

SPOT: Structured Prompting with Object-centric Tokens for open-world scene graphs

Anonymous CVPR submission

Paper ID 30

Abstract

001 *Scene graphs provide a compact and structured representation of visual scenes by capturing objects and their relationships, making them valuable for downstream tasks in vision-language reasoning and robotics. While early work focused on closed-vocabulary settings, newer efforts have shifted toward open-world scene graph generation (SGG) to better handle diverse real-world scenarios. Recent works explore leveraging VLMs and LLMs in open-world settings for their broad, open-vocabulary knowledge. However, existing approaches often rely on proprietary models like GPT-4o and are limited by the unstructured output behavior and weak spatial and object-level reasoning capabilities of pre-trained models. We introduce SPOT, a structured prompting framework that augments open-source VLMs with spatial reasoning abilities for scene graph generation with minimal training. By combining object-centric visual features with the model’s knowledge priors, SPOT achieves competitive or superior relation prediction compared to large proprietary models. Additionally, SPOT demonstrates strong cross-domain generalization. Our approach is built upon open-source models, offering a scalable and accessible framework for harnessing VLMs for SGG.*

023 1. Introduction

024 Understanding complex scenes has long been a core challenge in computer vision. Among various representations, scene graphs have emerged as a powerful paradigm to capture both semantic and spatial relationships between objects in a structured and interpretable format. By abstracting a scene into a graph with objects as nodes and their relationships as edges, scene graphs align closely with human perception and reasoning, providing a symbolic abstraction useful for downstream tasks [15]. In 2D settings, scene graphs enhance general-purpose vision-language models (VLMs) by enabling explicit relational understanding, benefiting applications such as visual question answering [13, 27, 31]. In robotics, scene graphs [1] have gained prominence as a compact and expressive modality for high-level percep-

tion, planning, and interaction within physical environments [3, 40, 46]. They offer a bridge between raw sensory data and the structured world, supporting effective decision-making and spatial reasoning. Consequently, Scene Graph Generation (SGG), the task of automatically constructing such graphs from sensory inputs like images, has become a foundational problem with growing attention in both vision and robotics communities. 038 039 040 041 042 043 044 045

The development of scene graph generation began with focus on improving prediction accuracy within domain-specific, closed-vocabulary datasets [26, 36, 41, 42]. While these efforts also addressed key challenges like the long-tail problem [4, 22], the fixed and predefined categories limit real-world applicability. To address this limitation, the field has shifted towards open-vocabulary scene graph generation [6, 10, 11, 19, 24, 47, 49] by leveraging open-vocabulary detectors [17, 28] to incorporate broad semantic knowledge into the pipeline. However, they are still constrained by predefined known relations or object categories when generalizing to unseen scenarios. Most recent works propose to leverage the VLMs, to achieve a more comprehensive open-world setting [5, 24]. Despite these advancements, these methods often depend on closed-source proprietary models, unconstrained prompting text, and insufficient visual reasoning, frequently producing vague or physically implausible relation predictions. 046 047 048 049 050 051 052 053 054 055 056 057 058 059 060 061 062 063

In this work, we introduce **Structured Prompting with Object-centric Tokens (SPOT)** for open-world scene graph generation. Our approach introduces three key components to address existing challenges of using VLMs for scene graph generation. First, we design a template-based structured prompt (in contrast to the free-form prompts of prior works [10, 24]) that more precisely guides the model to produce comprehensive scene graphs both out-of-the-box and after refinement through finetuning. Next, we encourage the model to consider the visual scene layout through integration of an object-centric visual feature when predicting relations. This additional signal improves upon the VLM’s standard processing of the image, increasing spatial alignment and relation accuracy. Finally, to enable open-world prediction 064 065 066 067 068 069 070 071 072 073 074 075 076 077

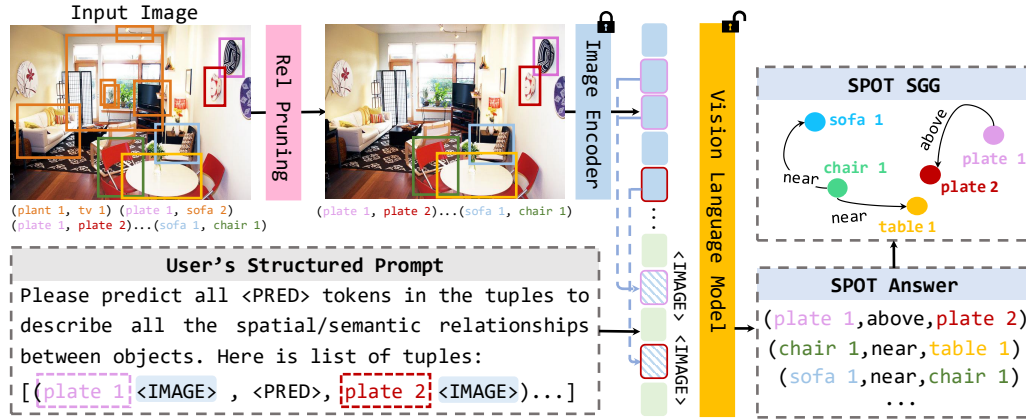


Figure 1. **SPOT framework.** Given an RGB image and object detections, our framework prunes spatially implausible relation pairs and constructs a structured prompt for scene graph generation. Our model further extracts object-centric visual feature embeddings by averaging over the corresponding image feature patches inside the bounding box areas and injecting these into the prompt.

078 with an external object detector and no pre-defined vocabulary, we propose leveraging spatially aware pruning and integrating flexibility during evaluation to minimize penalties for semantically similar predictions. This approach goes beyond the protocol which exhaustively constructs a fully-connected graph over all object pairs, which proves computationally expensive and suffers from redundant and irrelevant relation predictions.

086 Overall, our framework achieves superior performance on both in-domain (e.g., Visual Genome [20]) and cross-domain (e.g., PSG [43], 3DSSG [38]) evaluations compared to preceding works. The quantitative evaluation demonstrates the superiority, versatility, and broad applicability of our model in the physical environment.

