PolCLIP: A Unified Image-Text Word Sense Disambiguation Model via Generating Multimodal Complementary Representations

Anonymous ACL submission

Abstract

Word sense disambiguation (WSD) is divided into two subtasks: textual word sense disambiguation (Textual-WSD) and visual word sense disambiguation (Visual-WSD). They aim to identify the most semantically relevant senses or images to a given context containing ambiguous target words. However, existing WSD models seldom address these two subtasks jointly due to lack of images in Textual-WSD datasets or lack of senses in Visual-WSD datasets. To bridge this gap, we propose Pol-CLIP, a unified image-text WSD model. By employing an image-text complementarity strategy, it simulates stable diffusion to generate implicit visual representations for senses and imitates image captioning to provide implicit textual representations for images. Additionally, a disambiguation-oriented image-sense dataset is constructed for the training objective of learning multimodal polysemy representations. To the best of our knowledge, PolCLIP is the first model that can cope with both Textual-WSD and Visual-WSD. Extensive experimental results on benchmarks demonstrate the effectiveness of our method, achieving a 2.53% F1score increase over the state-of-the-art models on Textual-WSD and a 2.22% HR@1 improvement on Visual-WSD.

1 Introduction

004

800

011

012

014

018

023

034

042

Understanding and identifying the intended meaning of words with multiple senses (i.e., polysemy) is a significant challenge in natural language processing (Navigli, 2009). This promotes in-depth research on word sense disambiguation (WSD), which has recently been extended to multimodal downstream tasks (Bevilacqua et al., 2021). Techniques for WSD are critical for enhancing the accuracy and effectiveness of text understanding and information retrieval tasks such as machine translation (Raganato et al., 2019), image-text retrieval (Chen et al., 2020), and large language model inference (Kritharoula et al., 2023).





Figure 1: Illustration of the Multimodal-WSD task.

043

044

045

047

050

051

054

060

061

063

064

065

066

067

069

070

Theoretically, WSD can be divided into two subtasks: textual word sense disambiguation (Textual-WSD) (Bevilacqua et al., 2021) and visual word sense disambiguation (Visual-WSD) (Raganato et al., 2023). Given a context containing an ambiguous target word, the goal of Textual-WSD is to select the most semantically appropriate one from a set of candidate senses, while the goal of Visual-WSD is to choose the most semantically suitable one from a set of candidate images. Due to the distinct modalities, these two subtasks typically require specialized training datasets and methods (Bevilacqua and Navigli, 2020; Blevins and Zettlemoyer, 2020; Kwon et al., 2023). Nevertheless, they can be unified as a Multimodal-WSD task if the senses and images in existing WSD datasets are aligned. As shown in Figure 1, the task objective of Multimodal-WSD is to identify both the most semantically correct senses and images. Technically, developing a generic Multimodal-WSD model can realize the unification of WSD tasks and activate the potential of multimodal applications in understanding polysemy knowledge.

In the Textual-WSD datasets (Raganato et al., 2017), only textual senses serve as candidates (i.e., image-missing), while in the Visual-WSD datasets (Raganato et al., 2023), only images serve as candidates (i.e., sense-missing). This results in the

084

097

101

102

103

105

106

108

109

110

111

112

113

114

115

116

117

118

071

072

073

challenge of modality missing at the data level, limiting the unification of these two WSD subtasks. Furthermore, multimodal representations have been demonstrated to carry richer semantic information compared to unimodal representations in recent WSD works (Gella et al., 2016, 2019). However, constrained by model architecture, existing Textual-WSD models (Conia and Navigli, 2021; Maru et al., 2019; Huang et al., 2019) cannot supplement candidate senses with image information, and Visual-WSD models (Yang et al., 2023; Zhang et al., 2023; Dadas, 2023) cannot supplement candidate images with descriptions. This poses technical difficulties in developing a unified framework for Multimodal-WSD.

To address these issues, we propose PolCLIP, a unified image-text WSD model which is proficient in multimodal polysemy processing and is built upon CLIP (Radford et al., 2021) architecture. By employing an image-text complementarity strategy, it can simulate the stable diffusion (Ho et al., 2020) (generating images based on texts) and image captioning (Ramos et al., 2023) (generating descriptions based on images). The core idea of this strategy is to make PolCLIP initially focus on the key information of original unimodal senses or images, and then re-utilize the text or image encoder to generate implicit image-text complementary representations. Two widely used WSD datasets (Semcor (Miller et al., 1993) and VWSD-KB (Yang et al., 2023)) are integrated into a disambiguationoriented image-sense dataset for the training objective of learning aligned multimodal representations. Moreover, a fine-tuned GPT-3.5 model is utilized to generate lexical definitions for semantic enhancement in testing phase. The main contributions of this work can be summarized as follows:

- A unified image-text WSD model is proposed, which is the first model to jointly cope with Textual-WSD lacking images and Visual-WSD lacking senses.
- An image-text complementarity strategy is introduced to simulate stable diffusion and image captioning for addressing the modality missing issues in unimodal WSD datasets.
- A disambiguation-oriented image-sense dataset is constructed to provide a benchmark for the Multimodal-WSD task.

2 Related Work

Textual-WSD was mainly tackled by knowledgebased methods and supervised methods (Bevilacqua et al., 2021). Knowledge-based methods (Maru et al., 2019; Scozzafava et al., 2020) typically used external dictionary resources to provide sense lists for ambiguous words to resolve polysemy. Supervised methods (Huang et al., 2019; Wang and Wang, 2020) generally used pre-training language models to maximize the similarity probabilities between contexts and candidate senses in a feature space. BEM (Blevins and Zettlemoyer, 2020) and SACE (Wang and Wang, 2021) adopted biencoders and only retained the representations corresponding to ambiguous words. They achieved state-of-the-art results on English all-words benchmarks at that time. Moreover, the full utilization of visual features for verb sense disambiguation has attracted increasing interest (Gella et al., 2016, 2019). EViLBERT (Calabrese et al., 2020b) obtained better results by learning task-agnostic multimodal sense representations, compared to methods built solely on language models. Although these methods primarily leverage visual information to bolster performance on Textual-WSD, but they could not be applied to Visual-WSD straightforwardly.

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

Visual-WSD was introduced in SemEval-2023 Task 1 (Raganato et al., 2023). The Visual-WSD mainstream approaches employed Vision-Language Pre-training models (VLPs) for imagetext retrieval. FCLL (Yang et al., 2023) proposed a fine-grained image-text contrastive learning mechanism and won first place in SemEval-2023 Task 1. Moreover, large language models (LLMs) were widely used to enrich the semantic information of contexts (Ghahroodi et al., 2023). Calling APIs was a commonly adopted strategy, where simple prompts were designed to guide LLMs to return the interpretations of ambiguous target words in contexts (Kritharoula et al., 2023). However, Visual-WSD models depended on the prior knowledge of VLPs, which are pre-trained with objectives biased towards image-text understanding rather than WSD. This results in Visual-WSD models struggling to effectively resolve Textual-WSD.

Data is a key factor to unify these two WSD subtasks and develop a generic Multimodal-WSD model. WordNet (Miller et al., 1990) is a large lexicographic database and a standard inventory for English WSD. It contains approximately 120,000 synsets. BabelNet (Navigli and Ponzetto, 2010;



Figure 2: Overview framework of the proposed PolCLIP model.

