## 000 'NO' MATTERS: OUT-OF-DISTRIBUTION DETECTION IN MULTIMODALITY LONG DIALOGUE

**Anonymous authors** 

Paper under double-blind review

### ABSTRACT

Out-of-distribution (OOD) detection in multimodal contexts is essential for identifying deviations in combined inputs from different modalities, particularly in applications like open-domain dialogue systems or real-life dialogue interactions. This paper aims to improve the user experience that involves multi-round long dialogues by efficiently detecting OOD dialogues and images. We introduce a novel scoring framework named Dialogue Image Aligning and Enhancing Framework (DIAEF) that integrates the visual language models with the novel proposed scores that detect OOD in two key scenarios (1) mismatches between the dialogue and image input pair and (2) input pairs with previously unseen labels. Our experimental results, derived from various benchmarks, demonstrate that integrating image and multi-round dialogue OOD detection is more effective with previously unseen labels than using either modality independently. In the presence of mismatched pairs, our proposed score effectively identifies these mismatches and demonstrates strong robustness in long dialogues. This approach enhances domain-aware, adaptive conversational agents and establishes baselines for future studies.<sup>1</sup>

024 025 026

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

#### INTRODUCTION 1

027 028 029

In the regime of multimodal learning contexts, Out-Of-Distribution (OOD) detection involves identifying whether some unknown inputs from different modalities (e.g., text and images) deviate 031 significantly from the patterns in the previously seen data. Specifically, an OOD instance under the multimodal setting is defined as one that does not conform to a certain distribution of interest, either 033 by deviating in one modality or by showing the discrepancy across different modalities (Arora et al., 034 2021; Chen et al., 2021; Feng et al., 2022; Hsu et al., 2020). This is crucial in applications such as dialogue-image systems where the synergy between spoken or written language and visual elements is expected to adhere to certain semantic and contextual norms when identifying the In-Distribution (ID) pairs where they come from some known distribution. 037

Particularly in the dialogue system with inputs from different modalities, efficiently handling OOD queries/images can significantly improve user satisfaction and trust as the response quality hinges tightly on the understanding of the semantics from different modality inputs. Recognizing and 040 managing OOD queries — those that deviate from expected dialogue or image patterns and contents 041 — is essential for maintaining these systems' reliability and user experience, especially in dynamically 042 changing dialogue systems with real-life interaction from users with much noise (Gao et al., 2024a;b). 043 Taking three motivating examples as shown in Figure 1, we are given several dialogue-image pairs 044 for OOD detection where our ID label is 'cat'. We will then consider two typical OOD cases in dialogue systems where either: 1) the dialogue label and image labels are not matched, or 2) even if 046 the dialogue and image match, their labels might not be previously seen in the given data. 047

To effectively detect OOD samples in such a novel multi-modalities multi-round long dialogue 048 scenario, we introduce Dialogue Image Aligning and Enhancing Framework (DIAEF), a framework that incorporates a novel OOD score for taking the first attempt on dialogue-image OOD detection for long dialogue systems. We propose a new score design across these two modali-051 ties, enabling more targeted controls for misalignment detection and performance enhancement. 052

<sup>&</sup>lt;sup>1</sup>Code can be found in https://anonymous.4open.science/r/multimodal\_ood-E443/ README.md.

054 Such a framework could effectively boost anomaly 055 detection and give better response strategies in long 056 dialogue systems with interactive aims. This compre-057 hensive score framework not only advances the field 058 of multimodal conversations but also sets a new standard for domain-aware, adaptive long dialogue agent building for the future. To show the effectiveness of 060 the proposed framework, we constructed a dataset 061 consisting of over 120K dialogues in the multi-round 062 application Question-answering systems and open-063 domain real interactive dialogues (Seo et al., 2017; 064 Lee et al., 2021). Leveraging these dialogue datasets, 065 we apply our proposed framework and demonstrate 066 the effectiveness of the novel score design through 067 various experiments. These experiments establish 068 fundamental benchmarks and pave the way for future 069 explorations in such a novel dialogue setting. Further-



Figure 1: Motivating examples for ID, mismatched OOD and label OOD pair where the ID label is 'cat' and OOD label is 'sport'.

more, we integrate the crucial aspect of OOD detection, emphasizing its significance for enhancing
the robustness and applicability of multimodal dialogue systems (Dai et al., 2023; Dosyn et al., 2022;
Wu et al., 2024). To summarize, our contributions are listed as follows:

• We take the first attempt for OOD detection with the dialogues and propose a novel framework that enhances the OOD detection in cross-modal contexts, particularly focusing on scenarios where dialogue-image pairs either do not match with the semantics or even match, but their semantic labels are outside the known set, which matters in long-dialogue context for users.

• Our framework incorporates a novel scoring method by combining both dialogues and images to enhance the OOD detection while recognizing the mismatch pairs with the dialogue-image similarity.

• We demonstrate the practical application of our methods with the real-world multi-round long dialogue dataset, showcasing improvements in user experience and system reliability upon response. Further, our work establishes foundational benchmarks and methodologies that can serve as baseline standards for future research in the field of cross-modal detection on interactive dialogue systems.

083 084 085

073

074

075

076 077

078

079

081

# 2 RELATED WORK

086 087

880 **OOD Detection in Dialogue Systems.** Dialogue systems have become fundamental in applications 089 ranging from virtual assistants and customer service bots to educational platforms with continuous 090 multi-rounds (Feng et al., 2022; Kottur et al., 2019; Seo et al., 2017; Yu et al., 2019; Gao & Wang, 091 2024). The evolution of dialogue systems has seen a progression from rule-based and template-based 092 approaches to statistical and machine learning methods (Zheng et al., 2020; Lang et al., 2023; Deka 093 et al., 2023; Mei et al., 2024; Arora et al., 2021). Modern systems, particularly those based on deep learning models like BERT and GPT, have set new performance benchmarks (Yuan et al., 094 2024; Hendrycks et al., 2020; Yang et al., 2022; Ye et al., 2023). However, the complexity of these systems introduces challenges in understanding context and handling ambiguous semantic 096 queries, necessitating more sophisticated approaches to maintain dialogue coherence and accuracy in interactive dialogue contexts, especially for real-life cases. OOD detection is a critical aspect of 098 dialogue systems, ensuring their robustness and reliability in generating responses to user queries (Niu & Zhang, 2021; Chen et al., 2022). When dialogue systems encounter inputs that deviate from 100 the training data distribution in long multi-round data, they risk generating incorrect or nonsensical 101 answers, leading to user frustration and decreased trust. Effective OOD detection helps identify such 102 anomalous queries, allowing the system to gracefully handle or reject them, thereby maintaining the 103 quality and consistency of responses (Li & Lin, 2021), including softmax probability thresholding 104 (Liu et al., 2023; Dhuliawala et al., 2023), auxiliary models (Wang et al., 2024; Zheng et al., 2024; 105 Ramé et al., 2023), generative models (Cai & Li, 2023; Ktena et al., 2024; Graham et al., 2023), and self-supervised learning (Azizi et al., 2023; Wallin et al., 2024; Liu et al., 2021). The integration of 106 effective OOD detection mechanisms is crucial for the continued advancement and trustworthiness of 107 QA dialogue systems (Salvador et al., 2017; Feng et al., 2022).

