

Spatio-Structural Image Captioning via LLM Fine-Tuning with Depth-Enhanced Scene Graphs

Le Vinh Thuan¹, Nguyen Minh Khoa², Nguyen Vinh Thanh³, and Nguyen Thi Dinh⁴,

^{1,2,3}Faculty of Information Technology, University of Science, VNU-HCM,
Ho Chi Minh City, Vietnam,

⁴Ho Chi Minh City University of Industry and Trade, Ho Chi Minh City, Vietnam

¹lvthuan23@apcs.fitus.edu.vn, ²nmkhoa23@apcs.fitus.edu.vn,

³23120012@student.hcmus.edu.vn, ⁴dinhnt@huit.edu.vn

Abstract. The abstract should briefly summarize the contents of the paper in 150–250 words.

Keywords: First keyword · Second keyword · Another keyword.

1 Introduction

The automatic generation of natural language descriptions from images remains one of the most challenging and important tasks in computer vision and natural language processing [1, 13]. This interdisciplinary field, commonly known as image captioning, requires machines to not only perceive objects and their relationships but also to generate coherent, contextually appropriate text that captures the scene’s semantic essence [1, 17]. Recent advances in deep learning, particularly the adoption of Transformer-based models and Large Language Models (LLMs) in Vision-Language Pre-training (VLP) [?, ?], have significantly improved descriptive fluency. However, a critical limitation persists: current VLP and image captioning methods primarily rely on 2D feature extraction and lack a deep, grounded understanding of the 3D structure and spatial relationships within a scene. Consequently, generated captions often fail to accurately describe position, depth, and structural context (e.g., "the car is in front of the building" without specifying the distance or depth plane), limiting their applicability in domains requiring precise spatial reasoning, such as robotics or complex scene analysis. To address this gap, Knowledge Graphs (KGs) offer a crucial pathway by providing structured semantic information and enabling reasoning capabilities [10, 16]. While previous work has used KGs to incorporate general semantic and factual knowledge [?, 19], few have effectively integrated explicit structural and geometric knowledge like scene depth and 3D spatial relationships. Furthermore, existing KG-enhanced methods often lack commonsense reasoning, which is vital for generating captions that reflect human-like understanding of causality and purpose [22]. This paper introduces a novel framework, the Depth- and Commonsense-Augmented Scene Graph Captioner (DASG-CS

Captioner), designed to generate fine-grained image descriptions with enhanced Spatio-Structural and Commonsense-Grounded awareness. We propose a comprehensive pipeline that moves beyond 2D perception by integrating state-of-the-art depth estimation into the knowledge representation stage, thereby enriching the generated captions with precise spatial details. Our method leverages the structured reasoning capabilities of Knowledge Graphs and the generative power of Fine-Tuned LLMs. The main contributions of this work are: A Novel Depth- and Commonsense-Augmented Scene Graph (DASG-CS): We propose a new KG architecture that seamlessly fuses visual features, quantitative depth information (from Depth Anything v2), and commonsense knowledge (from ConceptNet) to create a multi-faceted, structured representation of the scene. Spatially-Aware LLM Fine-Tuning Strategy: We develop an effective encoding and fine-tuning strategy for LLMs, enabling the model to explicitly utilize the Spatio-Structural and Commonsense information encoded in the DASG-CS, leading to generated captions that are demonstrably richer in depth and contextual reasoning. A Comprehensive Four-Stage Pipeline: We present a robust and reproducible pipeline that integrates feature extraction, multi-source knowledge fusion, structured data encoding, and natural language generation, setting a new benchmark for spatially and contextually grounded image captioning. Extensive Evaluation with Novel Metrics: We conduct a comprehensive evaluation demonstrating the effectiveness of the DASG-CS framework, including specialized metrics to quantify the improvement in Spatial Reasoning Accuracy and Contextual Richness of the generated descriptions.

2 Related Work

Image captioning has evolved significantly from simple encoder-decoder architectures [13] to sophisticated systems incorporating spatial awareness and external knowledge [1]. Our work builds upon three critical research directions: depth-aware image understanding, knowledge graph construction for vision-language tasks, and fine-tuning language models for spatial-aware captioning.