Navigli et al., 2021) is the most popular multilin-170 gual dictionary, which can be semi-automatically 171 mapped to other resources to acquire encyclope-172 dic terms. It covers over 500 languages and was upgraded to version 5.3 recently. Researchers can 174 access various possible resources about ambiguous 175 words by BabelNet, including example sentences, 176 parts of speech, textual senses, and images. By linking WordNet with Wikipedia through Babel-178 Net, BabelPic (Calabrese et al., 2020a) expanded non-concrete image-sense pairs, paving the way 180 181 for our work to construct larger disambiguationoriented multimodal datasets. 182

3 Method

183

201

3.1 Task formulation

Textual-WSD and Visual-WSD can be unified as a Multimodal-WSD task which is a token classifi-187 cation problem. A given context c generally contains at least one ambiguous target word w_{target} . For Textual-WSD, there is a set of word senses $S = \{s_1, s_2, \dots, \hat{s}, \dots, s_n\}$ as candidates, where \hat{s} denotes the most semantically relevant sense to 191 w_{target} . Following Eq. 1, a Textual-WSD model is 192 required to learn a similarity function F to retrieve 193 \hat{s} from candidate senses. For Visual-WSD, there is a group of images $I = \left\{i_1, i_2, \dots, \hat{i}, \dots, i_n\right\}$ as 195 candidates, where \hat{i} represents the most semanti-196 cally relevant image to w_{target} . Following Eq. 2, a 197 Visual-WSD model is required to learn a similarity 198 function F to retrieve i from candidate images. 199

$$\hat{s} = \arg\max F\left(c, w^{target}, S\right) \tag{1}$$

$$\hat{i} = \arg \max F(c, w^{target}, I)$$

3.2 The PolCLIP model

The framework of the PolCLIP model is shown in Figure 2. It utilizes an image-text complementarity strategy and is built upon CLIP architecture. It employs 12-layer transformers as the text encoder and 24-layer visual transformers as the image encoder. A context with an ambiguous target word w_{target} is input into the text encoder to generate a complete context representation $e_c = \{[CLS], e^{w_1}, \dots, e^{w_{\text{target}}}, \dots, e^{w_n}, [SEP]\},$ where $e^{w_{\text{target}}}$ is the representation corresponding to w_{target} . The candidate senses of w_{target} are input into the text encoder to output a sense vector v_{sense} . The candidate images of w_{target} are fed into the image encoder to output an image vector v_{image} . 203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

225

226

227

228

229

230

231

232

233

234

236

A Sense-to-Image Generation (SIG) module and an Image-to-Sense Generation (ISG) module are designed to generate implicit image-text complementary information, which is expressed in vectors. Specifically, the SIG module consists of a shared 2-layer self-attention module, a sense-toimage projector, and the last four layers of the image encoder. The sense vector v_{sense} is input to SIG. Its key information is condensed through the self-attention module and the linear-layer projector. After that, the highly compressed sense information is transformed into a sense-generated image vector $v_{\text{image}}^{\text{sense-gen}}$ with visual knowledge through the last four layers of the image encoder. Similarly, the ISG module is composed of the same shared 2-layer self-attention module, an image-to-sense projector, and the last four layers of the text encoder. The image vector v_{image} is input to ISG and then converted into an image-generated sense

(2)

vector $v_{\text{sense}}^{\text{image-gen}}$ with textual knowledge.

237

238

241

243

246

247

251

252

253

257

260

261

262

264

265

In order to make the generated implicit imagetext complementary information actually beneficial to semantic augment, following ALBEF (Li et al., 2021), a text momentum encoder and an image momentum encoder are employed to generate pseudotarget vectors. Specifically, the candidate senses and images from the same batch are input separately into the text and image momentum encoders to output a pseudo-target sense vector $v_{\text{sense}}^{\text{pseudo}}$ and a pseudo-target image vector $v_{\text{image}}^{\text{pseudo}}$. These two pseudo-target vectors can supervise $v_{sense}^{image-gen}$ and $v_{\text{image}}^{\text{sense-gen}}$ to be close to ground truth. The text and image momentum encoders can retain the prior knowledge of the backbone model to counteract the issue of catastrophic forgetting (Li et al., 2023). Therefore, the pseudo-target vectors gain improvement continuously with these two momentum encoders being optimized at a small pace. The similarity between $v_{\text{sense}}^{\text{image-gen}}$ and $v_{\text{sense}}^{\text{pseudo}}$ can be calculated by Eq. 3-4. Also, the similarity between $v_{\text{image}}^{\text{sense-gen}}$ and $v_{\text{image}}^{\text{pseudo}}$ can be calculated by Eq. 5-6. s is a similarity function. \mathcal{P}^{ISG} and \mathcal{P}^{SIG} are the softmax-normalized similarities used to supervise the ISG module and the SIG module.

$$s(S^{\text{gen}}, S^{\text{pse}}) = v_{\text{sense}}^{\text{image-gen}} \cdot (v_{\text{sense}}^{\text{pseudo}})^{\text{T}}$$
 (3)

$$\mathcal{P}^{ISG} = \frac{\exp\left(s\left(S^{\text{gen}}, S^{\text{pse}}\right)\right)}{\sum_{n=1}^{N} \exp\left(s\left(S^{\text{gen}}, S^{\text{pse}}\right)\right)} \quad (4)$$

$$s(I^{\text{gen}}, I^{\text{pse}}) = v_{\text{image}}^{\text{sense-gen}} \cdot \left(v_{\text{image}}^{\text{pseudo}}\right)^{\text{T}}$$
 (5)

$$\mathcal{P}^{SIG} = \frac{\exp\left(s\left(I^{\text{gen}}, I^{\text{pse}}\right)\right)}{\sum_{n=1}^{N} \exp\left(s\left(I^{\text{gen}}, I^{\text{pse}}\right)\right)} \qquad (6)$$

A shared 4-layer cross-attention module serves as a fusion module. It integrates the original unimodal sense/image representations and the generated implicit image/sense representations into semantically enriched multimodal representations. v_{sense} serves as Q and $v_{\text{image}}^{\text{sense-gen}}$ serves as K and 274 V. They are fed into the fusion module and then a 275 sense-guided multimodal vector $v_{\text{multi}}^{\text{sense-gui}}$ is calcu-276 lated by softmax $\left(\frac{QK^{T}}{\sqrt{d_{k}}}V\right)$, where d_{k} denotes the dimension of 768. This $v_{\text{multi}}^{\text{sense-gui}}$ achieves an effectation 277 tive interaction of the original sense representation with the implicit sense-generated image representation. Similarly, v_{image} serves as Q and $v_{\text{sense}}^{\text{image-gen}}$ 281 serves as K and V. They are fed into the fusion module and then an image-guided multimodal vector $v_{\rm multi}^{\rm image-gui}$ is output by the same cross-attention

calculation process. This $v_{\text{multi}}^{\text{image-gui}}$ achieves an effective interaction of the original image representation with the implicit image-generated sense representation.

