022)

- 448 [59] Zhang, W.E., Sheng, Q.Z., Alhazmi, A., Li, C.: Adversarial attacks on deep-learning models
449 in natural language processing: A survey. *ACM Transactions on Intelligent Systems and*
450 *Technology (TIST)* (2020)
- 451 [60] Zhao, P., Xie, L., Ju, C., Zhang, Y., Wang, Y., Tian, Q.: Bottom-up temporal action localization
452 with mutual regularization. In: *Proceedings of the European Conference on Computer Vision*
453 (2020)
- 454 [61] Zheng, H., Chen, X., Yao, J., Yang, H., Li, C., Zhang, Y., Zhang, H., Tsang, I., Zhou, J., Zhou,
455 M.: Contrastive attraction and contrastive repulsion for representation learning. *arXiv preprint*
456 *arXiv:2105.03746* (2021)
- 457 [62] Zhou, J., Dong, L., Gan, Z., Wang, L., Wei, F.: Non-contrastive learning meets language-
458 image pre-training. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern*
459 *Recognition* (2023)
- 460 [63] Zhou, Y., Jiang, J.Y., Chang, K.W., Wang, W.: Learning to discriminate perturbations for
461 blocking adversarial attacks in text classification. *arXiv preprint arXiv:1909.03084* (2019)

462 A Theoretical Analysis

463 A.1 Decoding Objective

464 At each step i , the decoding objective to find $\operatorname{argmax}_{t_i} p(t_i|\mathcal{T}_{i-1}, \mathcal{V})$. Note that, $p(\mathcal{T}_{i-1}, \mathcal{V})$ is same
465 for all possible t_i . As a result, our objective is written as:

$$\operatorname{argmax}_{t_i} p(t_i|\mathcal{T}_{i-1}, \mathcal{V}) = \operatorname{argmax}_{t_i} p(t_i|\mathcal{T}_{i-1}, \mathcal{V})p(\mathcal{T}_{i-1}, \mathcal{V}) \quad (11)$$

$$= \operatorname{argmax}_{t_i} p(t_i, \mathcal{T}_{i-1}, \mathcal{V}) \quad (12)$$

$$= \operatorname{argmax}_{t_i} \log p(t_i, \mathcal{T}_{i-1}, \mathcal{V}) \quad (13)$$

466 A.2 Discriminative Step

467 At the discriminative step, we choose the best set of \mathcal{V} that helps decode $t_{c,i}$ for each semantic-class
468 c . To understand why \mathcal{V}_c , the set of visual samples v_j whose labels \mathcal{Y}_j are assigned to semantic-class
469 c are those who help decode most efficiently, we first introduce a hidden discrete random variable
470 $z_j \sim Q_j$ for each v_j , indicating the index of class assignment. $z_j = c$ means that $\operatorname{argmax} \mathcal{Y}_j = c$.

471 Knowing that all visual samples are independent and using Jensen inequality:

$$\log p(t_i, \mathcal{T}_{i-1}, \mathcal{V}) = \sum_j \log p(t_i, \mathcal{T}_{i-1}, v_j) \quad (14)$$

$$= \sum_j \log \sum_{z_j} p(t_i, \mathcal{T}_{i-1}, v_j, z_j) \quad (15)$$

$$= \sum_j \log \sum_{z_j} Q_j(z_j) \frac{p(t_i, \mathcal{T}_{i-1}, v_j, z_j)}{Q_j(z_j)} \quad (16)$$

$$\geq \sum_j \sum_{z_j} Q_j(z_j) \log \frac{p(t_i, \mathcal{T}_{i-1}, v_j, z_j)}{Q_j(z_j)} \quad (17)$$