2.1 Object Detection and Spatial Feature Extraction

Traditional image captioning systems rely primarily on 2D visual features extracted from RGB images [1,17]. The YOLO family has evolved from YOLOv1’s single-stage detection [12] through multiple generations [2,3,14] to more sophisticated architectures. Modern object detectors like YOLOv11 provide real-time performance with improved accuracy, particularly for small objects and crowded scenes. While object detection has advanced significantly, most systems still focus on 2D bounding boxes without incorporating depth information for true 3D spatial understanding.

2.2 Scene Understanding and Spatial Relationships

Scene graph generation has emerged as a method to capture object relationships explicitly [19, 21]. These approaches construct structured representations of visual scenes by identifying objects and their relationships [23]. However, most scene graphs focus on semantic relationships rather than precise metric spatial relationships. Our approach aims to incorporate depth estimation to enable more accurate spatial relationship modeling.

2.3 Knowledge-Enhanced Image Captioning

Knowledge graphs provide structured representations of visual scenes and external world knowledge, enabling richer semantic understanding. Visual Genome [8] pioneered large-scale visual knowledge graphs with detailed annotations of objects, attributes, and relationships. Flickr30k [20] provide image-caption pairs widely used for training vision-language models.

Several works have explored incorporating external knowledge into image captioning. Knowledge graphs have been used for visual question answering [15] and image classification [10]. Entity linking techniques [4] help map visual detections to knowledge base entities. Knowledge-based approaches [22, 24] have shown improvements by incorporating common-sense and factual knowledge.

Most existing approaches integrate knowledge at the generation stage [16]. Our pipeline differs by constructing a comprehensive knowledge graph that fuses visual observations from datasets like Visual Genome and Flickr with external knowledge sources before the generation phase.

2.4 Natural Language Generation for Image Captioning

Natural language generation is the final component of image captioning systems. Early encoder-decoder models [7, 13] used RNNs to generate captions from visual features. Attention mechanisms [1, 17] improved performance by allowing models to focus on relevant image regions during generation.

More recent approaches employ transformer-based architectures [5, 9] for better long-range dependency modeling. The T5 model [11] treats all NLP tasks as text-to-text problems, offering a unified framework adaptable to various generation tasks including image captioning.

Our approach leverages text-to-text transformers for generating descriptions from structured semantic representations. By fine-tuning on data that explicitly encodes spatial relationships and depth information, we aim to generate descriptions that accurately reflect 3D spatial configurations.

2.5 Research Gaps and Our Contributions

Existing image captioning systems face several limitations:

1. **Limited Spatial Awareness:** Most systems [13, 17] rely on 2D visual features without depth information for true 3D spatial understanding.

2. **Knowledge Integration:** While some works integrate external knowledge [10, 22], comprehensive frameworks that systematically combine visual observations with external knowledge graphs remain limited.
3. **Spatial Description Quality:** Generated descriptions often lack precise spatial relationships and depth-aware positioning information.

Our proposed three-system pipeline addresses these gaps:

1. **System 1 - Spatial Feature Extraction:** We combine YOLOv11 for object detection with Depth Anything V2 for monocular depth estimation, enabling extraction of objects, their relationships, and depth information.
2. **System 2 - Knowledge Graph Framework:** We construct knowledge graphs based on Visual Genome [8] and Flickr [20] datasets, enriched with external knowledge from ConceptNet and DBpedia.
3. **System 3 - Depth-Aware Caption Generation:** We fine-tune language models to generate image descriptions with emphasis on spatial and depth features, producing captions that accurately reflect 3D spatial configurations.

The key contribution is the systematic integration of depth information throughout the pipeline, from detection through knowledge graph construction to caption generation, emphasizing spatial accuracy in generated descriptions.

3 Theoretical Basis

This section presents the theoretical foundations underlying our four-stage architecture, providing formal definitions and mathematical models that support the integration of computer vision, knowledge graph reasoning, and natural language generation.