092 2. Related works

093 **Scene Graph Generation.** SGG aims to represent visual scenes as a compact graph, identifying objects as nodes and their interrelations as edges. Prior works have tackled SGG by leveraging graph convolutional networks [42], end-to-end DETR-style architectures [23], or novel strategies [18, 34, 41]. Further works target mitigating the long-tail problem through unbiased learning [4, 7, 22, 37], aggregating more diverse visual concept contexts [36, 48, 50], and reducing annotation cost by harnessing language-captions and leveraging weak-supervised learning [25, 45, 53, 54]. However, these approaches focus on a closed-vocabulary setting, where both object and relationship classes are limited to a predefined set from datasets.

106 **Open-Vocabulary Scene Graph Generation.** Recent works have focused on extending SGG from closed-vocabulary to broader open-vocabulary settings. Most approaches approach generalizing to open-world settings by targeting either generalization to new object pairings [11, 19, 47, 49] or predicting unseen predicates between seen object enti-

ties [19, 49]. However, these methods are still constrained by either predefined known relations or object categories when generalizing to unseen scenarios. Another line of work jointly addresses generalization to both new objects and predicates [6, 52], but they still train on closed-vocabulary datasets, limiting true generalization.

118 3. Method

119 We introduce SPOT, a structured prompting framework that augments open-source VLMs with spatial reasoning abilities for open-world scene graph generation with minimal training. Our method, illustrated in Fig. 1, addresses key limitations that arise while leveraging VLMs for this task and builds robust open-world relation prediction through three key components.

126 3.1. Structured Prompt over Selected Relations

127 Our first contribution is the design of a structured prompt to query the VLM to produce output relations over only the most likely relevant object-relation-object triplets. Often many of these objects themselves may be of little relevance for understanding a scene and moreover the relationships between unrelated and spatially disjoint objects may be irrelevant for solving a downstream task.

134 **Relation Proposal.** To mitigate these issues, we decouple the process of relation selection from relation prediction. Specifically, we apply a rule-based filtering mechanism to identify object pairs that are likely to interact in a meaningful way based on spatial proximity and visibility. To construct a comprehensive scene graph, one could consider all pairwise combinations of detected objects as candidate relation triplets. We then apply a spatial filtering strategy based on object distance, scaled by object sizes. This filtering step improves the efficiency and quality of scene graph generation.

Table 1. **Quantitative comparison of free-form prompting vs. our structured prompt template.** We compare the effect of using *Free-form* versus *Structured* prompts across four VLMs: LLaVA-OV-7B, Qwen2.5-VL-7B-Instruct, InternVL3-8B, and GPT-4o. Results are reported on the VG and PSG datasets, using Recall@50/100 and Recall with Similarity@50/100 as evaluation metrics.

| Model Backbone | Prompt Type | VG (1K) | | PSG (1K) | |
|-------------------|-------------|--------------------|---------------------|--------------------|--------------------|
| | | Structured Prompt | R@50/100↑ | R w/ Sim@50/100↑ | R@50/100↑ |
| LLaVA-OV-7B [21] | – | 0.145 / 0.145 | 2.73 / 2.73 | 1.72 / 1.72 | 8.56 / 8.56 |
| | ✓ | 3.17 / 3.44 | 17.2 / 19.13 | 2.47 / 2.75 | 23.4 / 25.6 |
| Qwen2.5-VL-7B [2] | – | 1.40 / 1.40 | 7.88 / 7.88 | 1.17 / 1.17 | 6.97 / 6.97 |
| | ✓ | 3.71 / 4.23 | 18.6 / 21.1 | 3.37 / 3.71 | 21.6 / 24.0 |
| InternVL3-8B [39] | – | 6.01 / 6.01 | 13.0 / 13.0 | 5.56 / 5.56 | 12.3 / 12.3 |
| | ✓ | 8.65 / 9.45 | 23.6 / 26.4 | 11.6 / 12.8 | 27.8 / 30.8 |
| GPT-4o [29] | – | 4.55 / 4.55 | 12.1 / 12.1 | 4.70 / 4.70 | 17.7 / 17.7 |
| | ✓ | 13.5 / 15.2 | 34.8 / 39.4 | 7.74 / 8.52 | 29.7 / 34.2 |

Table 2. **Exploring the value of language and visual priors.** When finetuning using our structured prompt, we evaluate the model’s reliance on language (object names) and visual priors.

| Input to VLM | | VG | | PSG | |
|--------------|--------------|-------------|-------------|--------------|--|
| vision | object names | Acc.↑ | Acc.↑ | Acc. w/ Sim↑ | |
| ✓ | – | 40.5 | 20.4 | 39.4 | |
| – | ✓ | 44.5 | 19.2 | 54.5 | |
| ✓ | ✓ | 45.7 | 22.3 | 57.9 | |

Structured Prompt. Prior methods for scene graph generation with VLMs [10, 24, 32] prompt the VLMs in a free-form manner; i.e., lists a set of objects and ask for the model to directly generate a list of relation triplets (o_i, r_{ij}, o_j) between all objects from an input image. Although straightforward, this free-form prompting often results in low coverage of relation pairs in the output scene graph, with the model omitting important relation pairs while including spatially irrelevant ones. We instead propose to prompt the VLM with a fill-in-the-blank format to predict the relation predicates for only the pairs filtered above. Using this structured prompt template effectively constrains the model’s output space, resolving the issues of incomplete relation pairs and degenerate scene graphs, and achieving a higher recall score with more complete graphs.

To demonstrate the effectiveness of the proposed template, we evaluate four different pretrained VLMs (LLaVA-OV, Qwen2.5-VL-7B-Instruct, InternVL3-8B, and GPT-4o) without finetuning on the VG [20] and PSG dataset [43] with free-form [10] vs. our proposed prompt template in Tab. 1. Our structured prompt uniformly brings significant boosts to the base models, even without any finetuning.