472 Equality is attained at $Q_j(z_j) \propto p(t_i, \mathcal{T}_{i-1}, v_j, z_j)$. Since $\sum_{z_j} Q_j(z_j) = 1$, to maximize the lower
473 bound, we have:

$$Q_j(z_j) = \frac{p(t_i, \mathcal{T}_{i-1}, v_j, z_j)}{\sum_{z_j} p(t_i, \mathcal{T}_{i-1}, v_j, z_j)} \quad (18)$$

$$= \frac{p(t_i, \mathcal{T}_{i-1}, v_j, z_j)}{p(t_i, \mathcal{T}_{i-1}, v_j)} \quad (19)$$

$$= p(z_j|t_i, \mathcal{T}_{i-1}, v_j) \quad (20)$$

$$= p(z_j|\mathcal{T}_i, v_j) \quad (21)$$

474 Given class texts and visual samples, the best estimation is:

$$\mathbb{P}(z_j = c|\mathcal{T}_i, v_j) = \begin{cases} 1 & c = \operatorname{argmax}_c \max_k \frac{\exp(S(v_j, \mathcal{T}_{c,i}^k))}{\sum_{k'} \exp(S(v_j, \mathcal{T}_{c,i}^{k'}))} \\ 0 & \text{otherwise} \end{cases} \quad (22)$$

475 Note that, Q_j is well defined because:

$$\lim_{Q_j(z_j) \rightarrow 0^+} Q_j(z_j) \log \frac{p(t_i, \mathcal{T}_{i-1}, v_j, z_j)}{Q_j(z_j)} = 0 \quad (23)$$

476 With Q_j defined in this way, we find the discriminative step to be identical to how Φ_{OVAR} assigns
 477 labels. We have $Q_j(c) = 1$ only for $\{j|v_j \in \mathcal{V}_c\}$:

$$\log p(t_i, \mathcal{T}_{i-1}, \mathcal{V}) \geq \sum_j \sum_{z_j} Q_j(z_j) \log \frac{p(t_i, \mathcal{T}_{i-1}, v_j, z_j)}{Q_j(z_j)} \quad (24)$$

$$= \sum_c \sum_{j, v_j \in \mathcal{V}_c} \sum_{z_j} Q_j(z_j) \log \frac{p(t_i, \mathcal{T}_{i-1}, v_j, z_j)}{Q_j(z_j)} \quad (25)$$

$$= \sum_c \sum_{j, v_j \in \mathcal{V}_c} \log p(t_i, \mathcal{T}_{i-1}, v_j, z_j = c) \quad (26)$$

$$= \sum_c \log p(t_{c,i}, \mathcal{T}_{c,i-1}, \mathcal{V}_c) \quad (27)$$

$$(28)$$

478 A.3 Generative Step

479 We optimize $t_{c,i}$ for each semantic-class:

$$\operatorname{argmax}_{t_{c,i}} \log p(t_{c,i}, \mathcal{T}_{c,i-1}, \mathcal{V}_c) = \operatorname{argmax}_{t_{c,i}} p(t_{c,i}, \mathcal{T}_{c,i-1}, \mathcal{V}_c) \quad (29)$$

$$= \operatorname{argmax}_{t_{c,i}} \prod_{v_j \in \mathcal{V}_c} p(t_{c,i}, \mathcal{T}_{c,i-1}, v_j) \quad (30)$$

$$= \operatorname{argmax}_{t_{c,i}} \prod_{v_j \in \mathcal{V}_c} p(\mathcal{T}_{c,i-1} | t_{c,i}, v_j) p(t_{c,i} | v_j) p(v_j) \quad (31)$$

$$= \operatorname{argmax}_{t_{c,i}} \prod_{v_j \in \mathcal{V}_c} p(\mathcal{T}_{c,i-1} | t_{c,i}, v_j) p(t_{c,i} | v_j) \quad (32)$$