285

286

287

289

290

291

294

295

296

297

298

299

302

303

304

305

306

307

308

310

311

312

313

314

315

316

317

318

319

320

321

325

326

327

330

331

To avoid the key information of the ambiguous target word w_{target} being smoothed out, we directly select the representation of w_{target} as the anchor vector for retrieval, instead of simply averaging the complete context representation e_c or taking [CLS]. Following Eq. 7-8, this anchor $e^{w_{\text{target}}}$ is used to calculate the similarity with the sense-guided multimodal vector $v_{\text{multi}}^{\text{sense-gui}}$. The similarity between the anchor $e^{w_{\text{target}}}$ and the image-guided multimodal vector $v_{\text{multi}}^{\text{sense-gui}}$ is calculated by Eq. 9-10. \mathcal{P}^{W2S} and \mathcal{P}^{W2I} are the softmax-normalized anchor-to-sense similarity and the anchor-to-image similarity. The PolCLIP model can identify the most semantically appropriate senses and images based on these two similarities.

$$s(W,S) = e^{w_{\text{target}}} \cdot \left(v_{\text{multi}}^{\text{sense-gui}}\right)^{\text{T}}$$
(7)

$$\mathcal{P}^{W2S} = \frac{\exp(s(W,S))}{\sum_{n=1}^{N} \exp(s(W,S))}$$
(8)

$$s(W, I) = e^{w_{\text{target}}} \cdot \left(v_{\text{multi}}^{\text{image-gui}}\right)^{\text{T}}$$
(9)

$$\mathcal{P}^{W2I} = \frac{\exp(s(W, I))}{\sum_{n=1}^{N} \exp(s(W, I))}$$
(10)

Four contrastive losses (Hadsell et al., 2006) are defined to optimize four training objectives jointly, comprising two generation-based objectives (SIG loss and ISG loss) and two understanding-based objectives (W2S loss and W2I loss).

The two generation-based objectives make the generated implicit image-text complementary information close to ground truth, to ensure that the enriched multimodal representations are semantically correct. The contrastive loss \mathcal{L}_{SIG} is defined as a cross-entropy \mathcal{H} between the sense-generated image vector $v_{image}^{sense-gen}$ and the pseudo-target image vector $v_{image}^{sense-gen}$:

$$\mathcal{L}_{SIG} = \mathbb{E}_{(I^{\text{gen}}, I^{\text{pse}}) \sim D} \mathcal{H} \left(\mathcal{Y}^{SIG}, \mathcal{P}^{SIG} \right) \quad (11)$$

 \mathcal{Y} indicates the ground-truth multi-label one-hot similarity, where negative pairs have a probability of 0 and the positive pairs have a probability of 1. Similarly, the image-generated sense vector $v_{\text{sense}}^{\text{image-gen}}$ and the pseudo-target sense vector $v_{\text{sense}}^{\text{pseudo}}$ are used to calculate \mathcal{L}_{ISG} :

$$\mathcal{L}_{ISG} = \mathbb{E}_{(S^{\text{gen}}, S^{\text{pse}}) \sim D} \mathcal{H} \left(\mathcal{Y}^{ISG}, \mathcal{P}^{ISG} \right) \quad (12)$$

ŀ	Algorithm 1: Pseudocode of PolCLIP Inference				
	input : a augmented context c with an ambiguous target word w_{target} ;				
	the candidates Cand with senses or images;				
	output: the ranked candidates Cand _{ranked} ;				
	the semantically optimal sense or image O_{best} ;				
1	$e_c \leftarrow \text{Text-Encoder}(c)$; # the complete context representations				
2	$e^{w_{\text{target}}} \leftarrow e_c$; # the anchor vector based on w_{target}				
3	if only senses in Cand then				
4	# for Textual-WSD				
5	$v_{\text{sense}} \leftarrow \text{Text_Encoder}(Cand);$				
6	# the sense-generated image representations				
7	$v_{\text{image}}^{\text{sense-gen}} \leftarrow \text{SIG}(v_{\text{sense}});$				
8	# the sense-guided multimodal representations				
9	$v_{\text{multi}} \leftarrow \text{Fusion}(v_{\text{sense}}, v_{\text{image}}^{\text{sense-gen}});$				
10	else				
11	# for Visual-WSD				
12	$v_{\text{image}} \leftarrow \text{Image_Encoder}(Cand);$				
13	# the image-generated sense representations				
14	$v_{\text{sense}}^{\text{image-gen}} \leftarrow \text{ISG}(v_{\text{image}});$				
15	# the image-guided multimodal representations				
16	$v_{\text{multi}} \leftarrow \text{Fusion}(v_{\text{image}}, v_{\text{sense}}^{\text{image-gen}});$				
17	end				
18	$similarity \leftarrow e^{w_{\text{target}}} \cdot (v_{\text{multi}})^T;$				
19	$Cand_{ranked} \leftarrow top_k(similarity); # k is the number of candidates$				
20	$O_{\text{best}} \leftarrow \arg \max(Cand_{\text{ranked}}); #$ the semantically optimal sense or image				

The two understanding-based objectives ensure that the PolCLIP model can accurately identify the semantically optimal senses and images. The anchor $e^{w_{\text{target}}}$ and the sense-guided multimodal vector $v_{\text{multi}}^{\text{sense-gui}}$ are used to calculate \mathcal{L}_{W2S} :

$$\mathcal{L}_{W2S} = \mathbb{E}_{(W,S)\sim D} \mathcal{H} \left(\mathcal{Y}^{W2S}, \mathcal{P}^{W2S} \right)$$
(13)

Also, the anchor $e^{w_{\text{target}}}$ and the image-guided multimodal vector $v_{\text{multi}}^{\text{image-gui}}$ are used to calculate \mathcal{L}_{W2I} :

$$\mathcal{L}_{W2I} = \mathbb{E}_{(W,I)\sim D} \mathcal{H} \left(\mathcal{Y}^{W2I}, \mathcal{P}^{W2I} \right)$$
(14)

Finally, the full training objective of PolCLIP is:

$$\mathcal{L} = \mathcal{L}_{SIG} + \mathcal{L}_{ISG} + \mathcal{L}_{W2S} + \mathcal{L}_{W2I} \qquad (15)$$

The pesudocode of training the PolCLIP model is provided in Appendix A.

3.3 Inference of the PolCLIP model

To further stimulate the potential of the PolCLIP model in understanding polysemy text, a semantic enhancement is implemented for contexts during the inference procedure. Different to methods that call APIs, we develop a disambiguation-oriented GPT-3.5 (D-GPT) to generate intended lexical definitions of a word in contexts. Fine-tuning on a random selection of 50,000 data from SemCor, D-GPT is developed based on the gpt-3.5-turbo-1106 model which is one of the latest fine-tunable GPT models released by OpenAI¹. More fine-tuning details are provided in Appendix B. During inference,

> ¹https://platform.openai.com/docs/guides/ fine-tuning

Item types	Image-Enhanced SemCor	VWSD-KB	
# of instances	226,036	48,469	
# of ambiguous target words	33,657	24,989	
# of senses	39,201	31,306	
# of images	181,123	111,575	

Table 1: The statistical details of Image-Enhanced Sem-Cor and VWSD-KB.

the original contexts in WSD test sets are concatenated with the lexical definitions generated by D-GPT, to create semantically augmented contexts. These augmented contexts are subsequently fed into the PolCLIP model. The trained SIG and ISG modules can support the PolCLIP model to address Multimodal-WSD even when any modality is missing. The inference procedure of the PolCLIP model can be abstracted in Algorithm 1. 358

359

360

361

362

363

364

365

367

368

369

370

371

372

373

374

375

376

378

379

380

381

382

383

385

387

388

390

391

392

393

394

395

396

397

399

3.4 Training data

The PolCLIP model relies on large-scale aligned image-sense pairs to learn multimodal polysemy knowledge. Thus, we construct a disambiguationoriented image-sense dataset by integrating Sem-Cor (Miller et al., 1993) and VWSD-KB (Yang et al., 2023), to achieve the training objective of unified image-text WSD. SemCor is the most prevalent dataset for training Textual-WSD models. VWSD-KB contains multimodal data such as words, senses, and images. The offline BabelNet v5.2 and BabelPic are used to collect relevant images to the senses in SemCor and VWSD-KB. The detailed construction process is provided in Appendix C.