108 Dialogue-based Multimodality OOD Detection. Due to the complexity of dialogue in multi-turns 109 and information complexity embedded in the connection of preceding turns within the long dialogue, 110 successfully detecting whether the information from the dialogue and images are within the same 111 domain stands as a technical challenge in OOD detection (Fort et al., 2021; Basart et al., 2022). 112 Previous works attempted to evaluate the generated pseudo-OOD samples' impact on the OOD section in dialogue settings (Marciniak, 2020), which improved OOD detection performance after introducing 113 the generated dialogues when utilizing unlabeled data, making the model practical and effective for 114 real-world applications (Zheng et al., 2020). Later studies used the information bottleneck principle to 115 extract robust representations by filtering out irrelevant information for multi-turn dialogue contexts 116 (Lang et al., 2023). Furthermore, the crucial aspect of OOD detection in multimodal long dialogue is 117 still under investigation, emphasizing the significance of multimodal conversational user experience 118 in question-answering systems. 119

Multi-label OOD Detection. While numerous studies have improved approaches for multi-class 120 OOD detection tasks, investigating multi-label OOD detection tasks has been notably limited. A 121 recent advancement is the introduction of Spectral Normalized Joint Energy (SNoJoE) (Mei et al., 122 2024), a method that consolidates label-specific information across multiple labels using an energy-123 based function. Later on, the sparse label co-occurrence scoring (SLCS) leverages these properties 124 by filtering low-confidence logits to enhance label sparsity and weighting preserved logits by label 125 co-occurrence information (Wang et al., 2022). Considering the vision-language information as input 126 in models like CLIP (Radford et al., 2021), traditional vision-language prompt learning methods face 127 limitations due to ID-irrelevant information in text embeddings. To address this, the Local regularized 128 Context Optimization (LoCoOp) approach enhances OOD detection by leveraging CLIP's local 129 features in one-shot settings (Miyai et al., 2024). However, previous approaches majorly implied the limitation only in computer vision tasks without focus on dialogue or Natural Language Processing 130 tasks(Wei et al., 2015; Zhang & Taneva-Popova, 2023; Wang et al., 2021; 2022). 131

132 133

#### **3 PROBLEM FORMULATION**

134

#### 5 FROBLEM FORMULATION

135 To formally define the cross-modal OOD problem, we focus on the detection with dialogue and image 136 pairs within a multi-class classification framework. Specifically, we have a batch of N pairs of images and dialogues, along with their labels, denoted by  $\{(i_n, t_n), \mathbf{y}_n\}_{n=1}^N$  where  $i_n \in \mathcal{I}$  and  $t_n \in \mathcal{T}$  denote the input image and dialogues and  $\mathcal{I}$  and  $\mathcal{I}$  are the input image. 137 138 the input image and dialogues and  $\mathcal{I}$  and  $\mathcal{T}$  are the image and dialogue spaces, respectively. Here, the 139 instance pair may be associated with multiple labels  $\mathbf{y}_n$  with  $\mathbf{y}_n = \{y_{n,1}, y_{n,2}, \cdots, y_{n,K}\} \in [0,1]^K$ where  $y_{n,k} = 1$  if the dialogue-image pair is associated with k-th label and K denotes the total 140 number of in-domain categories. Our proposed score function enhances the ability to distinguish 141 between ID and OOD data cross-joint detection for image and dialogue, making it applicable in 142 multimodality scenarios. Based on this setup, the goal of the OOD detection is to define a decision 143 function G such that: 144

145

146 147

148

149

150

151

152

153

154

 $G(i, t, \mathbf{y}) = \begin{cases} 0 & \text{if } (i, t, \mathbf{y}) \sim \mathcal{D}_{\text{out}}, \\ 1 & \text{if } (i, t, \mathbf{y}) \sim \mathcal{D}_{\text{in}}. \end{cases}$ (1)

**Remark 1** Different from unimodal OOD detection (Lee et al., 2018; Basart et al., 2022; Hendrycks & Gimpel, 2016; Du et al., 2022; Wu et al., 2023), in the cross-modal detection scenarios, we need to additionally consider whether the image and dialogue come from the same distribution, i.e., whether the image and dialogue are semantically matched in the interaction context. In particular, we will consider several scenarios for detecting OOD samples: 1). the image and dialogue do not match (e.g., in terms of content or description), and 2). the in-domain sample does not contain any out-of-domain labels, meaning previously unseen labels appear, or 3). both cases occur simultaneously.

To determine G in practice, we may need to consider the relationship between dialogue and images additionally. To this end, let  $M : \mathcal{I} \cup \mathcal{T} \to \mathbb{R}^d$  be a vision-language model that could encode the image  $i_n$  with the image embedding  $x_{i,n} \in \mathbb{R}^d$ , and the dialogues with the text embedding  $x_{t,n} \in \mathbb{R}^d$ in the same latent space as in the image. To classify the relevance of an image to a dialogue according to the label  $\mathbf{y}_n$ , we first use a scoring function  $s : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ , which evaluates the similarity or relevance between the image and text embeddings from M. We then further compare these two embeddings with the label  $\mathbf{y}_n$  with the dialogue score function  $s_T : \mathbb{R}^d \times [0, 1]^K \to \mathbb{R}$  and image score function  $s_I : \mathbb{R}^d \times [0, 1]^K \to \mathbb{R}$ . For simplicity, we use  $s_I$  (or  $s_T$ ) interchangeably with 162  $s_I(x, y)$  throughout the paper. Finally, we could conduct a fusion on the three scores  $g(s, s_T, s_I)$  for 163 some fusion function g and check if the numeric value exceeds  $\lambda$  to determine whether it is in-domain 164 or out-of-domain. Given the above definition, given a dialogue-image data pair (i, t), we will examine 165 whether it is ID or OOD per dialogue-image pair in the given label set  $\mathcal{Y}$  with the following criterion.

**Definition 1 (Cross-Modal OOD Detection)** We use the following detection criterion for out-ofdomain samples.

- In-domain: given both embeddings  $x_i$  from the images and  $x_t$  from the dialogue, and a certain label y. We say the image is in-domain with the dialogue if  $g(s(x_i, x_t), s_t(x_t, \mathbf{y}), s_i(x_i, \mathbf{y})) \geq \lambda$ .
- *Out-of-domain:* given both embeddings  $x_i$  from the images and  $x_t$  from the dialogue, we say the image is out-of-domain with the dialogue if  $g(s(x_i, x_t), s_t(x_t, \mathbf{y}), s_i(x_i, \mathbf{y})) < \lambda$ .

for some fusion function g and some threshold  $\lambda$ .

176 177

166

167

168 169

170

171

172 173

174

175

178 179

199

## 4 DIALOGUE IMAGE ALIGNING AND ENHANCING FRAMEWORK

To intuitively demonstrate our framework, we draw the overall workflow in Figure 2. The workflow 181 consists of three parts: in the first stage, we will employ a vision language model, such as CLIP 182 (Radford et al., 2021) and BLIP (Li et al., 2022), to derive meaningful descriptors or feature embed-183 dings from images and dialogues, respectively. Note that the model we used here would map the 184 image and dialogue into the same latent space so that the similarity between the two can be easily 185 calculated. These processes yield embeddings  $x_I$  for images and  $x_T$  for dialogues. Then, utilizing these embeddings, we apply a scoring function  $s(x_I, x_T)$  to assess the relevance between an image and a dialogue. The outcome of this function helps us determine whether the dialogue-image pair falls 187 within the categories, indicating a high relevance in semantics, or the out-of-distribution categories 188 with mismatches, suggesting low or no relevance. 189

In addition to this score, we will further train two label extractors to compare the whole pair with the label set to determine if the pair is in-distribution or out-of-distribution using  $s_I(x_I, \mathbf{y})$  that evaluates the similarity between the image and the label and  $s_T(x_T, \mathbf{y})$  that evaluates the similarity between the text and the label. We will use conventional methods to combine these two scores and determine whether the pair is ID or OOD based on the threshold  $\lambda$ .

This paper aims to enhance the detection of OOD samples by combining dialogues and images and identifying the misalignments between them. To this end, we naturally propose the DIAEF score function in general:

$$g(x_T, x_I, \mathbf{y}; s, s_T, s_I) = s(x_T, x_I)^{\gamma} (\alpha s_I(x_I, \mathbf{y}) + (1 - \alpha)s_T(x_T, \mathbf{y})),$$
(2)

200 where the first part  $s(x_T, x_I)^{\gamma}$ , which we call 201 the alignment term, controls the similarity between the image and the dialogue. If the image 202 and dialogue are highly similar, this term will 203 be large and vice versa. This allows us to iden-204 tify the misalignment between images and dia-205 logues in a long dialogue system. The second 206 part  $(\alpha s_I + (1 - \alpha) s_T)$ , namely the enhancing 207 term, enhances the detection of OOD samples 208 by linearly combining the dialogue and image 209 scores, where the weighting hyperparameter  $\alpha$ 

Table 1: OOD Scores for  $s_I/s_T$ .