480 Noting that $p(\mathcal{T}_{c,i-1})$ is the same for any possible $t_{c,i}$:

$$\operatorname{argmax}_{t_{c,i}} p(\mathcal{T}_{c,i-1} | t_{c,i}, v_j) = \operatorname{argmax}_{t_{c,i}} p(\mathcal{T}_{c,i-1} | t_{c,i}) \quad (33)$$

$$= \operatorname{argmax}_{t_{c,i}} \frac{p(t_{c,i} | \mathcal{T}_{c,i-1}) p(\mathcal{T}_{c,i-1})}{p(t_{c,i})} \quad (34)$$

$$= \operatorname{argmax}_{t_{c,i}} \frac{p(t_{c,i} | \mathcal{T}_{c,i-1})}{p(t_{c,i})} \quad (35)$$

481 It is possible to optimize with prior $p(t_{c,i})$ by considering that the more a word is frequent, the less it
 482 is likely to be misspelled in real-world scenarios. In this paper, for simplicity, we assume the $t_{c,i}$ to
 483 be uniform:

$$\operatorname{argmax}_{t_{c,i}} p(\mathcal{T}_{c,i-1} | t_{c,i}, v_j) = \operatorname{argmax}_{t_{c,i}} p(t_{c,i} | \mathcal{T}_{c,i-1}) \quad (36)$$

484 B Additional Experiments

485 B.1 DENOISER vs. Adversarial Training

486 Fig. 5 studies how adversarial training might mitigate the noise in text descriptions. We first train
 487 ActionCLIP ViT-B/32-8F from scratch on K400 by randomly injecting noise in its text labels, then
 488 test the model’s zero-shot performance on UCF101 under different noise rate scenarios. We find that
 489 adversarial training, though promising under closed-set scenarios in previous studies, is relatively
 490 ineffective under open-vocabulary settings. Specifically, training with more noise lowers significantly
 491 the model’s performance under low noise rate. Additionally, its added value is limited under heavy
 492 noise rate. These phenomena are probably related to the domain gap between datasets. By training
 493 on noisy text descriptions, the model tends to overfit the noise pattern, jeopardizing its zero-shot
 494 performance. We conclude that noisy text descriptions are better solved in testing time rather than
 495 during training stage. Our DENOISER framework shows a significant advantage over the adversarial
 496 training.

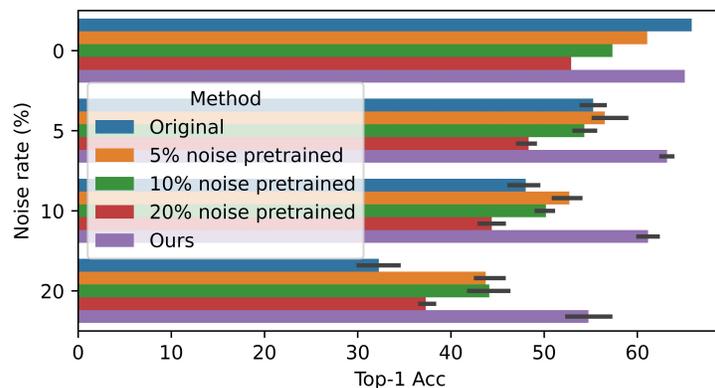


Figure 5: **Comparison to Adversarial Training.** Adversarial training is not efficient, especially in low-noise scenarios, even leading to a lower performance compared to the original model. It also falls behind our method by a significant margin.

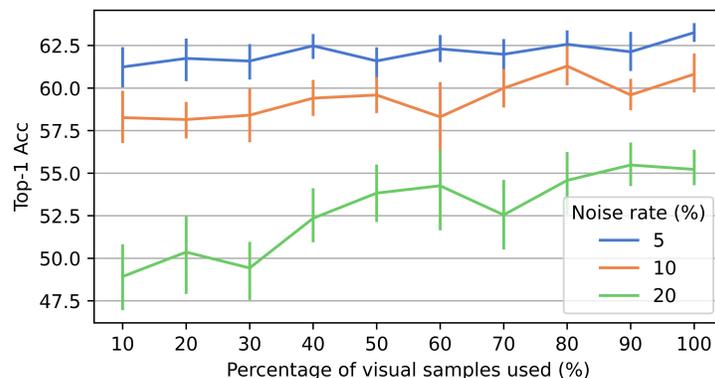


Figure 6: **Ablation Study on the Number of Visual Samples.** When fewer visual samples are used in Φ_{inter} , our method shows a drop in performance. The bigger the noise rate, the larger the drop, showing that Φ_{inter} plays a role of increasing importance when the noise is larger.