After construction, all senses in SemCor and VWSD-KB to be aligned with at least one image and at most five images. We filter out the collected images that are pornographic, violent, or invalid and conduct a manual validation to ensure there is no data leakage. The SemCor associated with images is called Image-Enhanced SemCor. Table 1 displays the statistical details of Image-Enhanced SemCor and VWSD-KB, comprising a total of 274,505 English instances (each instance includes a context with at least one ambiguous target word). An example of the disambiguation-oriented imagesense dataset is shown in Figure 1. Given a context "He had seen the Andromeda tree in Japan", there are four candidate senses for the ambiguous target word "Andromeda", and each sense corresponds to five images. The most semantically relevant sense and images to this context are Sense 4 and its five associated images.

346

347

348

351

354

		Textual-WSD				Visual-WSD				
Training Data	Models		SE3	SE7	SE10	SE13	SE15	ALL	5	SE23
			F1-score(%)					HR@1(%)	MRR@10(%)	
	Openai/CLIP-ViT-L/14 (ICML 2021)	53.07	47.19	35.60	37.50	57.18	53.52	49.40	57.45	72.60
Zero shot	Laion/CLIP-ViT-L/14 (NeurIPS 2022)	49.27	43.88	31.80	33.70	53.38	49.96	45.73	56.87	70.28
Zero-snot	Laion/CLIP-ViT-H/14 (NeurIPS 2022)	51.47	46.45	36.52	39.10	58.78	51.92	49.21	60.70	75.68
	UVWSD (ACL 2023)	-	-	-	-	-	-	-	80.50	87.60
	Openai/CLIP-ViT-L/14 (ICML 2021)	70.01	66.48	61.35	71.96	72.35	69.11	69.41	76.67	84.20
	Laion/CLIP-ViT-L/14 (NeurIPS 2022)	71.01	65.96	62.34	72.92	71.44	68.37	69.50	75.38	84.02
	Laion/CLIP-ViT-H/14 (NeurIPS 2022)	71.00	69.63	65.24	74.53	75.36	71.61	71.83	77.04	84.37
Imaga Enhanced	BEM (ACL 2020)	78.29	75.96	66.64	80.71	<u>81.38</u>	81.72	78.53	-	-
Semcor	SACE (ACL 2021)	80.29	<u>78.67</u>	<u>70.57</u>	<u>82.71</u>	80.86	<u>83.73</u>	<u>80.30</u>	-	-
Senicor	Z-Reweighting (ACL 2022)	79.98	77.04	67.72	82.01	79.94	82.81	79.32	-	-
	FCLL (SemEval 2023)	-	-	-	-	-	-	-	<u>80.13</u>	<u>87.41</u>
	PolCLIP _{base}	82.22	79.89	70.56	85.22	82.79	85.66	82.06	79.48	85.00
	PolCLIP _{base} with D-GPT	83.74	81.41	72.09	86.16	84.31	87.18	83.49	82.94	88.55
	Openai/CLIP-ViT-L/14 (ICML 2021)	71.90	68.23	63.21	73.76	74.12	67.80	70.85	75.98	83.93
	Laion/CLIP-ViT-L/14 (NeurIPS 2022)	69.86	67.37	60.46	72.85	73.24	70.00	69.93	77.88	84.68
Image Enhanced	Laion/CLIP-ViT-H/14 (NeurIPS 2022)	73.90	70.37	65.24	75.85	76.24	70.80	73.04	77.46	84.60
Somoor	BEM (ACL 2020)	78.09	75.03	67.65	80.01	78.45	80.09	77.46	-	-
Senicor	SACE (ACL 2021)	<u>81.93</u>	<u>79.71</u>	<u>71.71</u>	<u>84.35</u>	79.22	<u>84.87</u>	<u>81.09</u>	-	-
	Z-Reweighting (ACL 2022)	79.53	77.03	68.92	81.07	<u>82.50</u>	82.09	79.53	-	-
V WSD-KD	FCLL (SemEval 2023)	-	-	-	-	-	-	-	<u>81.37</u>	<u>87.69</u>
	PolCLIP _{large}	82.76	82.39	71.11	85.18	85.29	86.20	82.60	82.28	87.98
	PolCLIP _{large} with D-GPT	84.66	82.43	72.97	83.39	84.91	86.40	83.62	83.59	90.07

Table 2: Comparison with state-of-the-art methods on Multimodal-WSD benchmark test sets. Bold numbers indicate results of the SOTA model, and underlined numbers denote results of the second best model.

4 Experiments and Results

4.1 Datasets

Due to the training objective of unified imagetext WSD, we opt not to use the validation sets from Textual-WSD or Visual-WSD. We allocate 80% of the combined Image-Enhanced Semcor and VWSD-KB datasets as the training set and reserve the remaining 20% as the validation set. XL-WSD (Pasini et al., 2021), an extra-large evaluation framework for Textual-WSD, is employed to evaluate the model performance on Textual-WSD. XL-WSD is widely used since it encompasses six English all-words Textual-WSD benchmark datasets, including SensEval-2 (SE2, (Palmer et al., 2001)), SensEval-3 (SE3, (Snyder and Palmer, 2004)), SemEval-2007 (SE7, (Navigli et al., 2007)), SemEval-2010 (SE10, (Agirre et al., 2010)), SemEval-2013 (SE13, (Navigli et al., 2013)), and SemEval-2015 (SE15, (Moro and Navigli, 2015)). These six benchmark datasets comprise a total of 8,517 English instances for testing. SemEval-2023 (SE23, (Raganato et al., 2023)) is used to assess the model performance on Visual-WSD, as it is currently the most widely used Visual-WSD benchmark containing 463 English instances.

4.2 Settings

Our model is implemented on Pytorch 2.0.1 and 4 RTX 4090 GPUs. Both the text encoder and

image encoder are initialized by CLIP-ViT-L/14 (Radford et al., 2021). All parameters of the text encoder are optimized, while the image encoder is completely frozen. The sense batch size is set to 50, the image batch size is set to 250 and the epoch is set to 20. Following ALBEF (Li et al., 2021), the momentum is set to 0.005. AdamW is applied to optimize model parameters with a learning rate of 1e-04 and weight decay of 0.05. The image resolution is specified as 224×224, and the maximum text length is set to 77. F1-score is used to evaluate the model performance on Textual-WSD. Hit Rate at 1 (HR@1, i.e., accuracy) and Mean Reciprocal Rank at 10 (MRR@10) are used to assess the model performance on Visual-WSD. 428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

4.3 Baseline models

We train PolCLIP_{base} using Image-Enhanced Semcor and PolCLIP_{large} using the combination of Image-Enhanced Semcor and VWSD-KB. In the testing phase, both PolCLIP_{base} and PolCLIP_{large} are integrated with D-GPT for semantic enhancement. Our model is compared with recent stateof-the-art (SOTA) methods including (1) SOTA models in Textual-WSD: BEM (Blevins and Zettlemoyer, 2020), SACE (Wang and Wang, 2021) and Z-Reweighting (Su et al., 2022), (2) SOTA models in Visual-WSD: FCLL (Yang et al., 2023) and UVWSD (Kwon et al., 2023), (3) SOTA models in image-text learning tasks: Openai/CLIP-VIT-L/14

418

419

420

421

422

423

424

425

426

427

400

401

402

403

404

(Radford et al., 2021), Laion/CLIP-VIT-L/14 and
Openai/CLIP-VIT-H/14 (Schuhmann et al., 2022).
More details about baseline models are provided
in Appendix D. For a fair comparison, these baseline models are retrained using Image-Enhanced
Semcor and VWSD-KB.