3.1 Object Detection and Spatial Analysis Framework

YOLOv11 Detection Model The YOLOv11 object detection framework can be formalized as a function $f_{YOLO} : \mathbb{R}^{H \times W \times 3} \rightarrow \mathcal{D}$, where the input image $I \in \mathbb{R}^{H \times W \times 3}$ is mapped to a detection set $\mathcal{D} = \{d_1, d_2, \dots, d_n\}$. Each detection d_i is represented as:

$$d_i = (bbox_i, c_i, s_i) \quad (1)$$

where $bbox_i = (x_i, y_i, w_i, h_i)$ represents the bounding box coordinates, $c_i \in \mathcal{C}$ is the class label from the predefined class set \mathcal{C} , and $s_i \in [0, 1]$ is the confidence score.

Depth Estimation and Spatial Reasoning We incorporate the Depth-Anything-V2 model [18] to estimate depth information, formally defined as:

$$f_{depth} : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}^{H \times W} \quad (2)$$

The depth map $D = f_{depth}(I)$ provides spatial context that enhances object relationship understanding. Combined with the Shapely geometry library [6], we compute spatial relationships between detected objects using geometric operations:

$$R_{spatial}(d_i, d_j) = \{distance(d_i, d_j), overlap(d_i, d_j), relative_position(d_i, d_j)\} \quad (3)$$

3.2 Knowledge Graph Construction Theory

Entity Linking and Knowledge Fusion Given the detection set \mathcal{D} from Stage I, we define the entity linking function as:

$$f_{link} : \mathcal{D} \times \mathcal{KB} \rightarrow \mathcal{E}_{linked} \quad (4)$$

where \mathcal{KB} represents the external knowledge base and \mathcal{E}_{linked} is the set of linked entities. The linking process employs semantic similarity scoring:

$$sim(d_i, e_k) = \alpha \cdot sim_{text}(c_i, label(e_k)) + \beta \cdot sim_{context}(context(d_i), context(e_k)) \quad (5)$$

where α and β are weighting parameters, and sim_{text} and $sim_{context}$ represent textual and contextual similarity measures respectively.

Knowledge Graph Formalization The constructed knowledge graph is formally defined as a directed graph $KG = (V, E, R, A)$ where $V = \{v_1, v_2, \dots, v_m\}$ is the set of entity vertices, $E \subseteq V \times V$ represents the edges between entities, R is the set of relation types, and $A : V \rightarrow \mathcal{P}(\mathcal{A})$ maps entities to their attribute sets.

3.3 T5-based Natural Language Generation Framework

Text-to-Text Transfer Learning Our approach leverages the T5 (Text-to-Text Transfer Transformer) model [11] for generating natural language descriptions from structured semantic representations. The T5 framework treats all NLP tasks as text-to-text problems:

$$f_{T5} : \mathcal{S}_{semantic} \rightarrow \mathcal{T}_{text} \quad (6)$$

where $\mathcal{S}_{semantic}$ represents the semantic input derived from the knowledge graph and \mathcal{T}_{text} is the generated natural language output.

Semantic Representation Encoding The knowledge graph entities and relationships are encoded into a structured semantic representation $S_{semantic}$ that serves as input to the T5 model:

$$S_{semantic} = encode(KG, R_{spatial}, C_{context}) \quad (7)$$

where $C_{context}$ represents the contextual information derived from the enriched data analysis.

Attention Mechanism for Semantic Focus The T5 model employs multi-head self-attention to focus on relevant semantic elements:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (8)$$

where Q , K , and V represent query, key, and value matrices respectively, derived from the semantic representation.

3.4 Contextual Reasoning and Inference Theory

Enriched Data Analysis Our approach performs contextual reasoning using enriched data from the knowledge graph without explicit relationship modeling. The contextual inference function is defined as:

$$f_{inference} : \mathcal{E}_{enriched} \rightarrow \mathcal{I}_{context} \quad (9)$$

where $\mathcal{E}_{enriched}$ represents the enriched entity set and $\mathcal{I}_{context}$ denotes the inferred contextual information.