Method	Score
Probability	$P_y(x)$
MSP (Hendrycks & Gimpel, 2016)	$\max_{y \in \mathcal{Y}} \frac{f_y(x)}{\sum_y f_y(x)}$
Logits (Hendrycks et al., 2019)	$f_y(x)$
Energy (Wang et al., 2021)	$\log(1 + e^{fy(x)})$
ODIN (Liang et al., 2017)	$f_y(x+\epsilon\Delta)/T$
Mahalanobis (Lee et al., 2018)	$(x-\mu_y)^T \Sigma_y^{-1} (x-\mu_y)$

controls the relative importance of the image: if  $\alpha$  is selected to be large, we rely more on images for OOD detection; conversely for a small  $\alpha$ , we rely more on the dialogue. The purpose of using a multiplicative combination of the alignment and enhancing terms is: (1) identifying mismatched OOD pairs where either the image or dialogue might have high relevance to the label, making the enhancing term potentially large (depending on  $s_I$  or  $s_T$ ). To identify these pairs as OOD samples, we naturally multiply the enhancing term by  $s(x_T, x_I)$ ; (2) identifying matched pairs with OOD labels where  $s(x_T, x_I)$  may be large, but the enhancement term is likely to be small since the image



Figure 2: The workflow for three motivating examples for cross-modal OOD detection, including ID pair, mismatched OOD pair, and label OOD pair. The workflow consists of three main parts: the dialogue and image will be firstly processed and passed into a visual language model to get the image and dialogue embeddings; then two label extractors will be trained on both the image and dialogue embeddings for predictions and score calculations; finally the score function s,  $s_T$  and  $s_I$  are aggregated to determine the threshold  $\lambda$  at recall rate of 95%. The FPR95% is reported to demonstrate that combining images and dialogue outperforms using images or dialogue alone.

237 238 239

240

241

232

233

234

235

236

and dialogue have low relevance to the label. To show this mathematically, we give a theoretical justification for the proposed score in Appendix **B**.

The choice of the functions  $s(x_I, x_T)$  depends on the selection of the trained visual language model. 242 For example, in CLIP, contrastive loss is used to measure the similarity between images and text 243 (dialogue) based on cosine similarity (Radford et al., 2021). Similarly, BLIP employs image-text 244 matching loss and leverages cosine similarity to align the representations of images and text (Li et al., 245 2022). With those two models, selecting cosine similarity as an appropriate score for  $s(x_I, x_T)$  is 246 natural. Regarding  $s_I$  and  $s_T$ , which measure the scores between embeddings and labels, various 247 potential choices and aggregation methods are available. For example, one direct way is to use 248 the probability of the model output  $P_y(x)$  as the score for the category y with the input x, and 249 we could further aggregate the probability over all categories using sum or max methods to derive 250 our final DIEAF score. More complicated scores would involve some probability transformation, 251 such as the logits  $f_y(x)$  used in (Hendrycks et al., 2019) or the normalized version called MSP as used in (Hendrycks & Gimpel, 2016). Some other effective scores would involve the pre-trained 253 models, such as the ODIN method proposed in (Liang et al., 2017) modifies the input by adding a gradient-based perturbation, or the method proposed in (Lee et al., 2018) computes the Mahalanobis 254 distance between the embeddings from the pre-trained model and the class conditional distributions 255 in the feature space. Table 1 shows the list of possible scores that could fit in our framework. 256

257 258

259 260

261 262

263

#### 5 EXPERIMENTS

In this section, we evaluate DIAEF and other baselines using several datasets.

#### 5.1 EXPERIMENTAL SETUP

Datasets. In this section, we utilize the Visdial dataset (Das et al., 2017) and Real MMD dataset
 (Lee et al., 2021) for OOD detection in long dialogue systems. The Visdial dataset comprises over
 120K images sourced from the COCO image dataset (Lin et al., 2014), coupled with collected
 multi-round dialogues in a one-to-one mapping format between modalities. We constructed a testing
 multi-round question-answering dataset with full semantic context to evaluate our OOD detection
 methods, including all dialogue-image pairs and an additional 10K mismatched pairs. Each entry
 in this dataset contains an image, a full conversation, and a set of labels with 80 specific categories.

270 The dataset is further organized into 12 higher-level supercategories: animal, person, kitchen, food, 271 sports, electronic, accessory, furniture, indoor, appliance, vehicle, and outdoor. Another related 272 dataset, called the Real MMD dataset, contains images sourced from COCO (Lin et al., 2014) and 273 texts from different sources such as DailyDialog (Li et al., 2017) and Persona-Chat (Miller et al., 274 2017), meaning they may not be perfectly matched but instead have a certain degree of similarity. The dataset statistics are presented in Table 7 and Table 8 in Appendix A. 275

276 **OOD Label Selection.** In our study, we propose a label selection score 277 function for selecting OOD labels that effectively combines semantic dis- Table 2: Top 5 Labels 278 tance (Huang et al., 2008; Kadhim et al., 2014; Li & Han, 2013; Rahutomo 279 et al., 2012; Lahitani et al., 2016) and ontological hierarchy via the WordNet 280 path calculation (Aminu et al., 2021; Dosyn et al., 2022; Fellbaum, 2010; Marciniak, 2020; Martin, 1995). The score function integrates multiple cri-281 teria to enhance the robustness and accuracy of OOD detection. Semantic 282 distance is quantified using cosine similarity between vector representations 283 of candidate labels and the remaining labels in the label set. We compute 284 the maximum cosine similarity to any ID label for each candidate OOD label 285 and select those with values below a predefined threshold, ensuring semantic

Label	S(c)
Animal	5.12
Person	5.01
Sports	4.89
Vehicle	4.80
Outdoor	4.79

286 distinctiveness. Additionally, we leverage ontological hierarchies, such as WordNet, to measure 287 the path length between candidate labels and ID labels. Candidates with a minimum path length 288 exceeding a specified threshold are selected, ensuring they are not closely related in the hierarchy. 289 This dual-criteria approach ensures that selected OOD labels are both semantically distant and onto-290 logically distinct from ID labels, enhancing the efficacy of the OOD detection system. By integrating 291 these methods, our score function effectively mimics real scenarios where the OOD labels generally differ from the ID labels<sup>2</sup>. Therefore, we define the selection score as: 292

$$S(c) = w_1 \sum_{y \in \mathcal{Y} \setminus c} \left(1 - S_{\text{COS}}(M(c), M(y))\right) + w_2 \sum_{y \in \mathcal{Y} \setminus c} \left(1 - S_{\text{PATH}}(c, y)\right),\tag{3}$$

296 where

293

295

297 298 299

$$S_{\rm cos}(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|}, \quad S_{\rm PATH}(y_1, y_2) = \frac{1}{1 + \ell_d(y_1, y_2)}.$$
 (4)

300 Here, M(c) and M(y) are the vector representations of the candidate OOD labels c and the ID label 301 y with the encoder M, respectively,  $w_1$  and  $w_2$  are the weights assigned to each criterion, and  $\mathcal{Y}$ 302 represents the total valid label set. The term  $S_{cos}$  measures the semantic distance, and  $S_{PATH}(y_1, y_2)$ measures the ontological distance between labels with the path distance  $\ell_d(y_1, y_2)$  between the words 303  $y_1$  and  $y_2$  in the WordNet. We conduct the score selection on the Visdial dataset with  $w_1 = w_2 = 0.5$ , 304 and the top five scores with the most distinguished labels are shown in Table 2. To ensure that the 305 selected OOD labels are both semantically distant and ontologically distinct from ID labels, we select 306 candidates c where the score S(c) is the highest. 307