3.2. Object-centric Visual Feature Embedding.

Language models (and VLMs by extension) are known to possess powerful priors about common world knowledge. So, it naturally follows that when asking a VLM to predict an object-object relationship, the output prediction will be heavily biased by the prior relations observed between those objects. For example, even without seeing a visual input, a VLM is biased to predict that the relationship between “cup” and “table” is “on” to indicate that the cup rests on the

Table 3. **Closed-set SGG (PredCLS) results on VG150.**

| Type | SGG model | PredCLS | | |
|--------|-------------|-------------|-------------|-------------|
| | | R@50 | R@100 | Acc. |
| No-VLM | VCTree [35] | 50.1 | 52.5 | – |
| | GCA [18] | 51.2 | 53.4 | – |
| | EBM [34] | 52.8 | 54.9 | – |
| | SVRP [12] | 54.4 | 56.4 | – |
| | OvSGTR [6] | 60.5 | 61.9 | 60.8 |
| VLM | PGSG [24] | 33.8 | 40.2 | 49.2 |
| | SPOT | 58.2 | 63.4 | 70.3 |

Table 4. **Closed-set SGG SGGDet results on the VG150 test set.**

| Type | SGG model | SGGDet | |
|--------|------------|--------------------|-------------|
| | | R@50/100 | mR@50/100 |
| No-VLM | SGTR [23] | 24.6 / 28.4 | 12.0 / 15.2 |
| | ISG [16] | 29.7 / 32.1 | 8.0 / 8.8 |
| | RelTR [8] | 21.2 / 27.5 | 6.8 / 10.8 |
| | VS3 [52] | 34.5 / 39.2 | – |
| | SVRP [12] | 31.8 / 35.8 | 10.5 / 12.8 |
| | OvSGTR [6] | 36.4 / 42.4 | 7.2 / 8.8 |
| VLM | PGSG [24] | 20.3 / 23.6 | 10.5 / 12.7 |
| | SPOT | 38.1 / 44.8 | 10.3 / 12.7 |

table. While this textual bias contributes positively to overall relation prediction performance, it diminishes visual grounding. This limitation becomes critical in scenarios where spatial ground truths deviate from common expectations, for example, a cup actually being “under” a table.

To address this issue, we improve spatial grounding by embedding implicit coordinate-aware visual features into the prompt. While the standard VLM input may be a set of encoded image tokens together with text tokens from the encoded prompt, we additionally add an object-specific visual token after each text object reference. Rather than relying solely on object names, we extract local visual features for each object by mapping its 2D bounding box locations back onto the spatial feature map of the pretrained vision encoder SigLIP [51], as illustrated in Fig. 1.

Specifically, given visual encoder features $F \in \mathbb{R}^{H \times W \times D}$ containing positional information, we identify the set of patches $P_o \subset F$ that fall within the object’s bounding box, and compute the object embedding as $f_o =$

Table 5. Cross-domain SGG results on PSG and 3DSSG datasets. We evaluate several SGG models under the PredCLS setting.

| SGG model | PSG | | | 3DSSG Per Frame | | |
|------------------------------|-------------|-------------|--------------------|-----------------|-------------|--------------------|
| | Acc. | Acc w/ sim | PredCLS R@50/100 | Acc. | Acc w/ sim | PredCLS R@50/100 |
| Llava-OV-7B (Free-form) [21] | - | - | 1.7 / 1.7 | - | - | 5.4 / 5.4 |
| GPT-4o (Free-form) [29] | - | - | 4.7 / 4.7 | - | - | 18.5 / 18.5 |
| GPT-4o (Our template) [29] | 16.6 | 52.5 | 15.8 / 17.3 | 47.8 | 75.1 | 35.2 / 36.1 |
| OvSGTR [6] | 6.1 | 39.6 | 5.3 / 5.8 | 2.91 | 49.7 | 1.23 / 1.67 |
| SPOT | 23.6 | 59.5 | 15.9 / 18.0 | 39.1 | 71.2 | 26.9 / 29.0 |

195 $\frac{1}{|P_o|} \sum_{p \in P_o} p$. This average-pooled patch feature f_o captures
 196 both localized visual appearance and spatial positioning. We
 197 then concatenate f_o with the object’s text label within the in-
 198 put prompt, thereby encouraging the model to attend to both
 199 semantic and spatial positioning cues. By providing a more
 200 direct binding between vision and language for each object,
 201 we steer the model away from relying on texture priors alone
 202 and toward leveraging the visual evidence. As demonstrated
 203 in Tab. 2, this coordinate-aware visual enhancement leads to
 204 an improvement in relation prediction accuracy.

205 3.3. Frame-wise Scene Graph Generation at Infer- 206 ence Time

207 During training, we assume access to ground-truth relation
 208 triplets and corresponding 2D object coordinates, which
 209 enables direct supervision of the relation prediction module.
 210 However, at inference time, ground-truth annotations are
 211 unavailable, and the model must generate both object and
 212 relation proposals directly from raw images.

213 We adopt GroundingDINO [28] as our open-vocabulary
 214 object detector, and it outputs 2D object bounding boxes
 215 and corresponding category labels for each detected object.
 216 To suppress detection noise from repeated detections with
 217 similar object names, we apply standard and cross-category
 218 NMS on the results. We then leverage our filtering step
 219 to reduce the set of all possible object pairings into a smaller
 220 set of plausible object pairings for which the VLM will
 221 predict relations. This filtering step is essential for reducing
 222 computation overhead.

223 4. Experiments

224 4.1. In-Domain Evaluation

225 Since SPOT is trained on a 2D image dataset, it can be di-
 226 rectly applied to standard 2D scene graph generation tasks.
 227 We evaluate its performance on the Visual Genome bench-
 228 mark using both PredCLS and SGDet settings with prior
 229 methods in Tab. 3 and Tab. 4. To ensure a fair comparison,
 230 we adopt the object detector [28] that has been trained on
 231 Visual Genome, eliminating the potential distribution gap.
 232 In the PredCLS setting, ground-truth object boxes and class
 233 labels are used to isolate relation prediction performance. In
 234 the SGDet setting, the model operates end-to-end with de-
 235 tected objects. As shown in Tab 3 & 4, SPOT demonstrates

strong in-domain performance. It outperforms prior mod-
 els, particularly other VLM-based approaches like PGSG,
 indicating superior relation prediction capabilities, as evi-
 denced by PredCLS R@100 and accuracy. Overall, SPOT
 achieves robust relation prediction accuracy, highlighting
 the effectiveness of our object-centric structured prompting
 design.