Semantic Coherence Optimization To ensure semantic coherence in the generated descriptions, we optimize the following objective function:

$$\mathcal{L}_{total} = \mathcal{L}_{generation} + \lambda_1 \mathcal{L}_{coherence} + \lambda_2 \mathcal{L}_{semantic} \quad (10)$$

where $\mathcal{L}_{generation}$ is the standard language modeling loss, $\mathcal{L}_{coherence}$ measures textual coherence, $\mathcal{L}_{semantic}$ ensures semantic consistency with the input knowledge graph, and λ_1 , λ_2 are regularization parameters.

3.5 Multi-Modal Integration Theory

Cross-Modal Alignment The integration of visual features, spatial information, and semantic knowledge requires cross-modal alignment. We define the alignment function as:

$$f_{align} : \mathcal{F}_{visual} \times \mathcal{F}_{spatial} \times \mathcal{F}_{semantic} \rightarrow \mathcal{F}_{unified} \quad (11)$$

where \mathcal{F}_{visual} , $\mathcal{F}_{spatial}$, and $\mathcal{F}_{semantic}$ represent visual, spatial, and semantic feature spaces respectively, and $\mathcal{F}_{unified}$ is the unified representation space.

Information Fusion Strategy The multi-modal information fusion employs weighted combination of feature representations:

$$\mathcal{F}_{unified} = \sum_{i=1}^3 w_i \cdot \phi_i(\mathcal{F}_i) \quad (12)$$

where ϕ_i represents the transformation function for each modality, and w_i are learned weights that adapt to the relative importance of each information source.

This theoretical framework provides the mathematical foundation for our four-stage architecture, ensuring principled integration of computer vision, knowledge representation, and natural language generation components.

4 Proposed Approach

4.1 Overall Architecture

Our proposed architecture consists of three main stages: (1) Spatial Feature Extraction, (2) Knowledge Graph Construction, and (3) Depth-Aware Caption Generation. Each stage is designed to progressively enrich the input image data with spatial and semantic information, culminating in the generation of detailed and contextually relevant captions.

4.2 Entities extraction

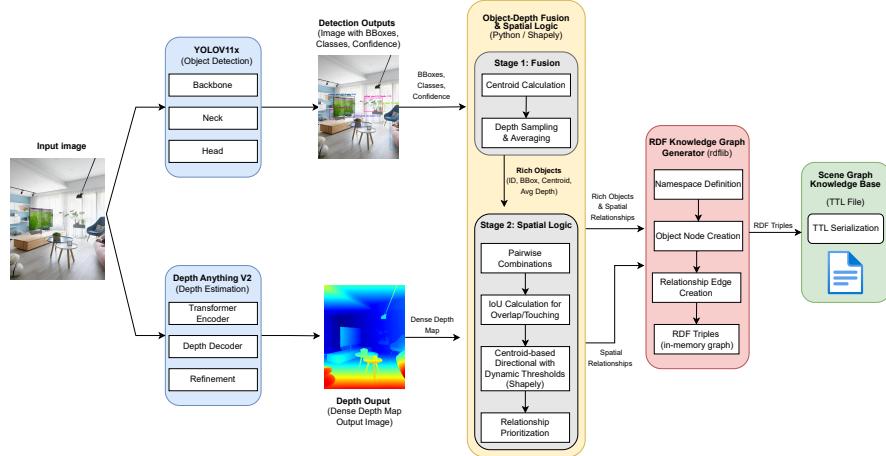


Fig. 1. Overall Architecture of the Proposed Approach

5 Evaluation

5.1 Environment and Experiment Data

6 Future Work

Acknowledgments

The authors would like to thank the Faculty of Information Technology, Ho Chi Minh City University of Industry and Trade, and University of Science, VNU-HCM, which are sponsors of this research. We also thank anonymous reviewers for their helpful comments on this paper.