308 Experiment Details. Based on Table 2, we select the label 'animal' as the OOD label to show the framework's effectiveness. We will have 95K ID pairs and 37K OOD pairs for QA dataset and 12.7K 309 ID pairs and 12.2K OOD pairs for the Real MMD dataset. We will use the 8:2 train-test split, yielding 310 77K/54K and 10.2K/14.7K train/test pairs, respectively. For encoders for images and dialogues, 311 we use CLIP ViT-B/32 (Radford et al., 2021) throughout the experiments, and we trained the label 312 extractors with the ID training sample with a 5-layer fully connected network. More details are given 313 in Appendix A. Additionally, we use sum and max aggregation methods for the above methods. The 314 sum aggregation method combines the scores across all considered classes or components, providing 315 an overall score that reflects the cumulative effect. The max aggregation method selects the maximum

<sup>316</sup> 

<sup>317</sup> <sup>2</sup>For tuning label selections, we list the table below using the cosine similarity (Descending order): [sports, outdoor, animal, fashion, electronics, person, bedroom, vehicle, appliance, kitchen, food, furniture]. With 318 wordnet only: [person, animal, vehicle, furniture, appliance, kitchen, food, bedroom, fashion, electricity, outdoor, 319 sports]. Using only cosine similarity, labels skewed towards broad, abstract categories like "sports" and "outdoor", 320 reflecting a focus on general semantic similarities (complex context where more background information is 321 needed). Comparably, using only WordNet similarity emphasized specific, taxonomically grounded categories 322 like "person" and "animal", highlighting hierarchical relationships (suitable when labels are short descriptors). 323 Adaptive weighting or context-specific tuning could be explored for future refinements where weights are dynamically adjusted regarding dataset characteristics or task requirements.

Table 3: The comparison of OOD detection performance with QA dataset under CLIP extraction and different scores. Bold numbers are superior results for each DIAEF score and aggregation method.
Metrics reported in % include FPR95 (↓ indicates the lower the better), AUROC, and AUPR (↑ indicates the higher the better).

	FPR95↓ / AUROC↑ / AUPR↑				
OOD Scores	Aggregation	Baseline w/	Baseline w/ OOD Scores		
OOD Scoles	Aggregation	Image	Dialogue	w/ OOD Score	
MSP	Max	84.4/ 64.8/ 49.0	76.9/ 66.5/ 48.8	73.4/ 73.2/ 53.	
Droh	Max	60.0 / 75.6 / <b>57.9</b>	67.9 / 73.5 / 56.1	55.3 / 78.8 / 57	
F100	Sum	<b>70.7</b> / 68.3 / 49.0	91.9 / 62.3 / 45.7	72.8 / <b>73.6</b> / <b>56</b>	
Lacita	Max	60.0 / 75.6 / 57.9	67.9 / 73.5 / 56.1	57.2 / 82.6 / 72	
Logits	Sum	91.2 / 59.2 / 43.6	98.6 / 44.1 / 36.0	97.2 / 49.9 / 37	
ODIN	Max	<b>59.1</b> / 75.4 / 57.6	72.1 / 73.2 / 55.5	59.6 / <b>78.9</b> / <b>58</b>	
ODIN	Sum	<b>71.2</b> / 68.0 / 48.8	91.9/61.6/45.2	73.0 / <b>73.2</b> / <b>56</b>	
Mahalanahia	Max	<b>49.2</b> / 81.3 / 62.9	66.0 / 75.8 / 56.8	49.7 / <b>83.2</b> / <b>67</b>	
Manalanoois	Sum	88.5 / 75.5 / 57.5	78.6 / 68.6 / 50.0	75.0 / 76.2 / 60	
LointEnergy	Max	60.0 / 75.6 / 57.9	67.9 / 73.5 / 56.1	57.6 / 82.5 / 72	
JohntEnergy	Sum	58.3 / 75.8 / 58.0	67.0 / 74.1 / 57.1	55.9 / 82.3 / 72	
Average	Max	62.1 / 74.7 / 57.2	69.8 / 72.7 / 54.9	58.8 / 79.9 / 63	
	Sum	76.0 / 69.4 / 51.4	85.6 / 62.1 / 46.8	74.8 / 71.0 / 56	

score among all classes or components, highlighting the strongest single match. These aggregation methods allow us to assess the performance and robustness of each scoring function comprehensively. We use the cosine similarity for  $s(x_I, x_T)$  for CLIP embeddings and set  $\gamma = 1$  and  $\alpha = 0.5$  as the hyperparameter default values. To ensure the consistency and reliability of our results, all experiments were executed on a system featuring a single NVIDIA RTX 2080 Super GPU.

351 Adopted OOD Scores. Throughout the experiments, we used several general OOD scores to evaluate 352 the effectiveness of the framework, which includes Probability (Prob), Maximum Softmax Probability 353 (MSP) (Hendrycks & Gimpel, 2016), Logits (Hendrycks et al., 2019), Joint Energy (Wang et al., 2021), ODIN (Liang et al., 2017) and Mahalanobis distance (Lee et al., 2018). These baseline 354 methods provide a diverse set of techniques for OOD detection, each leveraging different aspects 355 of the model's output and feature embeddings. Then, we included two baselines with the DIEAF 356 scores in our evaluation. The first baseline, with image only, utilizes the score function  $s_I(x_I, \mathbf{y})$  to 357 determine the score threshold for OOD. The second baseline, with dialogue only, employs a similar 358 approach, using the score function  $s_T(x_T, \mathbf{y})$  instead. All methods are evaluated with the metrics 359 FPR95, AUROC and AUPR as previously mentioned in Section 4. 360

Evaluation. We include the following metrics in our evaluation for OOD detection: FPR95, AUROC 361 and AUPR. FPR95 measures the rate at which false positives occur when the true positive rate is 362 fixed at 95%. This metric indicates how often the model incorrectly classifies a negative instance as 363 positive when it correctly identifies 95% of all positive instances; a lower FPR95 value signifies a 364 better-performing model. AUROC evaluates the overall ability of a model to discriminate between positive and negative classes across all possible classification thresholds. It involves plotting the 366 ROC curve with the true positive rate against the false positive rate at various threshold settings. A 367 higher AUROC value denotes a better-performing model. AUPR, similar to AUROC, focuses on the 368 precision-recall curve, which plots precision against recall. This metric is particularly useful in class 369 imbalance scenarios. A better AUPR indicates a better model's performance.

370 371

372

328

5.2 MAIN RESULTS

With the aforementioned experimental settings, we evaluate various DIAEF scores on the given QA
 and Real MMD datasets and report the performance results in Table 3 and 4. The tables show that our
 proposed framework generally outperforms the results obtained using only images or dialogue across
 most metrics. In particular, the joint energy and Mahalanobis scores with the sum or max aggregation
 consistently perform well across most metrics. In addition, the naive probability and ODIN scores
 also show competitive performance. Interestingly, the max aggregation method tends to be more

FPR95↓ / AUROC↑ / AUPR↑					
OOD Scores	Aggregation	Baseline w/	Baseline w/ OOD Scores		
OOD Scoles	Aggregation	Image	Dialogue	w/ OOD Scores	
MSP	Max	91.1/56.2/19.4	91.1/56.2/19.4 94.5/52.5/18.9		
Droh	Max	79.2/64.1/22.9	93.4/53.7/19.3	75.8/74.0/36.0	
FIOD	Sum	90.6/58.2/21.1	94.4/51.9/18.7	83.0/69.7/33.1	
Logits	Max	<b>79.2</b> /64.1/22.9	93.4/53.7/19.3	84.8/ <b>70.9/34.1</b>	
Logits	Sum	94.5/49.0/17.8	97.3/47.6/17.2	98.8/38.6/14.3	
ODIN	Max	79.6/64.3/23.4	94.0/53.4/19.3	75.3/74.4/36.9	
	Sum	91.1/57.1/20.8	94.9/51.3/18.5	82.2/69.2/32.0	
Mahalanobis	Max	<b>54.9</b> /69.9/26.1	93.5/51.2/17.6	63.5/ <b>76.8/36.2</b>	
	Sum	93.3/66.0/25.2	94.2/49.1/16.9	86.6/73.3/36.5	
JointEnergy	Max	<b>79.2</b> /64.1/22.9	93.4/53.7/19.3	83.5/ <b>71.6/34.3</b>	
	Sum	<b>79.5</b> /64.9/24.4	93.6/54.2/19.5	80.5/ <b>72.8/37.4</b>	
Augrago	Max	<b>77.2</b> /63.8/22.9	93.7/53.0/19.0	78.1/ <b>72.8/35.1</b>	
Average	Sum	89.8/59.0/21.9	94.9/50.8/18.2	86.2/64.7/30.7	

Table 4: The comparison of OOD detection performance with Real MMD dataset under CLIP extraction and different scores.