4.2. Cross-Domain Evaluation

243 We further evaluate the generalization ability of SPOT across
 244 datasets in Tab. 5 compared with previous methods and GPT-
 245 4o. We consider two benchmarks: PSG [43] and the single-
 246 frame (2D) version of 3DSSG [38]. Importantly, no data
 247 from these datasets is seen during our model’s training. PSG
 248 shares visual and categorical similarity with VG150, while
 249 3DSSG poses a more significant domain gap, targeting in-
 250 door scenes and including unseen relations such as "bigger
 251 than." In order to handle vocabulary mismatches in cross-
 252 domain evaluation, we adopt 2 rules for determining if rela-
 253 tions are equivalent: we consider predicted spatial relations
 254 correct if they contain the same spatial words as the ground
 255 truth (for instance, "standing on" and "on" would refer to
 256 the same spatial relation), and we use an LLM as a judge
 257 to calculate Acc. w/ sim to credit relations that are semanti-
 258 cally equivalent to the ground truth. Results show that SPOT
 259 maintains strong generalization and outperforms prior mod-
 260 els, including GPT-4o with/without our proposed structured.
 261 This result validates that our framework effectively adapts
 262 the VLM for improved scene graph relation reasoning, while
 263 leveraging the VLM’s general knowledge from large-scale
 264 pretraining for generalization.
 265

266 5. Conclusion

267 In this work, we presented SPOT, a structured prompting
 268 framework designed to enhance open-world scene graph
 269 generation with vision-language models. By combining
 270 template-based prompts, object-centric visual embeddings,
 271 and relation pruning, our approach addresses the limitations
 272 of free-form prompting and improves both spatial grounding
 273 and semantic coherence. Extensive experiments demonstrate
 274 that SPOT consistently outperforms prior methods across
 275 in-domain and cross-domain benchmarks, achieving com-
 276 petitive or superior results compared to large proprietary
 277 systems such as GPT-4o.

278

References

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

- [1] Iro Armeni, Zhi-Yang He, JunYoung Gwak, Amir R Zamir, Martin Fischer, Jitendra Malik, and Silvio Savarese. 3d scene graph: A structure for unified semantics, 3d space, and camera. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5664–5673, 2019. 1
- [2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. 3
- [3] Yun Chang, Luca Ballotta, and Luca Carlone. D-lite: Navigation-oriented compression of 3d scene graphs under communication constraints. *arXiv preprint arXiv:2209.06111*, 1(2), 2022. 1
- [4] Tianshui Chen, Weihao Yu, Riquan Chen, and Liang Lin. Knowledge-embedded routing network for scene graph generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6163–6171, 2019. 1, 2
- [5] Zuyao Chen, Jinlin Wu, Zhen Lei, Zhaoxiang Zhang, and Changwen Chen. Gpt4sgg: Synthesizing scene graphs from holistic and region-specific narratives. *arXiv preprint arXiv:2312.04314*, 2023. 1
- [6] Zuyao Chen, Jinlin Wu, Zhen Lei, Zhaoxiang Zhang, and Changwen Chen. Expanding scene graph boundaries: Fully open-vocabulary scene graph generation via visual-concept alignment and retention. In *European Conference on Computer Vision (ECCV)*, pages 108–124, 2024. 1, 2, 3, 4, A5
- [7] Meng-Jiun Chiou, Henghui Ding, Hanshu Yan, Changhu Wang, Roger Zimmermann, and Jiashi Feng. Recovering the unbiased scene graphs from the biased ones. In *Proceedings of the 29th ACM international conference on multimedia*, pages 1581–1590, 2021. 2
- [8] Yuren Cong, Michael Ying Yang, and Bodo Rosenhahn. Reltr: Relation transformer for scene graph generation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(9): 11169–11183, 2023. 3
- [9] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes, 2017. A4
- [10] Qiao Gu, Ali Kuwajerwala, Sacha Morin, Krishna Murthy Jatavallabhula, Bipasha Sen, Aditya Agarwal, Corban Rivera, William Paul, Kirsty Ellis, Rama Chellappa, et al. Concept-graphs: Open-vocabulary 3d scene graphs for perception and planning. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5021–5028. IEEE, 2024. 1, 3
- [11] Tao He, Lianli Gao, Jingkuan Song, and Yuan-Fang Li. Towards open-vocabulary scene graph generation with prompt-based finetuning. In *European Conference on Computer Vision*, pages 56–73. Springer, 2022. 1, 2
- [12] Tao He, Lianli Gao, Jingkuan Song, and Yuan-Fang Li. Towards open-vocabulary scene graph generation with prompt-based finetuning, 2022. 3
- [13] Marcel Hildebrandt, Hang Li, Rajat Koner, Volker Tresp, and Stephan Günnemann. Scene graph reasoning for visual question answering. *arXiv preprint arXiv:2007.01072*, 2020. 1
- [14] Jitesh Jain, Jianwei Yang, and Humphrey Shi. Vcoder: Versatile vision encoders for multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27992–28002, 2024. A1
- [15] Justin Johnson, Ranjay Krishna, Michael Stark, Li-Jia Li, David Shamma, Michael Bernstein, and Li Fei-Fei. Image retrieval using scene graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3668–3678, 2015. 1
- [16] Siddhesh Khandelwal and Leonid Sigal. Iterative scene graph generation, 2022. 3
- [17] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv:2304.02643*, 2023. 1
- [18] Boris Knyazev, Harm de Vries, Cătălina Cangea, Graham W Taylor, Aaron Courville, and Eugene Belilovsky. Generative compositional augmentations for scene graph prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15827–15837, 2021. 2, 3
- [19] Zihan Kong and Haiwei Zhang. Opensgen: Fine-grained relation-aware prompt for open-vocabulary scene graph generation. In *Proceedings of the 2025 International Conference on Multimedia Retrieval*, pages 634–643, 2025. 1, 2
- [20] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123(1):32–73, 2017. 2, 3, A2
- [21] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 3, 4
- [22] Lin Li, Long Chen, Yifeng Huang, Zhimeng Zhang, Songyang Zhang, and Jun Xiao. The devil is in the labels: Noisy label correction for robust scene graph generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18869–18878, 2022. 1, 2
- [23] Rongjie Li, Songyang Zhang, and Xuming He. Sgtr: End-to-end scene graph generation with transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19486–19496, 2022. 2, 3
- [24] Rongjie Li, Songyang Zhang, Dahua Lin, Kai Chen, and Xuming He. From pixels to graphs: Open-vocabulary scene graph generation with vision-language models, 2024. 1, 3, A3
- [25] Xingchen Li, Long Chen, Wenbo Ma, Yi Yang, and Jun Xiao. Integrating object-aware and interaction-aware knowledge for weakly supervised scene graph generation. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 4204–4213, 2022. 2