4.4 Multimodal-WSD results

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

503

505

506

The comparison results between PolCLIP and the baseline models on Multimodal-WSD benchmark test sets are shown in Table 2. The PolCLIP model achieves the state-of-the-art performance. SACE and FCLL are the second best models for Textual-WSD and Visual-WSD, respectively. Without D-GPT, PolCLIPlarge reaches an F1-score of 82.60% on all the Textual-WSD test data, which is 1.51% higher than SACE. It attains an HR@1 of 82.28% and a MRR@10 of 87.98% on Visual-WSD, which are 0.91% and 0.29% higher than FCLL respectively. With D-GPT, the performance of both PolCLIP_{base} and PolCLIP_{large} is enhanced thanks to the semantically augmented contexts with lexical definitions. In this situation, PolCLIP_{large} gains an F1-score of 83.59% on Textual-WSD, which is 2.53% higher than SACE. It attains an HR@1 of 83.59% and a MRR@10 of 90.07% on Visual-WSD, which are 2.22% and 2.38% higher than FCLL respectively.

Specifically, without training, Openai/CLIP-VIT-L/14, Laion/CLIP-VIT-L/14, and Openai/CLIP-VIT-H/14 (collectively called CLIPs) obtain zeroshot F1-scores below 50% on all the Textual-WSD test data, due to the pre-training data and goal of CLIPs do not target WSD. Conversely, UVWSD obtains over 80% zero-shot HR@1 on Visual-WSD. While only using the Image-Enhanced Semcor as the training set, PolCLIPbase outperforms CLIPs, BEM, SACE, and Z-Reweighting on Textual-WSD, even though its performance on Visual-WSD is slightly inferior to FCLL. When the combination of Image-Enhanced Semcor and VWSD-KB is used as the training set, PolCLIPlarge shows further improvement over PolCLIPbase and surpasses all baseline models. Additionally, the effectiveness of the disambiguation-oriented image-sense dataset is proven, with the performance of all the baseline models on Multimodal-WSD being bolstered.

4.5 Ablation study

An ablation study is conducted to reveal the contribution of each module and the results are reported in Table 3. For the two generation-based training

Models	Textual-WSD	Visual-WSD
widdels	ALL F1-score (%)	HR@1(%)
w/o-SIG	74.52 (-8.08)	81.21 (-1.07)
w/o-ISG	81.17 (-1.43)	76.46 (-5.82)
w/o-W2I	82.95 (+0.35)	9.29 (-72.99)
w/o-W2S	19.73 (-62.87)	82.72 (+0.44)

Table 3: Ablation study of $PolCLIP_{large}$ on the Multimodal-WSD benchmark test sets.

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

objectives, the Sense-to-Image Generation module is removed first, which corresponds to \mathcal{L}_{SIG} . In this scenario, the F1-score of PolCLIP_{large} on Textual-WSD decreases by 8.08% and the HR@1 on Visual-WSD drops by 1.07%. This indicates that the implicit image information generated by the SIG module can aid PolCLIP_{large} in acquiring enriched multimodal representations. Secondly, the Image-to-Sense Generation module is removed, which corresponds to \mathcal{L}_{ISG} . The HR@1 of PolCLIP_{large} on Visual-WSD decreases by 5.82% and the F1-score on Textual-WSD drops by 1.43%. This demonstrates that the implicit sense information generated by the ISG module can also facilitate PolCLIP_{large} learning deep polysemy knowledge.

Regarding the two understanding-based training objectives, the alignment process between the anchor focused on ambiguous target words and candidate images is eliminated first, which corresponds to \mathcal{L}_{W2I} . This means that PolCLIP_{large} exclusively trains for Textual-WSD. At this point, the HR@1 of PolCLIPlarge on Visual-WSD is only 9.29%. PolCLIP_{large} can be regarded as a model that makes random selections for Visual-WSD, since it is required to choose one from ten candidate images for each instance. However, the trade-off problem caused by multi-task training objectives allows its F1-score on Textual-WSD has a 0.35% improvement. Secondly, the alignment process between the anchor and candidate senses is eliminated, corresponding to \mathcal{L}_{W2S} . This means that PolCLIP_{large} exclusively trains for Visual-WSD. In this situation, the F1-score of PolCLIPlarge on Textual-WSD is only 19.73% but the HR@1 on Visual-WSD increases by 0.44%. The reasons can be similarly explained.

Overall, the two generation-based modules actually facilitate PolCLIP learning multimodal polysemy knowledge. The two understanding-based alignment processes are the most critical components, since they maximize the similarities between contexts and senses/images in a feature space.



Figure 3: The experimental results of $PolCLIP_{large}$ with different layer numbers for optimizing the image encoder.

An additional experiment, which investegates the generality of D-GPT for Multimodal-WSD, is provided in Appendix E.

4.6 Optimal layer number for optimizing image encoder

To reduce the computational cost, we opt not to optimize all parameters of the image encoder. With all parameters of the text encoder being optimized, the last 0/4/8/12 layers of the image encoder are separately optimized to investigate their impact on the model performance. The results of PolCLIP_{large} with different layer numbers for optimizing the image encoder are displayed in Figure 3. When zero layers are optimized (meaning the image encoder is completely frozen), PolCLIP_{large} has the smallest parameter size and the SOTA results on Multimodal-WSD. Therefore, this model configuration is selected as our best model (as reported in Table 2). It is interesting that when more layers are optimized, the model performance gradually improves in a small way for Textual-WSD, but drops significantly for Visual-WSD. This is contrary to our expectations. Theoretically, optimizing more layers of the image encoder should enhance the model ability to capture image knowledge. We speculate that redundant knowledge, introduced by some noisy image-sense pairs in the training set, increases the model's training burden. Thus, refining this disambiguation-oriented image-sense dataset would be valuable.

4.7 Analysis on model performance for different PoS

Since some words may present different parts of speech (PoS) in contexts, exploring the model performance for ambiguous target words with different PoS is beneficial to reveal the unique advan-

Models	NOUN	VERB	ADJ	ADV
SACE	82.84	74.23	84.77	81.90
DolCLID.	84.23	74.38	88.13	87.62
FOICLIPlarge	(+1.39)	(+0.15)	(+3.36)	(+5.72)

Table 4: The F1-score results of SACE and PolCLIP_{*large*} for ambiguous target words with different parts of speech.