References

1. Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., Zhang, L.: Bottom-up and top-down attention for image captioning and visual question answering. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 6077–6086 (2018)
2. Bochkovskiy, A., Wang, C.Y., Liao, H.Y.M.: Yolov4: Optimal speed and accuracy of object detection. arXiv preprint arXiv:2004.10934 (2020)
3. Cao, C., Song, Y., He, J.: An overview of yolo series: From yolov1 to yolov8. Journal of Physics: Conference Series **2685**(1), 012002 (2023)
4. Chen, Y., Ma, L., Huang, Y., Zhang, S., Qian, Y.: A survey of knowledge graph-based entity linking from text to knowledge graph. IEEE Access **9**, 13000–13020 (2021)
5. Cornia, M., Stefanini, M., Baraldi, L., Cucchiara, R.: Meshed-memory transformer for image captioning. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 10578–10587 (2020)
6. Gillies, S.: Shapely: manipulation and analysis of planar geometric objects (2007), python package, available at <https://pypi.org/project/Shapely/>
7. Karpathy, A., Fei-Fei, L.: Deep visual-semantic alignments for generating image descriptions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3128–3137 (2015)
8. Krishna, R., et al.: Visual genome: Connecting language and vision using crowd-sourced dense image annotations. International Journal of Computer Vision **123**(1), 32–73 (2017)
9. Li, G., Zhu, L., Liu, P., Yang, Y.: Entangled transformer for image captioning. In: IEEE/CVF International Conference on Computer Vision (ICCV). pp. 8928–8937 (2019)
10. Marino, K., Salakhutdinov, R., Gupta, A.: The more you know: Using knowledge graphs for image classification. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2673–2681 (2017)
11. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of Machine Learning Research **21**(140), 1–67 (2020)
12. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 779–788 (2016)

13. Vinyals, O., Toshev, A., Bengio, S., Erhan, D.: Show and tell: A neural image caption generator. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3156–3164 (2015)
14. Wang, C.Y., Bochkovskiy, A., Liao, H.Y.M.: Yolov7: Trainable bag-of-freebies for real-time object detection. arXiv preprint arXiv:2207.02696 (2022)
15. Wang, P., Wu, Q., Shen, C., van den Hengel, A.: Fvqa: Fact-based visual question answering. IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI) **40**(10), 2413–2427 (2018)
16. Wu, Q., Shen, C., Wang, P., Dick, A., van den Hengel, A.: Image captioning and visual question answering based on attributes and external knowledge. IEEE Transactions on Pattern Analysis and Machine Intelligence **40**(6), 1367–1381 (2017)
17. Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhutdinov, R., Zemel, R.S., Bengio, Y.: Show, attend and tell: Neural image caption generation with visual attention. In: Proceedings of the International Conference on Machine Learning (ICML). pp. 2048–2057 (2015)
18. Yang, L., Kang, B., Huang, Z., Zhao, Z., Xu, X., Feng, J., Zhao, H.: Depth anything v2. arXiv preprint arXiv:2406.09414 (2024)
19. Yao, T., Pan, Y., Li, Y., Mei, T.: Exploring visual relationship for image captioning. In: European Conference on Computer Vision (ECCV). pp. 684–699 (2018)
20. Young, P., Lai, A., Hodosh, M., Hockenmaier, J.: From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. Transactions of the Association for Computational Linguistics **2**, 67–78 (2014)
21. Zellers, R., Yatskar, M., Thomson, S., Choi, Y.: Neural motifs: Scene graph generation with global context. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5831–5840 (2018)
22. Zhao, W., Wu, B., Ma, S.: Knowledge enhanced fine-grained image captioning. In: ACM International Conference on Multimedia (MM). pp. 1631–1640 (2021)
23. Zhong, Y., Wang, L., Chen, J., Yu, D., Li, Y.: Comprehensive image captioning via scene graph decomposition. In: European Conference on Computer Vision (ECCV). pp. 211–229 (2020)
24. Zhu, Y., Yang, Z., Salakhutdinov, R., Xing, E.P.: Incorporating commonsense knowledge into image captioning via graph convolutional networks. In: IEEE/CVF International Conference on Computer Vision (ICCV). pp. 9008–9017 (2019)