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

- [26] Yikang Li, Wanli Ouyang, Bolei Zhou, Jianping Shi, Chao Zhang, and Xiaogang Wang. Factorizable net: an efficient subgraph-based framework for scene graph generation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 335–351, 2018. 1
- [27] Dayong Liang, Changmeng Zheng, Zhiyuan Wen, Yi Cai, Xiao-Yong Wei, and Qing Li. Seeing beyond the scene: Enhancing vision-language models with interactional reasoning. *arXiv preprint arXiv:2505.09118*, 2025. 1
- [28] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 1, 4
- [29] OpenAI et al. Gpt-4o system card, 2024. 3, 4
- [30] Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision, 2023. A1
- [31] Jungin Park, Jiyoun Lee, and Kwanghoon Sohn. Bridge to answer: Structure-aware graph interaction network for video question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15526–15535, 2021. 1
- [32] Jae Sung Park, Zixian Ma, Linjie Li, Chenhao Zheng, Cheng-Yu Hsieh, Ximing Lu, Khyathi Chandu, Quan Kong, Norimasa Kobori, Ali Farhadi, et al. Synthetic visual genome. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 9073–9086, 2025. 3
- [33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *CoRR*, abs/2103.00020, 2021. A1
- [34] Mohammed Suhail, Abhay Mittal, Behjat Siddiquie, Chris Broaddus, Jayan Eledath, Gerard Medioni, and Leonid Sigal. Energy-based learning for scene graph generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13936–13945, 2021. 2, 3
- [35] Kaihua Tang, Hanwang Zhang, Baoyuan Wu, Wenhan Luo, and Wei Liu. Learning to compose dynamic tree structures for visual contexts, 2018. 3
- [36] Kaihua Tang, Hanwang Zhang, Baoyuan Wu, Wenhan Luo, and Wei Liu. Learning to compose dynamic tree structures for visual contexts. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6619–6628, 2019. 1, 2
- [37] Kaihua Tang, Yulei Niu, Jianqiang Huang, Jiaxin Shi, and Hanwang Zhang. Unbiased scene graph generation from biased training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3716–3725, 2020. 2
- [38] Johanna Wald, Helisa Dharmo, Nassir Navab, and Federico Tombari. Learning 3d semantic scene graphs from 3d indoor reconstructions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 4
- [39] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, et al. Internvl3.5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025. 3
- [40] Abdelrhman Werby, Chenguang Huang, Martin Büchner, Abhinav Valada, and Wolfram Burgard. Hierarchical open-vocabulary 3d scene graphs for language-grounded robot navigation. In *First Workshop on Vision-Language Models for Navigation and Manipulation at ICRA 2024*, 2024. 1
- [41] Danfei Xu, Yuke Zhu, Christopher B Choy, and Li Fei-Fei. Scene graph generation by iterative message passing. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5410–5419, 2017. 1, 2, A3
- [42] Jianwei Yang, Jiasen Lu, Stefan Lee, Dhruv Batra, and Devi Parikh. Graph r-cnn for scene graph generation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 670–685, 2018. 1, 2
- [43] Jingkang Yang, Yi Zhe Ang, Zujin Guo, Kaiyang Zhou, Wayne Zhang, and Ziwei Liu. Panoptic scene graph generation. In *ECCV*, 2022. 2, 3, 4
- [44] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *arXiv:2406.09414*, 2024. A1
- [45] Yuan Yao, Ao Zhang, Xu Han, Mengdi Li, Cornelius Weber, Zhiyuan Liu, Stefan Wermtner, and Maosong Sun. Visual distant supervision for scene graph generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15816–15826, 2021. 2
- [46] Hang Yin, Xiuwei Xu, Zhenyu Wu, Jie Zhou, and Jiwen Lu. Sg-nav: Online 3d scene graph prompting for llm-based zero-shot object navigation. *Advances in Neural Information Processing Systems*, 37:5285–5307, 2024. 1
- [47] Qifan Yu, Juncheng Li, Yu Wu, Siliang Tang, Wei Ji, and Yueting Zhuang. Visually-prompted language model for fine-grained scene graph generation in an open world. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 21560–21571, 2023. 1, 2
- [48] Alireza Zareian, Svebor Karaman, and Shih-Fu Chang. Bridging knowledge graphs to generate scene graphs. In *European conference on computer vision*, pages 606–623. Springer, 2020. 2
- [49] Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14393–14402, 2021. 1, 2
- [50] Rowan Zellers, Mark Yatskar, Sam Thomson, and Yejin Choi. Neural motifs: Scene graph parsing with global context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5831–5840, 2018. 2
- [51] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training, 2023. 3, A1

- 508 [52] Yong Zhang, Yingwei Pan, Ting Yao, Rui Huang, Tao
509 Mei, and Chang-Wen Chen. Learning to generate language-
510 supervised and open-vocabulary scene graph using pre-trained
511 visual-semantic space. In *Proceedings of the IEEE/CVF Con-
512 ference on Computer Vision and Pattern Recognition (CVPR)*,
513 pages 2915–2924, 2023. [2](#), [3](#)
- 514 [53] Yong Zhang, Yingwei Pan, Ting Yao, Rui Huang, Tao
515 Mei, and Chang-Wen Chen. Learning to generate language-
516 supervised and open-vocabulary scene graph using pre-trained
517 visual-semantic space. In *Proceedings of the IEEE/CVF Con-
518 ference on Computer Vision and Pattern Recognition*, pages
519 2915–2924, 2023. [2](#)
- 520 [54] Yiwu Zhong, Jing Shi, Jianwei Yang, Chenliang Xu, and Yin
521 Li. Learning to generate scene graph from natural language
522 supervision. In *Proceedings of the IEEE/CVF International
523 Conference on Computer Vision*, pages 1823–1834, 2021. [2](#)
- 524 [55] Junsheng Zhou, Jinsheng Wang, Baorui Ma, Yu-Shen Liu,
525 Tiejun Huang, and Xinlong Wang. Uni3d: Exploring unified
526 3d representation at scale. In *International Conference on
527 Learning Representations (ICLR)*, 2024. [A1](#)

528

Appendix

529

A. Implementation Details for Architecture

530

A.1. Object-level Feature Extraction

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

An essential aspect of our feature extraction pipeline is object-level patch feature integration, which enriches relationship queries with fine-grained visual information. We explore three types of visual encoders: the CLIP image encoder, DINOv2-pretrained features, and the original SigLIP encoder used in LLaVA-OV. As shown in Table 6, the SigLIP features yield the most significant performance improvement, likely due to their stronger image-text alignment learned during pretraining. The benefit of using SigLIP is that it serves as the original vision encoder for LLaVA-OV, which not only ensures better feature alignment but also avoids introducing an additional vision encoder. For each object in a triplet, we use its downscaled ($384 \times 384 \rightarrow 27 \times 27$) 2D bounding box to query the corresponding image patch features. If a feature patch overlaps with the object’s 2D region, it is included in the computation. We then average the selected patch embeddings to form the object-level feature, which is concatenated with the corresponding text embedding.