380 381 382

effective than the sum method. This could be because we are dealing with a multi-label problem.
Adding up scores for all labels might introduce more noise, which confuses the OOD and ID scores and thus reduces detection performance. For dialogues, the performance is not as good as for images.
This is because dialogues contain some noises, such as stopwords, that are unrelated to the labels, whereas images with segmentation are more directly related to the labels. However, even though the dialogue alone may not perform well, combining it with images could significantly enhance the OOD detection performance. The results demonstrate that DIEAF performs effectively when combining dialogue and image scores, especially when introducing mismatched pairs.

406 407

5.3 ANALYSIS OF EXPERIMENTAL RESULTS

To gain deeper insights into the proposed framework, we conduct several ablation studies to examine the impact of mismatched pairs, the effectiveness of  $s(x_I, x_T)$ , and the choices of  $\alpha$  and  $\gamma$ .

411 **Effect of Mismatched Pairs.** To investigate the effect of the mismatched pairs, we conduct the 412 experiments with the same setting by excluding the mismatched pair in the testing set and report the 413 results in Table 5. Here, we only report FPR95 for simplicity and also compare the results by setting 414  $\gamma = 0$  without introducing the dialogue-image similarity.

415 From the table, it can be observed that when there are no mismatched pairs, setting  $\gamma$  to 1 can 416 actually harm our results to some extent. This is because, for OOD pairs without mismatched pairs, their similarity score  $s(x_I, x_T)$  can still be high. In such cases, multiplying by the similarity can 417 adversely affect OOD results. Setting  $\gamma$  to 0 in these situations improves FPR95 results for most 418 cases, indicating that simply combining image and dialogue modalities, even without mismatched 419 pairs, performs better than the unimodality. Additionally, comparing Table 3 and 5, we see that 420 introducing mismatched pairs generally leads to worse performance than having no mismatched pairs. 421 This demonstrates that mismatched pairs indeed pose a challenge for OOD detection. To achieve 422 better results, we will further study the impact of  $\gamma$  and  $\alpha$  to optimize OOD detection performance. 423

Effect of VLM models. We further tested the performance of the DIAEF score function with the
 BLIP model (Li et al., 2022) under the same setting as CLIP (also see details in Appendix A), and we
 report the results in Table 6. Even with BLIP, the pattern is still maintained as the proposed score
 achieves better performance compared to the single modality, and the framework handles mismatched
 and previously unseen OOD scenarios.

**Effect of**  $s(x_I, x_T)$ . We draw Figure 5 for image scores as an illustration that consists of three subplots showing the change of score distribution with  $s(x_I, x_T)$  introduced. Here Figures 5a and 5c present the distribution of  $s_I(x_T, x)$  and  $s_I(x_T, x)s(x_I, X_T)$  for both ID and OOD data with FPR95 highlighted, respectively. Figure 5b displays the joint distribution of  $P(s, s_I)$  for both ID and OOD

<sup>396</sup> 397

OOD Scores	Aggregation	Baseline Image Dialogue		DIAEF ( $\gamma = 0$ )	DIAEF ( $\gamma = 1$ )
MSP	Max	81.2	71.4	69.4	81.4
Prob	Max	49.7	59.9	<b>46.6</b>	64.4
	Sum	<b>63.6</b>	91.2	72.7	77.1
Logits	Max	49.7	59.9	<b>45.7</b>	47.6
	Sum	<b>90.1</b>	99.7	98.1	96.2
ODIN	Max	48.5	65.4	<b>48.5</b>	69.2
	Sum	64.2	91.2	72.4	79.3
Mahalanobis	Max	35.5	57.5	<b>34.3</b>	37.8
	Sum	86.4	73.7	68.1	<b>65.0</b>
JointEnergy	Max	49.7	59.9	46.7	48.7
	Sum	47.4	58.6	45.7	47.6
Average	Max	52.4	62.3	<b>48.5</b>	58.1
	Sum	<b>70.3</b>	82.9	71.4	73.0

Table 5: The comparison of FPR95 performance (the lower the better) in % with DIEAF framework
under different scores without any mismatched pairs. Bold numbers are superior results for each
DIAEF score and aggregation method.

Table 6: The comparison of OOD detection performance with QA dataset under BLIP extraction and different scores.

FPR95↓ / AUROC↑ / AUPR↑				
OOD Scores	Aggregation	Base Image	Baseline Image Dialogue	
MSP	Max	85.8/58.7/37.4	83.5/64.8/39.8	75.9/75.1/52.7
Proh	Max	<b>64.3</b> /71.2/45.1	80.5/67.1/42.2	67.0/ <b>78.7/56.5</b>
1100	Sum	78.8/64.4/39.3	96.8/55.9/35.9	74.2/72.7/51.2
Logita	Max	64.3/71.2/45.1	80.5/67.1/42.2	62.9/80.9/63.8
Logits	Sum	95.8/52.9/33.8	98.1/41.9/29.3	99.1/40.1/26.5
ODIN	Max	<b>63.9</b> /71.1/44.9	81.4/67.2/42.1	67.7/ <b>79.3/57.7</b>
	Sum	79.1/64.2/39.2	97.0/56.1/36.0	74.5/72.5/50.9
Mahalanahia	Max	<b>46.9</b> /77.7/50.6	81.0/66.9/40.5	52.6/ <b>87.7/75.4</b>
wianaianobis	Sum	79.7/71.5/46.2	92.5/59.0/35.9	67.6/78.7/61.0
IointEnergy	Max	64.3/71.2/45.1	80.5/67.1/42.2	62.8/81.0/63.8
JointEnergy	Sum	63.0/71.8/45.8	80.4/67.3/43.2	61.7/80.7/64.5
A	Max	64.9/70.2/44.7	81.2/66.7/41.5	64.8/80.5/61.2
Average	Sum	79.3/65.0/40.9	93.0/56.0/36.1	75.4/68.9/50.8

data, with the x-axis representing the similarity score  $s(x_I, x_T)$  and the y-axis representing the image score  $s_I(x_I, \mathbf{y})$ , with density indicated by colour intensity and marginal distributions shown as histograms. The figures show that without multiplying by  $s(x_I, x_T)$ , the distributions of ID and OOD are not well-separated, and the FPR95 is around 0.58. However, after applying the similarity score, the distributions of ID and OOD become more apart, and the FPR95 decreases to approximately 0.54. This occurs because, when examining the joint distribution, we find that for the ID data, most similarity values are around 0.25. In contrast, there are two peaks for the OOD data: one around 0.25 (for matched pairs) and another around 0.15 (for mismatched pairs). This indicates that if we multiply by this similarity, the mismatched OOD pairs would have lower scores, making distinguishing between ID and OOD easier. 

**Effect of**  $\gamma$ **.** Intuitively, when  $\gamma$  is smaller, similar and dissimilar dialogue-image pairs will have approximately the same alignment score. Conversely, when  $\gamma$  is larger, the score differences between similar and dissimilar pairs become more pronounced, emphasizing the role of dialogue-image similarity in OOD detection. Therefore, we selected several values of  $\gamma$  ranging from 0 (i.e., not using dialogue-image similarity) to 3 and plotted the curves under different score aggregation methods. Figure 3 shows that the optimal value of  $\gamma$  varies significantly depending on the choice of score and



Figure 5: An illustration of the effectiveness of  $s(x_I, x_T)$ 

aggregation method. For instance, with max aggregation, most methods show a trend where the FPR95 initially decreases with increasing  $\gamma$  and then rises again, with the optimal value around 1. However, the Energy and Logits methods show a trend of decreasing FPR95 as  $\gamma$  increases, indicating these methods are more sensitive to misalignment. On the other hand, for the sum aggregation method, changing the  $\gamma$  value has a limited effect on OOD detection. This could be because the sum method combines too much redundant label information, and the enhancement term plays a major role. If the enhancement term is not particularly effective, the impact of misalignment is minimal.