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

tages of the PolCLIP model. The F1-score results of SACE and PolCLIPlarae, trained on the combination of Image-Enhanced SemCor and VWSD-KB, are shown in Table 4. Compared with SACE, PolCLIP_{large} has improvements of 1.39%, 0.15%, 3.36% and 5.72% in NOUN, VERB, ADJ and ADV respectively. One of the challenges in WSD tasks is the difficulty of understanding non-concrete words accurately, which are often adjectives or adverbs. PolCLIP_{large} happens to have a more obvious improvement in adjectives and adverbs. Therefore, we believe that PolCLIP_{large} has favorable adaptability and flexibility for ambiguous target words with different PoS, due to its generation-based advantages. On the one hand, it can supplement tangible senses or images with semantically concrete images or descriptions. On the other hand, it can supplement non-concrete senses or images with semantically abstract images or descriptions. To further intuitively reveal the effectiveness and importance of the two generation-based modules, visualizations of the implicit image-text complementary information generated by the SIG and ISG modules for concrete and non-concrete examples are provided in Appendix F.

5 Conclusion

This paper proposes a unified image-text WSD model PolCLIP, which achieves state-of-the-art performance on Textual-WSD and Visual-WSD benchmark datasets. Extensive experimental results prove the effectiveness of our image-text complementarity strategy. A series of in-depth explorations of the model architecture demonstrate the Sense-to-Image Generation module and the Imageto-Sense Generation module can effectively simulate stable diffusion and image captioning, respectively. The disambiguation-oriented imagesense dataset empirically facilitates WSD models understanding of multimodal polysemy knowledge. This may provide a benchmark for the future Multimodal-WSD task.

574

576

577

581

584

549

550

722

723

724

725

726

727

728

729

730

731

732

733

734

735

Limitation

626

641

642

650

651

657

663

664

670

671 672

675

676

678

Although our method achieves state-of-the-art results on English all-words benchmarks, we do not explore it on multilingual data. The proposed image-text complementarity strategy can support the PolCLIP model to generating multimodal complementary representations for addressing the modality missing issues in unimodal WSD datasets. However, the PolCLIP model still lacks the ability to generate realistic senses and images that can be intuitively validated in terms of their semantics. In future work, we plan to further expand the disambiguation-oriented image-sense dataset to cover more languages. We will also develop a large generic model suitable for Multimodal-WSD.

References

- Eneko Agirre, Oier Lopez de Lacalle, Christiane Fellbaum, Shu-Kai Hsieh, Maurizio Tesconi, Monica Monachini, Piek Vossen, and Roxanne Segers. 2010.
 SemEval-2010 task 17: All-words word sense disambiguation on a specific domain. In *Proceedings of the 5th International Workshop on Semantic Evaluation*, pages 75–80, Uppsala, Sweden. Association for Computational Linguistics.
- Michele Bevilacqua and Roberto Navigli. 2020. Breaking through the 80% glass ceiling: Raising the state of the art in word sense disambiguation by incorporating knowledge graph information. In *Proceedings* of the 58th Annual Meeting of the Association for Computational Linguistics, pages 2854–2864, Online. Association for Computational Linguistics.
- Michele Bevilacqua, Tommaso Pasini, Alessandro Raganato, and Roberto Navigli. 2021. Recent trends in word sense disambiguation: A survey. In Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21, pages 4330–4338. International Joint Conferences on Artificial Intelligence Organization. Survey Track.
- Terra Blevins and Luke Zettlemoyer. 2020. Moving down the long tail of word sense disambiguation with gloss informed bi-encoders. In *Proceedings* of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1006–1017, Online. Association for Computational Linguistics.
- Agostina Calabrese, Michele Bevilacqua, and Roberto Navigli. 2020a. Fatality killed the cat or: BabelPic, a multimodal dataset for non-concrete concepts. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4680– 4686, Online. Association for Computational Linguistics.
- Agostina Calabrese, Michele Bevilacqua, Roberto Navigli, et al. 2020b. Evilbert: Learning task-agnostic

multimodal sense embeddings. In *IJCAI*, pages 481–487. International Joint Conferences on Artificial Intelligence Organization.

- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR.
- Simone Conia and Roberto Navigli. 2021. Framing word sense disambiguation as a multi-label problem for model-agnostic knowledge integration. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3269–3275, Online. Association for Computational Linguistics.
- Slawomir Dadas. 2023. OPI at SemEval-2023 task 1: Image-text embeddings and multimodal information retrieval for visual word sense disambiguation. In Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023), pages 155– 162, Toronto, Canada. Association for Computational Linguistics.
- Spandana Gella, Desmond Elliott, and Frank Keller. 2019. Cross-lingual visual verb sense disambiguation. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 1998– 2004, Minneapolis, Minnesota. Association for Computational Linguistics.
- Spandana Gella, Mirella Lapata, and Frank Keller. 2016. Unsupervised visual sense disambiguation for verbs using multimodal embeddings. In *Proceedings of the* 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 182–192, San Diego, California. Association for Computational Linguistics.
- Omid Ghahroodi, Seyed Arshan Dalili, Sahel Mesforoush, and Ehsaneddin Asgari. 2023. SUT at SemEval-2023 task 1: Prompt generation for visual word sense disambiguation. In *Proceedings of the 17th International Workshop on Semantic Evaluation* (*SemEval-2023*), pages 2160–2163, Toronto, Canada. Association for Computational Linguistics.
- Raia Hadsell, Sumit Chopra, and Yann LeCun. 2006. Dimensionality reduction by learning an invariant mapping. In 2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06), volume 2, pages 1735–1742. IEEE.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840– 6851.
- Luyao Huang, Chi Sun, Xipeng Qiu, and Xuanjing Huang. 2019. GlossBERT: BERT for word sense disambiguation with gloss knowledge. In *Proceedings*

736

- 761 764 770

- 779
- 781
- 783
- 784

- 790
- 792

of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3509-3514, Hong Kong, China. Association for Computational Linguistics.

- Anastasia Kritharoula, Maria Lymperaiou, and Giorgos Stamou. 2023. Large language models and multimodal retrieval for visual word sense disambiguation. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 13053-13077, Singapore. Association for Computational Linguistics.
- Sunjae Kwon, Rishabh Garodia, Minhwa Lee, Zhichao Yang, and Hong Yu. 2023. Vision meets definitions: Unsupervised visual word sense disambiguation incorporating gloss information. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1583-1598, Toronto, Canada. Association for Computational Linguistics.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. arXiv preprint arXiv:2301.12597.
- Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. 2021. Align before fuse: Vision and language representation learning with momentum distillation. Advances in neural information processing systems, 34:9694-9705.
- Marco Maru, Federico Scozzafava, Federico Martelli, SyntagNet: Chaland Roberto Navigli. 2019. lenging supervised word sense disambiguation with lexical-semantic combinations. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3534–3540, Hong Kong, China. Association for Computational Linguistics.
- George A Miller, Richard Beckwith, Christiane Fellbaum, Derek Gross, and Katherine J Miller. 1990. Introduction to wordnet: An on-line lexical database. International journal of lexicography, 3(4):235–244.
- George A. Miller, Claudia Leacock, Randee Tengi, and Ross T. Bunker. 1993. A semantic concordance. In Human Language Technology: Proceedings of a Workshop Held at Plainsboro, New Jersey, March 21-24, 1993.
- Andrea Moro and Roberto Navigli. 2015. SemEval-2015 task 13: Multilingual all-words sense disambiguation and entity linking. In Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015), pages 288–297, Denver, Colorado. Association for Computational Linguistics.
- Roberto Navigli. 2009. Word sense disambiguation: A survey. ACM computing surveys (CSUR), 41(2):1-69.