Table 6. **Ablation study of patch visual model variants.** We evaluate the impact of different visual models on the object patch features on relation prediction accuracy.

| Vision Encoder | VG (In-domain 5K) | | PSG (Cross-domain 1K) | |
|--------------------|-------------------|----------------|-----------------------|------------------------------|
| | Acc \uparrow | Acc \uparrow | Acc \uparrow | Acc w/ Similarity \uparrow |
| <i>CLIP</i> [33] | 46.2 | 22.8 | 57.3 | |
| <i>DINOv2</i> [30] | 46.2 | 22.5 | 57.7 | |
| <i>SigLIP</i> [51] | 46.3 | 23.6 | 59.2 | |

549

A.2. Depth Exploration

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

While our method primarily relies on 2D visual and semantic cues, we explore the possibility of incorporating depth information into the relation prompting process, motivated by recent works such as VCoder [14]. Unlike approaches that use explicit 3D positional embeddings or train a specialized 3D encoder [55], requiring large-scale supervision to align with language space, we instead adopt an efficient strategy that reuses the pretrained visual encoder already aligned with the text modality for depth, as Sec. A.1. Specifically, we estimate monocular depth maps from 2D images using Depth Anything v2 [44], and normalize the predicted depth values to the $[0,1]$ range for consistency. For inference, the model can be seamlessly used with ground-truth depth or estimated depth. After obtaining the depth map, we feed it directly into the frozen image encoder. The object-level depth features could be extracted by pooling visual encoder patch embeddings corresponding to each object’s region in

the depth map, as the image embedding. They are concatenated after the $\langle \text{DEPTH} \rangle$ placeholder for each object, as shown in Fig. 2, and then fed into the backbone.

567

568

569

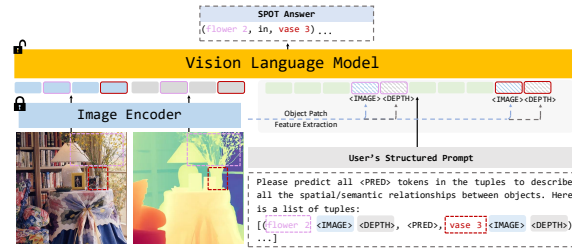


Figure 2. **SPOT method with depth for input**

Similarly to object-centric image embedding, we also explore different intuitive alternatives for depth embedding and compare quantitatively in Tab. 7. The exploration ranges from 3D positional embedding, an additional 3D encoder, and the reuse of the visual encoder, as in our model. For the 3D positional embedding, the depth is lifted into the point cloud to extract the bounding box minimum and maximum coordinates along the z-axis, and we use a fixed sinusoidal positional encoding to embed these 3D locations. For the additional 3D encoder, we use the same method to lift the point from 2D to 3D for the whole object point cloud, and use the pretrained Uni3D [55] encoder for embedding. An additional projection layer is trained to align the input space. As observed in Tab. 7, encoding the depth map with the VLM image encoder performs better. We hypothesize this is because the image encoder is already aligned with the LLM input space, and the depth patches are naturally aligned with the image patches since they are encoded by the same model.

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

Table 7. **Ablation study of different depth integration variants.**

| Depth Info Variant | VG (In-domain 5K) | | PSG (Cross-domain 1K) | |
|--------------------------------|-------------------|----------------|-----------------------|------------------------------|
| | Acc \uparrow | Acc \uparrow | Acc \uparrow | Acc w/ Similarity \uparrow |
| <i>3D Positional Embedding</i> | 46.3 | 23.2 | 58.4 | |
| <i>3D Encoder</i> [55] | 46.3 | 22.6 | 59.1 | |
| <i>SigLIP Depth</i> [51] | 46.4 | 23.6 | 59.5 | |

Empirically, given the above analysis, we observe that the inclusion of depth features yields marginal improvement in quantitative metrics. However, from a qualitative perspective, the model produces better spatial relations as shown in Fig. 10. While the improvements brought by depth are not quantitatively reflected in our current evaluation settings, we believe it presents a promising direction for future work on explicit spatial reasoning. Hence, we provide the depth integration option here for further development and investigation.

588

589

590

591

592

593

594

595

596

597

598 B. Implementation Details for Training

599 B.1. Data Preprocessing

600 In our experiments, we finetune the model on Visual Genome
601 150 [20], which serves as the primary dataset to ensure fair
602 comparison with other baseline methods presented in the
603 main paper. Input images are preprocessed to match the
604 expectations of the SigLIP vision encoder. All images are
605 resized to a resolution of 384×384 and normalized using
606 the standard mean and standard deviation values from the
607 model’s pretraining configuration. We do not apply any data
608 augmentation techniques such as random flipping during the
609 fine-tuning stage in order to maintain a consistent mapping
610 between visual cues and relational language. Bounding box
611 coordinates are extracted inside every object pair and nor-
612 malized to a range of $[0, 1]$ relative to the image dimensions
613 before being used for object-centric feature extraction.