**Effect of**  $\alpha$ . When  $\alpha$  is small, we place more emphasis on the image score along with the alignment term for OOD detection; conversely, when  $\alpha$  is large, we emphasize more on the dialogue score. We plotted the results for different score aggregations in Figure 4. From the max aggregation results, we observe that using only the image or dialogue scores is not the most effective approach. Instead, combining both and selecting a value around 0.5 yields the best results, demonstrating the effectiveness of our framework. In the sum aggregation plot, we see that for most methods (except for Mahalanobis), the performance in terms of FPR95 improves as  $\alpha$  increases. This indicates that images do not significantly contribute to recognition for the sum aggregation compared to dialogue.

526 527 528

510 511

# 6 CONCLUSION AND LIMITATION

529 This paper introduces a cross-modal OOD score framework, DIAEF, designed to expand OOD 530 detection in cross-modal QA systems by integrating images and dialogues. DIAEF combines 531 alignment scores between dialogue-image pairs with an enhancing term that leverages both the image 532 and dialogue. Experimental results demonstrate DIAEF's superiority over baseline methods with 533 general metrics such as FPR95 and show the framework's effectiveness. However, there are some 534 spaces for future work. First, due to the scarcity of datasets, we initially validated our framework on 535 VisDial and demonstrated its effectiveness. More dialogue-image datasets are worth exploring for 536 validation. Second, the existing scores have proven the effectiveness of this framework, but further 537 improvements could be achieved by applying some transformations or smoothing techniques to make the distributions of ID and OOD more distinct. Finally, this framework is applicable to more visual 538 language models, such as multimodal models like BLIP, and can further enhance OOD performance using various image-text matching criteria.

#### 540 REFERENCES 541

561

565

566 567

568

569

576

580

581

582

- Enesi Femi Aminu, Ishaq Oyebisi Oyefolahan, Muhammad Bashir Abdullahi, and Muham-542 madu Tajudeen Salaudeen. An enhanced wordnet query expansion approach for ontology based 543 information retrieval system. In Information and Communication Technology and Applications: 544 Third International Conference, ICTA 2020, Minna, Nigeria, November 24–27, 2020, Revised Selected Papers 3, pp. 675–688. Springer, 2021. 546
- 547 Udit Arora, William Huang, and He He. Types of out-of-distribution texts and how to detect them. 548 arXiv preprint arXiv:2109.06827, 2021.
- 549 Shekoofeh Azizi, Laura Culp, Jan Freyberg, Basil Mustafa, Sebastien Baur, Simon Kornblith, Ting 550 Chen, Nenad Tomasev, Jovana Mitrović, Patricia Strachan, et al. Robust and data-efficient 551 generalization of self-supervised machine learning for diagnostic imaging. Nature Biomedical 552 Engineering, 7(6):756–779, 2023. 553
- Steven Basart, Mazeika Mantas, Mostajabi Mohammadreza, Steinhardt Jacob, and Song Dawn. 554 Scaling out-of-distribution detection for real-world settings. In International Conference on 555 Machine Learning, 2022. 556
- Mu Cai and Yixuan Li. Out-of-distribution detection via frequency-regularized generative models. 558 In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 559 5521-5530, 2023.
- Jiefeng Chen, Yixuan Li, Xi Wu, Yingyu Liang, and Somesh Jha. Atom: Robustifying out-ofdistribution detection using outlier mining. In Machine Learning and Knowledge Discovery in 562 Databases. Research Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 563 13-17, 2021, Proceedings, Part III 21, pp. 430-445. Springer, 2021.
  - Long Chen, Yuhang Zheng, and Jun Xiao. Rethinking data augmentation for robust visual question answering. In European conference on computer vision, pp. 95–112. Springer, 2022.
  - Yi Dai, Hao Lang, Kaisheng Zeng, Fei Huang, and Yongbin Li. Exploring large language models for multi-modal out-of-distribution detection. arXiv preprint arXiv:2310.08027, 2023.
- Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José MF Moura, Devi Parikh, 570 and Dhruv Batra. Visual dialog. In Proceedings of the IEEE conference on computer vision and 571 pattern recognition, pp. 326–335, 2017. 572
- 573 Rimpa Deka, Palash Pratim Dutta, and Aparajita Dutta. Distributed feature representations for 574 out-of-domain detection in dialogue systems. In 2023 IEEE Guwahati Subsection Conference 575 (GCON), pp. 1-6. IEEE, 2023.
- Shehzaad Zuzar Dhuliawala, Mrinmaya Sachan, and Carl Allen. Variational classification: A 577 probabilistic generalization of the softmax classifier. Transactions on Machine Learning Research, 578 2023. 579
  - Dmytro Dosyn, Yousef Ibrahim Daradkeh, Vira Kovalevych, Mykhailo Luchkevych, and Yaroslav Kis. Domain ontology learning using link grammar parser and wordnet. In *MoMLeT*+ DS, pp. 14-36, 2022.
- 583 Xuefeng Du, Xin Wang, Gabriel Gozum, and Yixuan Li. Unknown-aware object detection: Learning 584 what you don't know from videos in the wild. In Proceedings of the IEEE/CVF Conference on 585 Computer Vision and Pattern Recognition, pp. 13678–13688, 2022. 586
- Christiane Fellbaum. Wordnet. In Theory and applications of ontology: computer applications, pp. 587 231–243. Springer, 2010. 588
- 589 Jiazhan Feng, Qingfeng Sun, Can Xu, Pu Zhao, Yaming Yang, Chongyang Tao, Dongyan Zhao, 590 and Qingwei Lin. Mmdialog: A large-scale multi-turn dialogue dataset towards multi-modal 591 open-domain conversation. arXiv preprint arXiv:2211.05719, 2022. 592
- Stanislav Fort, Jie Ren, and Balaji Lakshminarayanan. Exploring the limits of out-of-distribution detection. Advances in Neural Information Processing Systems, 34:7068–7081, 2021.