497 **B.2 Ablation Study on the Number of Visual Samples**

498 Fig. 6 ablates on the number of visual samples in Φ_{inter} . Our method shows a drop in performance
 499 when fewer visual samples are used in Φ_{inter} . The performance tends to converge towards that
 500 when solely Φ_{intra} is used. We hypothesize that fewer visual samples make Φ_{inter} harder to extract
 501 added value to Φ_{intra} . With the noise rate increasing, we find an increasingly large drop in perfor-
 502 mance, which shows conversely that Φ_{inter} is more important under large noise scenarios as textual
 503 information becomes more ambiguous and less informative.

504 **B.3 Qualitative Results**

505 Fig. 7 visualizes the embedding of (visual samples, text descriptions) from three semantic-classes:
 506 bird (green), ship (yellow), truck (blue) in CIFAR-10 using T-SNE. The first principal component of
 507 textual embedding is removed following ReCLIP[12] to prevent them from clustering at the same
 508 place. The Left shows that classification accuracy is low when text descriptions are noisy. Almost
 509 all visual samples are recognized as “bird”. The Middle shows the embeddings of proposed text
 510 candidates. Some of them remain at the same place, because they move perpendicular to this 2D space
 511 in the real semantic space. We assign the best set of visual samples for each semantic-class to help
 512 denoise, e.g., the blue dots are used to vote on the two candidates “trump” (red) and “truck” (purple)
 513 of “trunk”. The Right shows that the denoised text descriptions improve the OVAR performance.

514 Tab. 4 quantifies some good/bad cases. We find GPT 3.5 is better at understanding semantics of
 515 noisy text descriptions, e.g., “wal4ingm with a dog” \rightarrow “dogwalking”. However, its output is highly

Table 4: **Cases of Denoised Text Descriptions for GPT 3.5 and DENOISER.** The output from GPT 3.5 [1] tends to be unstable, and sometimes it’s a relatively high-level understanding of noisy text descriptions. Our *DENOISER* ensures a relatively faithful output in terms of spelling but could be slightly mistaken when two words are similar in terms of both semantics and spelling.

	Ground Truth	Noisy Text Descriptions	GPT 3.5 [1]	Ours
Good Case	walking with a dog baby crawling cutting in kitchen	wal4ingm with a dog babty crawling cutting i_aitnchen	dogwalking baby crying cutting	walking with a dog baby crawling cutting in kitchen
Bad Case	juggling balls	juggling ball_	juggling	juggling ball_

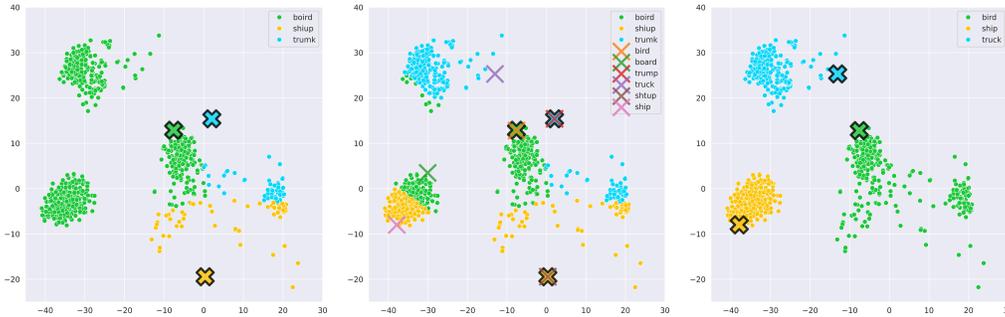


Figure 7: **Denoising Visualization.** **Left:** result with noisy text descriptions (crosses w black border). **Middle:** text candidates (crosses w/o black border), the visual samples (in dots) that are used to vote for candidates. **Right:** denoised class texts (crosses w black border) help for better classification.