614 B.2. Structured Prompt Details

615 Fig. 3 illustrates our structured prompt template with exam-
616 ple input. Instead of sequentially querying each object pair
617 for their relationship, we adopt a more efficient approach
618 by constructing a structured prompt template, illustrated in
619 Figure 3. Within this template, `<IMAGE>` (and `<DEPTH>`
620 added alternatively) are fixed special tokens that serve as
621 anchors for inserting object-level visual (and depth) fea-
622 tures, as described in Section A.1. These tokens are not
623 updated during training but are used to locate where the
624 visual embeddings should be injected. We also introduce
625 a special `<PRED>` token to indicate the position where the
626 model should predict the relationship between the object
627 pair. This design allows for explicit supervision and reduces
628 redundancy in the prediction process. Additionally, each ob-
629 ject name is followed by a numeric identifier to distinguish
630 between different instances of the same class.

```

{
  "from" "human",
  "value" "<image>\nYou are an agent specializing in identifying the physical and
spatial relationships in images for 3D mapping.\nYour task is to analyze the
images and output a list of tuples describing the physical relationships between
objects.\nNote that you are describing the physical relationships between
the objects inside the image. \nYou will also be given a text list of
relation tuples. The list will be in the format: [{"object 1", "<PRED>", "object
2"}, {"object 3", "<PRED>", "object 2"}].\nPlease predict all <PRED> in the
tuples.\nHere is the list of predicate tuples: [{"bench 0 <IMAGE>, <PRED>, grass
6 <IMAGE>}, {"sky 4 <IMAGE>, <PRED>, grass 6 <IMAGE>}, {"wall 8 <IMAGE>,
<PRED>, building 7 <IMAGE>}]. Please describe the spatial relationships between
the objects in all tuples."
}
{
  "from" "gpt",
  "value" "[{"bench 0, on, grass 6}, {"sky 4, over, grass 6}, {"sky 4, over,
building 7}, {"wall 8, enclosing, building 7}]"
}

```

Figure 3. SPOT structured prompt template.

631 B.3. Hyperparameter Tuning

632 During training, the vision encoder is kept frozen, while the
633 language model and adapter layers are trainable. We do not
634 apply AnyRes during either training or inference. The model
635 is optimized using Adam with a learning rate of 1×10^{-5} ,

a cosine learning rate scheduler, and a warm-up ratio of
0.03. Training is conducted on 8 NVIDIA A40 GPUs with a
batch size of 16 for 1 epoch. Due to out-of-memory issues
commonly encountered with long prompt sequences, we
limit the model’s maximum token length to 5000. All other
hyperparameters and training configurations follow those of
LLaVA-OV-7B.

636 C. Implementation Details for Inference

637 In this section, we specifically introduce how our model
638 works during inference to supplement the outline in Sec-
639 tion 3.3 in the main paper. There are two key design points
640 to meet our expectations: (1) The pipeline should work ro-
641 bustly for open-world applications with broad generalization
642 abilities. (2) The pipeline should have the ability to filter
643 and rank noisy relations. To satisfy these two requirements,
644 SPOT disentangles relation pruning and proposal from di-
645 rect relation prediction, leaving the VLM module only the
646 task for spatial reasoning. The relation proposals are pro-
647 vided by out-of-the box detection models, leveraging the
648 most recent advancements in this area. We observed that us-
649 ing this simple but effective approach, our framework could
650 detect fine-grained relations and could be seamlessly applied
651 to different scenarios, providing more comprehensive graphs
652 compared to directly using VLMs to process everything.
653 However, this framework raises three challenges:

- 654 • The quantity of detected bounding boxes is generally too
655 large.
 - 656 • The relation proposal sequences are overly long, making
657 it inefficient for VLM to output long texts.
 - 658 • There are many spatially implausible object pairs existing.
- 659 To mitigate these problems, we apply some strategies to
660 improve the overall performance

661 Over-detection Suppression. In detection results, multi-
662 ple bounding boxes may correspond to the same real-world
663 object but have different yet reasonable category labels. For
664 instance, a person may be detected as both “man” and “per-
665 son,” and both exist in the vocabulary. To mitigate this, we
666 apply cross-category Non-Maximum Suppression (NMS)
667 apart from the standard NMS to reduce overlapping boxes
668 across semantically similar labels. For each remaining ob-
669 ject, we aggregate all plausible labels into a label group.
670 During evaluation, if the ground-truth label is present in this
671 group, the object is considered correctly classified. For the
672 real application, any label works for usage.

673 Filter and Rank Relations. To address this problem, we
674 incorporate the classification probabilities from the object
675 detector as object-level confidence scores, selecting the first
676 100 objects of higher probabilities. For relation-level filter-
677 ing, we compute the distance between the centers of two
678 objects. We explored three different types of distance for
679 reference.

680 1 Standard normalized geometry distance: $d = \frac{\|c_i - c_j\|_2}{\sqrt{H^2 + W^2}}$, 687


```

Prediction: [<mouse3> <next to> <keyboard 4>]
Ground Truth: [<mouse3> <beside> <keyboard 4>]

Prompt: Given (mouse next to keyboard), (mouse beside keyboard),
whether these two phrases indicate similar spatial relations between
two objects inside the the tuple. Only return Yes or No without
further explanation.

Answer: Yes

```

Figure 6. LLM as a judge to evaluate the similarity between two relation words

the results are evaluated strictly based on the directional relation that is present in the annotations. If the model predicts the inverse (e.g., “B under A”), it is not considered correct unless it exactly matches the ground truth. While this setting ensures consistency with the prior work and fair comparison, we acknowledge that extending the evaluation protocol to account for equivalent relationships in the opposite direction is an interesting future work to explore.

Cross-Domain Evaluation. A common limitation in existing evaluation protocols for cross-domain datasets is that semantically similar but lexically different relations are treated as incorrect. For instance, synonymous expressions like “beside” and “next to” may convey the same meaning but only one may appear in the ground truth, leading to unfair penalization. To address this, we introduce a more semantically-aware evaluation method by leveraging a large language model (Qwen-7B) as a judge. Given a relation prediction and its corresponding ground-truth triplet, we prompt the LLM with a predefined template (shown in Figure 6) to determine whether the predicted relation is semantically equivalent to the ground truth. If the model answers yes, the prediction is considered correct. This enhanced evaluation metric is reported as “Recall w/ Similarity” or “Accuracy w/ Similarity”. This evaluation protocol aims to give model credit for correct relations expressed in an open vocabulary that are semantically equivalent to the ground truth.