594 595 596	Rena Gao, Carsten Roever, and Jey Han Lau. Interaction matters: An evaluation framework for interactive dialogue assessment on english second language conversations. <i>arXiv preprint arXiv:2407.06479</i> , 2024a.
597 598 599 600	Rena Gao, Jingxuan Wu, Carsten Roever, Xuetong Wu, Jing Wu, Long Lv, and Jey Han Lau. Cnima: A universal evaluation framework and automated approach for assessing second language dialogues. <i>arXiv preprint arXiv:2408.16518</i> , 2024b.
601 602 603	Wei Gao and Menghan Wang. Listenership always matters: active listening ability in l2 business english paired speaking tasks. <i>International Review of Applied Linguistics in Language Teaching</i> , (0), 2024.
605 606 607	Mark S Graham, Walter HL Pinaya, Petru-Daniel Tudosiu, Parashkev Nachev, Sebastien Ourselin, and Jorge Cardoso. Denoising diffusion models for out-of-distribution detection. In <i>Proceedings</i> of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2947–2956, 2023.
608 609 610	Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. <i>arXiv preprint arXiv:1610.02136</i> , 2016.
611 612 613	Dan Hendrycks, Steven Basart, Mantas Mazeika, Mohammadreza Mostajabi, Jacob Steinhardt, and Dawn Song. A benchmark for anomaly segmentation. <i>arXiv preprint arXiv:1911.11132</i> , 1(2):5, 2019.
614 615 616 617	Dan Hendrycks, Xiaoyuan Liu, Eric Wallace, Adam Dziedzic, Rishabh Krishnan, and Dawn Song. Pretrained transformers improve out-of-distribution robustness. <i>arXiv preprint arXiv:2004.06100</i> , 2020.
618 619 620	Yen-Chang Hsu, Yilin Shen, Hongxia Jin, and Zsolt Kira. Generalized odin: Detecting out-of- distribution image without learning from out-of-distribution data. In <i>Proceedings of the IEEE/CVF</i> <i>conference on computer vision and pattern recognition</i> , pp. 10951–10960, 2020.
621 622 623 624	Anna Huang et al. Similarity measures for text document clustering. In <i>Proceedings of the sixth new zealand computer science research student conference (NZCSRSC2008), Christchurch, New Zealand</i> , volume 4, pp. 9–56, 2008.
625 626 627	Ammar Ismael Kadhim, Yu-N Cheah, Nurul Hashimah Ahamed, and Lubab A Salman. Feature extraction for co-occurrence-based cosine similarity score of text documents. In 2014 IEEE student conference on research and development, pp. 1–4. IEEE, 2014.
628 629 630	Satwik Kottur, José MF Moura, Devi Parikh, Dhruv Batra, and Marcus Rohrbach. Clevr-dialog: A diagnostic dataset for multi-round reasoning in visual dialog. <i>arXiv preprint arXiv:1903.03166</i> , 2019.
632 633 634 635	Ira Ktena, Olivia Wiles, Isabela Albuquerque, Sylvestre-Alvise Rebuffi, Ryutaro Tanno, Abhijit Guha Roy, Shekoofeh Azizi, Danielle Belgrave, Pushmeet Kohli, Taylan Cemgil, et al. Generative models improve fairness of medical classifiers under distribution shifts. <i>Nature Medicine</i> , pp. 1–8, 2024.
636 637 638 639	Alfirna Rizqi Lahitani, Adhistya Erna Permanasari, and Noor Akhmad Setiawan. Cosine similarity to determine similarity measure: Study case in online essay assessment. In 2016 4th International conference on cyber and IT service management, pp. 1–6. IEEE, 2016.
640 641	Hao Lang, Yinhe Zheng, Binyuan Hui, Fei Huang, and Yongbin Li. Out-of-domain intent detection considering multi-turn dialogue contexts. <i>arXiv preprint arXiv:2305.03237</i> , 2023.
642 643 644 645	Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. <i>Advances in neural information processing systems</i> , 31, 2018.
646 647	Nyoungwoo Lee, Suwon Shin, Jaegul Choo, Ho-Jin Choi, and Sung-Hyun Myaeng. Constructing multi-modal dialogue dataset by replacing text with semantically relevant images. <i>arXiv preprint arXiv:2107.08685</i> , 2021.

648	Baoli Li and Lining Han. Distance weighted cosine similarity measure for text classification. In Intel
649	ligent Data Engineering and Automated Learning_IDFAL 2013: 14th International Conference
650	IDEAL 2013. Hefei, China, October 20-23, 2013. Proceedings 14, pp. 611–618. Springer, 2013.
651	
652	Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-
653	training for unified vision-language understanding and generation. In International conference on
654	<i>machine learning</i> , pp. 12888–12900. PMLR, 2022.
655	Minghan Li and Jimmy Lin. Encoder adaptation of dense passage retrieval for open-domain question
656	answering arXiv preprint arXiv:2110.01599 2021
657	
658	Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. Dailydialog: A manually
659	labelled multi-turn dialogue dataset. arXiv preprint arXiv:1710.03957, 2017.
660	Shivu Liang Viyuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution
661	image detection in neural networks arXiv preprint arXiv:1706.02690.2017
662	
663	Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
664	Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer Vision-
665	ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings,
666	<i>Part V 13</i> , pp. 740–755. Springer, 2014.
667	Hong Liu Jeff Z HaoChen Adrien Gaidon and Tengyu Ma Self-supervised learning is more robust
868	to dataset imbalance arXiv preprint arXiv:2110.05025, 2021
000	
670	Xixi Liu, Yaroslava Lochman, and Christopher Zach. Gen: Pushing the limits of softmax-based
671	out-of-distribution detection. In Proceedings of the IEEE/CVF Conference on Computer Vision
672	and Pattern Recognition, pp. 23946–23955, 2023.
672	Jacek Marciniak Wordnet as a backbone of domain and application conceptualizations in systems
674	with multimodal data. In <i>Proceedings of the LREC 2020 Workshop on Multimodal Wordnets</i>
675	( <i>MMW2020</i> ), pp. 25–32, 2020.
676	
677	Philippe Martin. Using the wordnet concept catalog and a relation hierarchy for knowledge acquisition.
678	In Proc. 4th Petrce Workshop, 1995.
679	Yihan Mei, Xinyu Wang, Dell Zhang, and Xiaoling Wang. Multi-label out-of-distribution detection
680	with spectral normalized joint energy. arXiv preprint arXiv:2405.04759, 2024.
681	
682	A. H. Miller, W. Feng, A. Fisch, J. Lu, D. Batra, A. Bordes, D. Parikh, and J. Weston. Parial: A
683	dialog research software platform. arxiv preprint arxiv:1703.06476, 2017.
684	Atsuyuki Miyai, Qing Yu, Go Irie, and Kiyoharu Aizawa. Locoop: Few-shot out-of-distribution
685	detection via prompt learning. Advances in Neural Information Processing Systems, 36, 2024.
686	
687	Yulei Niu and Hanwang Zhang. Introspective distillation for robust question answering. Advances in Naural Information Processing Systems 24:16202, 16204, 2021
688	Neurai Information Processing Systems, 54:10292–10504, 2021.
689	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
690	Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
691	models from natural language supervision. In International conference on machine learning, pp.
692	8748–8763. PMLR, 2021.
693	Faisal Rahutomo Tanuaki Kitasuka Masayoshi Aritsugi at al. Samantia cosina similarity. In The 7th
694	international student conference on advanced science and technology ICAST volume 4 pp. 1
695	University of Seoul South Korea 2012
696	
697	Alexandre Ramé, Kartik Ahuja, Jianyu Zhang, Matthieu Cord, Léon Bottou, and David Lopez-Paz.
698	Model ratatouille: Recycling diverse models for out-of-distribution generalization. In <i>International</i>
699	Conference on Machine Learning, pp. 28656–28679. PMLR, 2023.
700	Amaia Salvador, Nicholas Hynes, Yusuf Aytar, Javier Marin, Ferda Ofli, Ingmar Weber, and Antonio
701	Torralba. Learning cross-modal embeddings for cooking recipes and food images. In <i>Proceedings</i>
	of the IEEE conference on computer vision and pattern recognition, pp. 3020–3028, 2017.

702	Paul Hongsuck Seo, Andreas Lehrmann, Bohyung Han, and Leonid Sigal. Visual reference resolution
703	using attention memory for visual dialog. Advances in neural information processing systems, 30.
704	2017.
705	