516 affected by input prompts, and thus tends to be unstable: important text parts are sometimes omitted
 517 or misinterpreted, *e.g.*, “babty crawling” → “baby crying”. Such unstable outputs require manual
 518 cleaning, limiting its applications in reality. Our *DENOISER* remains faithful in terms of spelling,
 519 *e.g.*, “wal4ingm with a dog” → “walking with a dog” instead of “dogwalking”. While it may be
 520 mistaken when two words are similar in semantics and spelling (rare cases), *e.g.*, “ball” and “balls”.

521 C On the efficiency of DENOISER

522 Our model requires a trade-off between computational cost and performance. As shown in Fig. 4
 523 and Fig. 6, the performance of our *DENOISER* increases as the number of proposals K and the
 524 percentage of the visual samples used. Since the theoretical complexity of *DENOISER* increases
 525 linearly with K and the percentage of visual samples used, while the marginal contribution of a larger
 526 K or percentage is decreasing, a trade-off between computational cost and performance is necessary.

527 *DENOISER* requires only simple operations for each iteration. After having extracted the embedding
 528 of visual samples, *DENOISER* only requires recomputing the text embedding and doing a dot product
 529 with visual embeddings, which is extremely fast. Compared to other approaches that intend to align
 530 noisy text-image pairs or to train spell-checking models, *DENOISER* that denoises at evaluation time
 531 is extremely time-saving.

532 **NeurIPS Paper Checklist**

533 **1. Claims**

534 Question: Do the main claims made in the abstract and introduction accurately reflect the
535 paper's contributions and scope?

536 Answer: [\[Yes\]](#)

537 Justification: The main claims made in the abstract and introduction accurately reflect the
538 paper's contributions and scope.

539 Guidelines:

- 540 • The answer NA means that the abstract and introduction do not include the claims
541 made in the paper.
- 542 • The abstract and/or introduction should clearly state the claims made, including the
543 contributions made in the paper and important assumptions and limitations. A No or
544 NA answer to this question will not be perceived well by the reviewers.
- 545 • The claims made should match theoretical and experimental results, and reflect how
546 much the results can be expected to generalize to other settings.
- 547 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
548 are not attained by the paper.

549 **2. Limitations**

550 Question: Does the paper discuss the limitations of the work performed by the authors?

551 Answer: [\[Yes\]](#)

552 Justification: We discuss the limitation of our method at the end of the paper, and in the
553 appendix.

554 Guidelines:

- 555 • The answer NA means that the paper has no limitation while the answer No means that
556 the paper has limitations, but those are not discussed in the paper.
- 557 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 558 • The paper should point out any strong assumptions and how robust the results are to
559 violations of these assumptions (e.g., independence assumptions, noiseless settings,
560 model well-specification, asymptotic approximations only holding locally). The authors
561 should reflect on how these assumptions might be violated in practice and what the
562 implications would be.
- 563 • The authors should reflect on the scope of the claims made, e.g., if the approach was
564 only tested on a few datasets or with a few runs. In general, empirical results often
565 depend on implicit assumptions, which should be articulated.
- 566 • The authors should reflect on the factors that influence the performance of the approach.
567 For example, a facial recognition algorithm may perform poorly when image resolution
568 is low or images are taken in low lighting. Or a speech-to-text system might not be
569 used reliably to provide closed captions for online lectures because it fails to handle
570 technical jargon.
- 571 • The authors should discuss the computational efficiency of the proposed algorithms
572 and how they scale with dataset size.
- 573 • If applicable, the authors should discuss possible limitations of their approach to
574 address problems of privacy and fairness.
- 575 • While the authors might fear that complete honesty about limitations might be used by
576 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
577 limitations that aren't acknowledged in the paper. The authors should use their best
578 judgment and recognize that individual actions in favor of transparency play an impor-
579 tant role in developing norms that preserve the integrity of the community. Reviewers
580 will be specifically instructed to not penalize honesty concerning limitations.