The whole cross-domain evaluation pipeline is as follows, including some specific rules for simple cases:

- **Case 1:** prediction relation == ground truth relation → correct
- **Case 2:** prediction contains the same spatial word as the label but with the inclusion of “on”, “in”, or “from” → correct:
 - In our experiments on the 3DSSG dataset, we use “standing on/on”, “built in/in”, “hanging on/hanging from”.
 - While these rules account for common cases that are clearly correct, it is difficult to iterate all plausible rules comprehensively. Thus, for all other cases, we propose to use LLM as a judge in order to assess the accuracy of the model’s open vocabulary predictions against a closed vocabulary label space.
- **Case 3:** otherwise, we use LLM (Qwen-7B) to judge. The prompt is as Fig. 6:
 - We checked the results of LLM’s judgment to verify its

ability to judge semantically equivalent spatial phrases and found it to be accurate in practice. Given its strong performance and flexibility, we adopt the LLM-based judgement for all cross-domain tasks.

E. More Qualitative Results

In this section, we provide more qualitative results on different datasets. In Fig. 7, Fig. 8, Fig. 9, we separately display more visualizations on three cross-domain datasets: PSG, 3DSSG, and ScanNet [9], showing the generalization ability of our model. In Fig. 10, we also present results on PSG that illustrate how SPOT predicts richer spatial relations (e.g. “behind”), even though the ground-truth answers are different. This finding raises the need for a more robust and comprehensive evaluation method to quantify different perspectives of the scene graph prediction performance.

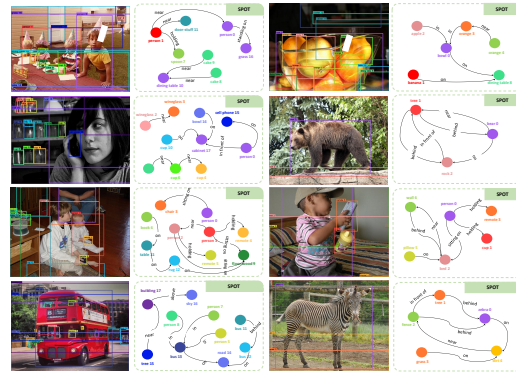


Figure 7. Additional cross-domain qualitative results on Panoptic Scene Graph Dataset. We visualize only the relations that exist in the ground truth.

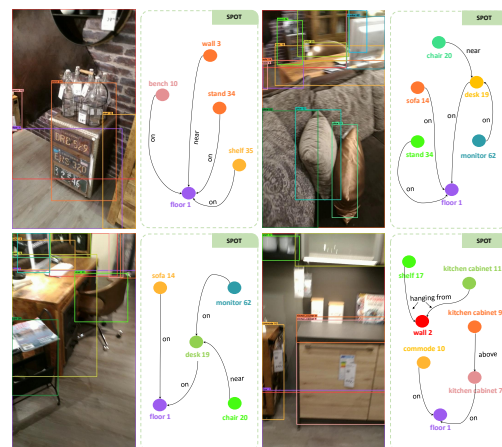


Figure 8. Additional cross-domain qualitative results on 3DSSG. We visualize only the relations that exist in the ground truth.

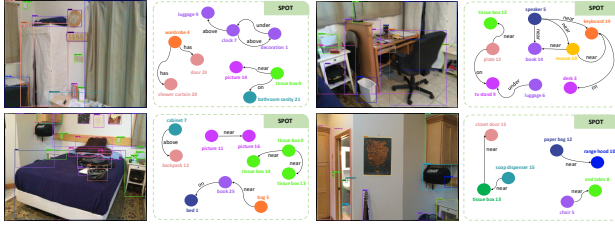


Figure 9. Qualitative cross-domain results on ScanNet. In contrast to PSG and 3DSSG, ScanNet is not a standard scene graph dataset. We include these results to show the generalization ability of SPOT and effectiveness in practical use with a real object detections.

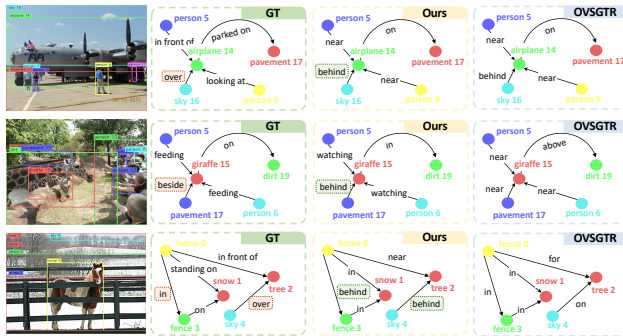


Figure 10. Additional cross-domain qualitative results on Panoptic Scene Graph Dataset. We observe that SPOT predicts richer spatial relations, which we highlight with green boxes.

823 F. Limitations

824 While SPOT demonstrates strong generalization by leveraging
 825 a vision-language model (VLM), it remains inherently
 826 constrained by the prior knowledge and reasoning abilities
 827 of the underlying language model. This reliance limits its
 828 ability to handle relationships or concepts that fall too differ-
 829 ent from the pretrained distribution. In terms of efficiency,
 830 SPOT trades off throughput for generalization. Unlike tradi-
 831 tional approaches such as OvSGTR [6], which utilize
 832 lightweight architectures like Grounding DINO and simple
 833 MLP classifiers, our VLM-based method incurs additional
 834 computational overhead. Conventional methods can predict
 835 thousands of relation pairs within seconds, whereas prompt-
 836 ing a large language model for the same quantity of relations
 837 introduces latency and resource demands. Another potential
 838 limitation arises from the length of the prediction sequence.
 839 As the number of predicted triplets grows, the prompt length
 840 increases accordingly, which can lead to degraded perfor-
 841 mance due to the model’s limited context window. Although
 842 we attempt to mitigate this by segmenting long prompts into
 843 smaller batches during inference, this issue of forgetting re-
 844 mains a broader challenge in both the scene graph generation
 845 downstream task and general language modeling.

G. Social impact & Safeguards

846 A key application of 3D scene graph models is their integra-
 847 tion into embodied agent systems for planning and interac-
 848 tion, where agents rely on spatial relation triplets and object
 849 locations to make decisions. Prediction of scene graphs
 850 for this purpose offers the potential for more explainable
 851 decision-making in these critical systems. However, cur-
 852 rent quantitative evaluation protocols for open-world 3D
 853 scene graph generation remain limited. To explore the full
 854 value of these systems to improve the safety of robotic ap-
 855 plications, further quantitative vetting becomes essential to
 856 ensure that both relationship predictions and object detec-
 857 tions faithfully represent the physical environment. Without
 858 such safeguards, erroneous predictions could result in agents’
 859 interaction with objects in unsafe or unintended ways, posing
 860 potential risks in real-world deployment. 861