Roberto Navigli, Michele Bevilacqua, Simone Conia, Dario Montagnini, and Francesco Cecconi. 2021. Ten years of babelnet: A survey. In IJCAI, pages 4559-4567.

793

794

795

796

797

798

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

- Roberto Navigli, David Jurgens, and Daniele Vannella. 2013. SemEval-2013 task 12: Multilingual word sense disambiguation. In Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013), pages 222-231, Atlanta, Georgia, USA. Association for Computational Linguistics.
- Roberto Navigli, Kenneth C. Litkowski, and Orin Hargraves. 2007. SemEval-2007 task 07: Coarse-grained English all-words task. In Proceedings of the Fourth International Workshop on Semantic Evaluations (SemEval-2007), pages 30-35, Prague, Czech Republic. Association for Computational Linguistics.
- Roberto Navigli and Simone Paolo Ponzetto. 2010. BabelNet: Building a very large multilingual semantic network. In Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, pages 216-225, Uppsala, Sweden. Association for Computational Linguistics.
- Martha Palmer, Christiane Fellbaum, Scott Cotton, Lauren Delfs, and Hoa Trang Dang. 2001. English tasks: All-words and verb lexical sample. In Proceedings of SENSEVAL-2 Second International Workshop on Evaluating Word Sense Disambiguation Systems, pages 21-24, Toulouse, France. Association for Computational Linguistics.
- Tommaso Pasini, Alessandro Raganato, and Roberto Navigli. 2021. Xl-wsd: An extra-large and crosslingual evaluation framework for word sense disambiguation. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pages 13648-13656.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In International conference on machine learning, pages 8748-8763. PMLR.
- Alessandro Raganato, Iacer Calixto, Asahi Ushio, Jose Camacho-Collados, and Mohammad Taher Pilehvar. 2023. SemEval-2023 task 1: Visual word sense disambiguation. In Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023), pages 2227–2234, Toronto, Canada. Association for Computational Linguistics.
- Alessandro Raganato, Jose Camacho-Collados, and Roberto Navigli. 2017. Word sense disambiguation: A unified evaluation framework and empirical comparison. In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, pages 99-110, Valencia, Spain. Association for Computational Linguistics.

Alessandro Raganato, Yves Scherrer, and Jörg Tiedemann. 2019. The MuCoW test suite at WMT 2019: Automatically harvested multilingual contrastive word sense disambiguation test sets for machine translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 470–480, Florence, Italy. Association for Computational Linguistics.

851

852

854

865

870

871

874

876

878

879

882

899

900

901

902

903

904 905

906

- Rita Ramos, Desmond Elliott, and Bruno Martins. 2023. Retrieval-augmented image captioning. In Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, pages 3666–3681, Dubrovnik, Croatia. Association for Computational Linguistics.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. Laion-5b: An open large-scale dataset for training next generation imagetext models. Advances in Neural Information Processing Systems, 35:25278–25294.
- Federico Scozzafava, Marco Maru, Fabrizio Brignone, Giovanni Torrisi, and Roberto Navigli. 2020. Personalized PageRank with syntagmatic information for multilingual word sense disambiguation. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, pages 37–46, Online. Association for Computational Linguistics.
 - Benjamin Snyder and Martha Palmer. 2004. The English all-words task. In Proceedings of SENSEVAL-3, the Third International Workshop on the Evaluation of Systems for the Semantic Analysis of Text, pages 41–43, Barcelona, Spain. Association for Computational Linguistics.
 - Ying Su, Hongming Zhang, Yangqiu Song, and Tong Zhang. 2022. Rare and zero-shot word sense disambiguation using Z-Reweighting. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4713–4723, Dublin, Ireland. Association for Computational Linguistics.
- Ming Wang and Yinglin Wang. 2020. A synset relationenhanced framework with a try-again mechanism for word sense disambiguation. In *Proceedings of the* 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6229–6240, Online. Association for Computational Linguistics.
- Ming Wang and Yinglin Wang. 2021. Word sense disambiguation: Towards interactive context exploitation from both word and sense perspectives. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 5218–5229, Online. Association for Computational Linguistics.

- Qihao Yang, Yong Li, Xuelin Wang, Shunhao Li, and 907 Tianyong Hao. 2023. TAM of SCNU at SemEval-908 2023 task 1: FCLL: A fine-grained contrastive 909 language-image learning model for cross-language 910 visual word sense disambiguation. In Proceedings of 911 the 17th International Workshop on Semantic Eval-912 uation (SemEval-2023), pages 506-511, Toronto, 913 Canada. Association for Computational Linguistics. 914
- Xudong Zhang, Tiange Zhen, Jing Zhang, Yujin Wang, and Song Liu. 2023. SRCB at SemEval-2023 task
 1: Prompt based and cross-modal retrieval enhanced visual word sense disambiguation. In Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023), pages 439–446, Toronto, Canada. Association for Computational Linguistics.

A Pseudocode of training PolCLIP

Our PolCLIP model, built upon the CLIP (Radford et al., 2021) architecture, leverages an image-text complementarity strategy to excel in processing multimodal polysemy. It is trained on aligned image-sense pairs, enhancing its generalization abilities on WSD. The PolCLIP model optimizes four distinct objectives through four contrastive losses, encompassing two generation-based modules and two understanding-based alignment processes. The overall training procedure of the Pol-CLIP model can be abstracted in Algorithm 2.

ł	Algorithm 2: Pseudocode of Training PolCLIP				
	data : a context c with an ambiguous target word w_{target} ; the candidate senses S and the candidate images I ;				
1	while	e c, w_{target} , S, I do			
2	<u> </u>	$e_c \leftarrow \text{Text-Encoder}(c)$; # the complete context representations			
3		$e^{w_{\text{target}}} \leftarrow e_c$; # the anchor vector based on w_{target}			
4		# the original unimodal sense/image representations			
5		$v_{\text{sense}} \leftarrow \text{Text-Encoder}(S);$			
6		$v_{\text{image}} \leftarrow \text{Image-Encoder}(I);$			
7		# the generated implicit image/sense representations			
8		$v_{\text{image}}^{\text{sense-gen}} \leftarrow \text{SIG}(v_{\text{sense}});$			
9		$v_{\text{sense}}^{\text{image-gen}} \leftarrow \text{ISG}(v_{\text{image}});$			
10		# the semantically enriched multimodal representations			
11		$v_{\text{multi}}^{\text{sense-gui}} \leftarrow \text{Fusion}(v_{\text{sense}}, v_{\text{image}}^{\text{sense-gen}});$			
12		$v_{\text{multi}}^{\text{image-gui}} \leftarrow \text{Fusion}(v_{\text{image}}, v_{\text{sense}}^{\text{image-gen}});$			
13		# the anchor-to-sense and anchor-to-image similarities			
14		$sim(W2S) \leftarrow e^{w_{\text{target}}} \cdot (v_{\text{multi}}^{\text{sense-gui}})^T;$			
15		$sim(W2I) \leftarrow e^{w_{\text{target}}} \cdot (v_{\text{multi}}^{\text{tange-gul}})^T;$			
16					
17		# the pseudo sense/image representations			
18		$v_{\text{sense}}^{\text{pseudo}} \leftarrow \text{Text-Momentum-Encoder}(S);$			
19		$v_{\text{image}}^{\text{pseudo}} \leftarrow \text{Image-Momentum-Encoder}(I);$			
20		# the SIG and IGS similarities			
21		$sim(SIG) \leftarrow v_{\text{image}}^{\text{sense-gen}} \cdot (v_{\text{image}}^{\text{pseudo}})^T;$			
22		$sim(ISG) \leftarrow v_{sense}^{image-gen} \cdot (v_{sense}^{pseudo})^T;$			
23		# the two generation-based loss			
24		$L_{SIG} \leftarrow \text{CrossEntropyLoss}(sim(SIG), labels(SIG));$			
25		$L_{ISG} \leftarrow \text{CrossEntropyLoss}(sim(ISG), labels(ISG));$			
26		# the two understanding-based loss			
27		$L_{W2S} \leftarrow \text{CrossEntropyLoss}(sim(W2S), labels(W2S));$			
28		$L_{W2I} \leftarrow \text{CrossEntropyLoss}(sim(W2I), labels(W2I));$			
29	end				

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931



Figure 4: The format of message-style data for finetuning D-GPT. The red GT denotes the ground-truth sense.