581 **3. Theory Assumptions and Proofs**

582 Question: For each theoretical result, does the paper provide the full set of assumptions and
583 a complete (and correct) proof?

584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636

Answer: [Yes]

Justification: We provide detailed derivation in the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We detail the proposed algorithm and the setting of experiments. Additionally, we provide source code.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

637 Question: Does the paper provide open access to the data and code, with sufficient instruc-
638 tions to faithfully reproduce the main experimental results, as described in supplemental
639 material?

640 Answer: [Yes]

641 Justification: We provide source code. Datasets are publicly accessible.

642 Guidelines:

- 643 • The answer NA means that paper does not include experiments requiring code.
- 644 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
645 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 646 • While we encourage the release of code and data, we understand that this might not be
647 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
648 including code, unless this is central to the contribution (e.g., for a new open-source
649 benchmark).
- 650 • The instructions should contain the exact command and environment needed to run to
651 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
652 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 653 • The authors should provide instructions on data access and preparation, including how
654 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 655 • The authors should provide scripts to reproduce all experimental results for the new
656 proposed method and baselines. If only a subset of experiments are reproducible, they
657 should state which ones are omitted from the script and why.
- 658 • At submission time, to preserve anonymity, the authors should release anonymized
659 versions (if applicable).
- 660 • Providing as much information as possible in supplemental material (appended to the
661 paper) is recommended, but including URLs to data and code is permitted.

662 6. Experimental Setting/Details

663 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
664 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
665 results?

666 Answer: [Yes]

667 Justification: We specify all settings of experiments in the experiments section.

668 Guidelines:

- 669 • The answer NA means that the paper does not include experiments.
- 670 • The experimental setting should be presented in the core of the paper to a level of detail
671 that is necessary to appreciate the results and make sense of them.
- 672 • The full details can be provided either with the code, in appendix, or as supplemental
673 material.

674 7. Experiment Statistical Significance

675 Question: Does the paper report error bars suitably and correctly defined or other appropriate
676 information about the statistical significance of the experiments?

677 Answer: [Yes]

678 Justification: We report confidence intervals.

679 Guidelines:

- 680 • The answer NA means that the paper does not include experiments.
- 681 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
682 dence intervals, or statistical significance tests, at least for the experiments that support
683 the main claims of the paper.
- 684 • The factors of variability that the error bars are capturing should be clearly stated (for
685 example, train/test split, initialization, random drawing of some parameter, or overall
686 run with given experimental conditions).
- 687 • The method for calculating the error bars should be explained (closed form formula,
688 call to a library function, bootstrap, etc.)

- 689
- The assumptions made should be given (e.g., Normally distributed errors).
 - 690 • It should be clear whether the error bar is the standard deviation or the standard error
691 of the mean.
 - 692 • It is OK to report 1-sigma error bars, but one should state it. The authors should
693 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
694 of Normality of errors is not verified.
 - 695 • For asymmetric distributions, the authors should be careful not to show in tables or
696 figures symmetric error bars that would yield results that are out of range (e.g. negative
697 error rates).
 - 698 • If error bars are reported in tables or plots, The authors should explain in the text how
699 they were calculated and reference the corresponding figures or tables in the text.

700 8. Experiments Compute Resources

701 Question: For each experiment, does the paper provide sufficient information on the com-
702 puter resources (type of compute workers, memory, time of execution) needed to reproduce
703 the experiments?