- Erik Wallin, Lennart Svensson, Fredrik Kahl, and Lars Hammarstrand. Improving open-set semisupervised learning with self-supervision. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 2356–2365, 2024.
- Haoran Wang, Weitang Liu, Alex Bocchieri, and Yixuan Li. Can multi-label classification networks
  know what they don't know? *Advances in Neural Information Processing Systems*, 34:29074–29087, 2021.
- Lei Wang, Sheng Huang, Luwen Huangfu, Bo Liu, and Xiaohong Zhang. Multi-label out-of-distribution detection via exploiting sparsity and co-occurrence of labels. *Image and Vision Computing*, 126:104548, 2022.
- Qizhou Wang, Zhen Fang, Yonggang Zhang, Feng Liu, Yixuan Li, and Bo Han. Learning to augment distributions for out-of-distribution detection. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yunchao Wei, Wei Xia, Min Lin, Junshi Huang, Bingbing Ni, Jian Dong, Yao Zhao, and Shuicheng Yan. Hcp: A flexible cnn framework for multi-label image classification. *IEEE transactions on pattern analysis and machine intelligence*, 38(9):1901–1907, 2015.
- Xuetong Wu, Jonathan H Manton, Uwe Aickelin, and Jingge Zhu. A bayesian approach to (online)
   transfer learning: Theory and algorithms. *Artificial Intelligence*, 324:103991, 2023.
- Xuetong Wu, Jonathan H Manton, Uwe Aickelin, and Jingge Zhu. On the generalization for transfer
   learning: An information-theoretic analysis. *IEEE Transactions on Information Theory*, 2024.
- Linyi Yang, Shuibai Zhang, Libo Qin, Yafu Li, Yidong Wang, Hanmeng Liu, Jindong Wang, Xing Xie, and Yue Zhang. Glue-x: Evaluating natural language understanding models from an out-of-distribution generalization perspective. *arXiv preprint arXiv:2211.08073*, 2022.
- Hai Ye, Yuyang Ding, Juntao Li, and Hwee Tou Ng. Robust question answering against distribution shifts with test-time adaptation: An empirical study. *arXiv preprint arXiv:2302.04618*, 2023.
- Xintong Yu, Hongming Zhang, Yangqiu Song, Yan Song, and Changshui Zhang. What you see is what
   you get: Visual pronoun coreference resolution in dialogues. *arXiv preprint arXiv:1909.00421*, 2019.
- Lifan Yuan, Yangyi Chen, Ganqu Cui, Hongcheng Gao, Fangyuan Zou, Xingyi Cheng, Heng Ji,
   Zhiyuan Liu, and Maosong Sun. Revisiting out-of-distribution robustness in nlp: Benchmarks,
   analysis, and Ilms evaluations. *Advances in Neural Information Processing Systems*, 36, 2024.
- Dell Zhang and Bilyana Taneva-Popova. A theoretical analysis of out-of-distribution detection in multi-label classification. In *Proceedings of the 2023 ACM SIGIR International Conference on Theory of Information Retrieval*, pp. 275–282, 2023.
- Haotian Zheng, Qizhou Wang, Zhen Fang, Xiaobo Xia, Feng Liu, Tongliang Liu, and Bo Han.
   Out-of-distribution detection learning with unreliable out-of-distribution sources. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yinhe Zheng, Guanyi Chen, and Minlie Huang. Out-of-domain detection for natural language understanding in dialog systems. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:1198–1209, 2020.

751

752

753

- 754
- 755

# 756 A EXPERIMENT DETAILS

#### The dataset stats are summarized as follows:

759		01101101				
760	Table 7:	: Statistics o	f Visdial QA da	ataset		
761						
762	Stats	Matched	Mismatched	ID	OOD	
763	# Pair	122K	10K	95K	37K	
764	# Train	77K	0	77K	0	
765	# Test	45K	10K	18K	37K	
766	# Turn per dialog		10			
767	# Categories		80			
768	# Supercategories		12			
769						
770	Table 8:	: Statistics o	f Real MMD da	ataset		
772	Stats	Matched	Mismatched	ID	OOD	
773	# Pair	17K	8K	12.7K	12.2K	
774	# Train	10.2K	0	10.2K	0	
775	# Test	14.7K	8K	2.5K	12.2K	
776	# Turn per dialog		$5 \sim 15$			
777	# Categories		80			
779	# Supercategories		12			
781 782	Ta'	ble 9: Exper	imental Details			
703	Parameters	Conf	igurations			
704	$\gamma$	1				
700	$\alpha$	0.5			( D	
700	Image Encoder Dialogue Encoder		VI-1 B/32 or E	SLIP III DI ID ITN	A Base	
707	$s(x_{\rm L}, x_{\rm T})$	CLIP	vi-1 D/32 01 f	DLIF III	vi Dase	
780	Label Extractor	5-La	ver DNN with	size [51	2/256. 25	6, 128, 64,
700		111	)		_,,	.,,,
701	Activation Function	Relu	& Sigmoid			
702	Batch Size	32	C			
792	Learning Rate	0.001	l			
793	Optimizer	Adan	n			
794	ID label	perso	on, kitchen, food	, sports,	electronic,	accessory,
795	0001111	furni	ture, indoor, ap	pliance,	vehicle, or	ıtdoor
707	ood label		ai			
708	T in ODIN	0.001	L			
700	Image Features in Mahalanohis		P/BI IP(Image)			
800	Text Features in Mahalanobis	CLIP	/BLIP(Dialogu	e)		
500		2211	gu	- /		

# **B** THEORETICAL JUSTIFICATION

**Assumption 1** We denote ID distribution as  $P(x_I, x_T, y)$  and OOD distribution as  $\tilde{P}(x_I, x_T, y)$  where  $\tilde{P}$  may differ from P in terms of the following assumptions.

• *Case 1: The image and text match, but labels are out of the set, namely:* 

 $\mathbb{E}_{P(x_I, x_T)} \left[ \log s(x_I, x_T) \right] = \mathbb{E}_{\tilde{P}(x_I, x_T)} \left[ \log s(x_I, x_T) \right],$ 

810	For every pair $x_I$ , $x_T$ and any $\alpha$ ,
811	$\mathbb{E}_{\mathcal{D}(x)} = \sum \left[ \log(\alpha s_t(x_t, y) + (1 - \alpha) s_{\mathcal{D}}(x_{\mathcal{D}}, y)) \right] > \mathbb{E}_{\mathcal{D}(x)} = \sum \left[ \log(\alpha s_t(x_t, y) + (1 - \alpha) s_{\mathcal{D}}(x_{\mathcal{D}}, y)) \right]$
813	$\mathbb{E}_{P(y x_{I},x_{T})}\left[\log(\alpha \sigma_{I}(x_{I},y) + (1 - \alpha)\sigma_{I}(x_{I},y))\right] \ge \mathbb{E}_{P(y x_{I})}\left[\log(\alpha \sigma_{I}(x_{I},y) + (1 - \alpha)\sigma_{I}(x_{I},y))\right],$
814	which means that the ID pairs $x_I$ and $x_T$ should have stronger expressity about the ID label
815	y than OOD pairs.
816	• Case 2: The image and text do not match, which we assume:
817 818	$\mathbb{E}_{P(x_I, x_T)} \left[ \log s(x_I, x_T) \right] > \mathbb{E}_{\tilde{P}(x_I, x_T)} \left[ \log s(x_I, x_T) \right],$
819 820	which means the ID pairs should have higher similarity than OOD pairs in this case. For every pair $x_I$ , $x_T$ and any $\alpha$ ,
821 822	$\mathbb{E}_{P(y x_{I},x_{T})}\left[\log(\alpha s_{I}(x_{I},y) + (1-\alpha)s_{T}(x_{T},y))\right] = \mathbb{E}_{\tilde{P}(y x_{I},x_{T})}\left[\log(\alpha s_{I}(x_{I},y) + (1-\alpha)s_{T}(x_{T},y))\right],$
823 824	which means that some ID pairs $x_I$ and $x_T$ may have the same expressity about the label y compared with the OOD pairs.
825 826	• Case 3: The image or text does not match with the labels, which we assume:
827 828	$\mathbb{E}_{P(y x_I, x_T)} \left[ \log(\alpha s_I(x_I, y) + (1 - \alpha) s_T(x_T, y)) \right] > \mathbb{E}_{\tilde{P}(y x_I, x_T)} \left[ \log(\alpha s_I(x_I, y) + (1 - \alpha) s_T(x_T, y)) \right].$
829	<b>Theorem 1</b> With Assumption 1, we can show that the proposed DIEAF score satisfies the following:
830 831	$\mathbb{E}_{ ilde{P}(x_I,x_T,y)}[\log g(x_I,x_T,y)] < \mathbb{E}_{P(x_I,x_T,y)}[\log g(x_I,x_T,y)].$
832 833	<b>Proof 1</b> It is easy to write that:
834	$\mathbb{E}[\log g(x_I, x_T, y)] = \gamma \mathbb{E}_{x_I, x_T}[\log s(x_I, x_T)] + \mathbb{E}_{x_I, x_T} \mathbb{E}_{y x_I, x_T}[\log([\alpha s_I(x_I, y) + (1 - \alpha)s_T(x_T, y)])].$
836 837 838 839	The proof simply follows the assumptions we made for each case. Note that this score only works for positive scores, but sometimes, we may encounter negative scores, and the log may be ill-posed. As a surrogate score function, we eliminate the log and maintain $g(x_I, x_T, y)$ for the same intuition as the above theorem.
840	
841	
842	
843	
844	
845	
040 9/17	
848	
849	
850	
851	
852	
853	
854	
855	
856	
857	
858	
859	
860	
100	
863	
500	