B Fine-tuning D-GPT

934

935

937

938

939

940

943

944

947

951

953

955

956

961

962

963

965

966

967

969

970

A fine-tuned disambiguation-oriented GPT-3.5 (D-GPT) is developed to generate lexical definitions for ambiguous target words during the inference procedure. The gpt-3.5-turbo-1106 model is chosen as the backbone model, since it is one of the latest fine-tunable GPT models released by OpenAI and outperforms several LLMs with smaller parameter sizes in terms of inference capabilities. Constrained by fine-tuning costs, we randomly collect 50,000 data from SemCor to serve as the fine-tuning dataset. Each data consists of a context with an ambiguous target word and its ground-truth sense. Following OpenAI's guidelines, the format of message-style data for fine-tuning D-GPT is shown in Figure 4.

C Construction of the disambiguation-oriented image-sense datasets

The disambiguation-oriented image-sense dataset is constructed based on SemCor (Miller et al., 1993) and VWSD-KB (Yang et al., 2023) datasets. Specially, the offline version of BabelNet² v5.2is employed to collect a list of relevant images based on each sense in these two datasets. If there are more than five available images in the list, the top five are selected; otherwise, all images are retained. However, a minority of the senses fail to associate with any obtain any image through BabelNet, as they are typically non-concrete, like expressing sadness. Thus, BabelPic (Calabrese et al., 2020a), an image-text dataset for non-concrete concepts, is utilized to find images for a part of nonconcrete senses based on babel-ids. Furthermore, those senses that we have collected relevant images are set as an internal knowledge base. For each sense s_i that is not included in either Babel-



Figure 5: The evaluation results of SACE, FCLL, and PolCLIP_{large} on the benchmark test sets after integrating D-GPT.

Net or BabelPic, RoBERTa³ is used to identify the three senses most semantically similar to s_i within this internal knowledge base, based on text similarity. The first image from each of these three most similar senses is aggregated as the set of images corresponding to s_i . The entire construction process enables all senses in SemCor and VWSD-KB to be aligned with at least one image and at most five images.

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1005

1006

D Baselines

The details of baselines are as follows:

SOTA models in Textual-WSD: BEM (Blevins and Zettlemoyer, 2020) adopted two text encoders and focuses on the representations of ambiguous target words rather than the complete context representations. SACE (Wang and Wang, 2021) employed an interactive context exploitation method and selects similar sentences from the same document to enhance context representations. Z-Reweighting (Su et al., 2022) utilized a strategy for adjusting training on imbalanced datasets at the word level. These three models obtain outstanding performance on Textual-WSD benchmarks when training exclusively on SemCor.

SOTA models in Visual-WSD: FCLL (Yang et al., 2023) employed a fine-grained image-text contrastive learning mechanism and benefited from VWSD-KB. It won first place in SemEval-2023 Task 1. UVWSD (Kwon et al., 2023) did not necessitate training but achieved remarkable performance by employing Bayesian inference to incorporate sense definitions.

SOTA models in image-text learning tasks: Openai/CLIP-VIT-L/14 (Radford et al., 2021), Laion/CLIP-VIT-L/14 and Openai/CLIP-VIT-H/14 (Schuhmann et al., 2022) all employ a dual-stream architecture to learn image-text knowledge, simi-

²https://babelnet.org/

³https://huggingface.co/FacebookAI/ roberta-large



Figure 6: Visualizations of the implicit image-text complementary information generated by the SIG and ISG modules for concrete and non-concrete examples.

lar to our PolCLIP model. The first model is pretrained on over 400 million image-text pairs, and the latter two are pre-trained on the English subset of LAION-5B (Schuhmann et al., 2022), a large publicly available image-text dataset.

E Generality of D-GPT

1008

1011

1013

An additional experiment is conducted to explore 1014 the generality of D-GPT for Multimodal-WSD. Specifically, D-GPT is integrated with SACE and 1016 FCLL, and its impact on the performance of these two models and PolCLIP_{large} is illustrated in Figure 5. D-GPT indeed enhances the performance of 1019 these three WSD models. This indicates that us-1020 ing lexical definitions generated by D-GPT to create semantically augmented contexts is a general-1023 purpose and convenient pipeline for Multimodal-WSD. It can be applied to various WSD models. 1024 Furthermore, compared to the evaluation results 1025 without D-GOT, PolCLIP_{large} shows an improvement of 1.02% F1-score and 1.31% HR@1. These 1027

are respectively higher than the 0.44% F1-score1028increase of SACE on Textual-WSD and the 0.64%1029HR@1 increase of FCLL on Visual-WSD. This1030also leads us to believe that PolCLIP could gain1031more when dealing with contexts that are semanti-1032cally more accurate, thanks to the image-text complementarity strategy.1034

F Effectiveness of the SIG and ISG modules

In order to further intuitively reveal the effective-1037 ness and importance of the two generation-based modules (i.e., SIG and ISG), the generated implicit 1039 image-text complementary information is visual-1040 ized. Two groups of concrete and non-concrete examples are collected from the test set. Each 1042 group of examples contains a sense and an im-1043 age. They are fed into the trained SIG and ISG 1044 modules respectively, and then a sense-generated 1045 image vector and an image-generated sense vector are output. All senses and images in the training 1047

set are transformed into vectors by text and image 1048 encoders, serving as two separate candidate pools. 1049 By calculating vector similarity, the top-5 most 1050 similar images can be identified from the image 1051 candidate pool based on the sense-generated image 1052 vector. Also, the top-5 most similar senses can be 1053 identified from the sense candidate pool based on 1054 the image-generated sense vector. Visualizations 1055 of these two groups of concrete and non-concrete 1056 examples are shown in Figure 6. For the concrete 1057 example, the top-3 images are semantically con-1058 sistent with Sense 1. Even if the last two images 1059 do not depict the shape of a bell, they are related 1060 to music or sound, which is one of the functions 1061 of bells. Based on Image 1 related to milk, the 1062 retrieved five senses are all semantically correct. 1063 For the non-concrete example, the top-3 images are 1064 semantically relevant to Sense 2 related to beauty. 1065 Even based on Image 2, which primarily shows a 1066 man's face expressing pleasure, the top-3 senses 1067 accurately capture concepts of relaxation, laughter, 1068 and beard. 1069

In summary, thanks to the two generation-based training objectives, the SIG and ISG modules are reliable in effectively imitating stable diffusion and image captioning.

1070

1071

1072