704 Answer: [Yes]

705 Justification: We report information of computer resources.

706 Guidelines:

- 707 • The answer NA means that the paper does not include experiments.
- 708 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
709 or cloud provider, including relevant memory and storage.
- 710 • The paper should provide the amount of compute required for each of the individual
711 experimental runs as well as estimate the total compute.
- 712 • The paper should disclose whether the full research project required more compute
713 than the experiments reported in the paper (e.g., preliminary or failed experiments that
714 didn't make it into the paper).

715 9. Code Of Ethics

716 Question: Does the research conducted in the paper conform, in every respect, with the
717 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

718 Answer: [Yes]

719 Justification: We conduct in the paper conform, in every respect, with the NeurIPS Code of
720 Ethics.

721 Guidelines:

- 722 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 723 • If the authors answer No, they should explain the special circumstances that require a
724 deviation from the Code of Ethics.
- 725 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
726 eration due to laws or regulations in their jurisdiction).

727 10. Broader Impacts

728 Question: Does the paper discuss both potential positive societal impacts and negative
729 societal impacts of the work performed?

730 Answer: [Yes]

731 Justification: Our model helps users better leverage the existing Open-Vocabulary models in
732 a more robust way.

733 Guidelines:

- 734 • The answer NA means that there is no societal impact of the work performed.
- 735 • If the authors answer NA or No, they should explain why their work has no societal
736 impact or why the paper does not address societal impact.
- 737 • Examples of negative societal impacts include potential malicious or unintended uses
738 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
739 (e.g., deployment of technologies that could make decisions that unfairly impact specific
740 groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All the assets are properly cited. License and terms of use are properly respected.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

794 • If this information is not available online, the authors are encouraged to reach out to
795 the asset’s creators.

796 **13. New Assets**

797 Question: Are new assets introduced in the paper well documented and is the documentation
798 provided alongside the assets?

799 Answer: [Yes]

800 Justification: We provided well-documented source code.

801 Guidelines:

- 802 • The answer NA means that the paper does not release new assets.
- 803 • Researchers should communicate the details of the dataset/code/model as part of their
804 submissions via structured templates. This includes details about training, license,
805 limitations, etc.
- 806 • The paper should discuss whether and how consent was obtained from people whose
807 asset is used.
- 808 • At submission time, remember to anonymize your assets (if applicable). You can either
809 create an anonymized URL or include an anonymized zip file.

810 **14. Crowdsourcing and Research with Human Subjects**

811 Question: For crowdsourcing experiments and research with human subjects, does the paper
812 include the full text of instructions given to participants and screenshots, if applicable, as
813 well as details about compensation (if any)?

814 Answer: [NA]

815 Justification: The paper does not involve crowdsourcing nor research with human subjects.

816 Guidelines:

- 817 • The answer NA means that the paper does not involve crowdsourcing nor research with
818 human subjects.
- 819 • Including this information in the supplemental material is fine, but if the main contribu-
820 tion of the paper involves human subjects, then as much detail as possible should be
821 included in the main paper.
- 822 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
823 or other labor should be paid at least the minimum wage in the country of the data
824 collector.

825 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
826 Subjects**

827 Question: Does the paper describe potential risks incurred by study participants, whether
828 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
829 approvals (or an equivalent approval/review based on the requirements of your country or
830 institution) were obtained?

831 Answer: [NA]

832 Justification: The paper does not involve crowdsourcing nor research with human subjects.

833 Guidelines:

- 834 • The answer NA means that the paper does not involve crowdsourcing nor research with
835 human subjects.
- 836 • Depending on the country in which research is conducted, IRB approval (or equivalent)
837 may be required for any human subjects research. If you obtained IRB approval, you
838 should clearly state this in the paper.
- 839 • We recognize that the procedures for this may vary significantly between institutions
840 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
841 guidelines for their institution.
- 842 • For initial submissions, do not include any information that would break anonymity (if
843 applicable), such as the institution conducting the review.