# Beyond Generation: Multi-Hop Reasoning for Factual Accuracy in Vision-Language Models

**Shamima Hossain**[1,2]

[1]Department of Computer Science, Brac University, Bangladesh
[2]bKash Limited, Dhaka, Bangladesh
shamima.hossain@g.bracu.ac.bd, shamima.alma@bkash.com

## Abstract

Visual Language Models (VLMs) are powerful generative tools but often produce factually inaccurate outputs due to a lack of robust reasoning capabilities. While extensive research has been conducted on integrating external knowledge for reasoning in large language models (LLMs), such efforts remain underexplored in VLMs, where the challenge is compounded by the need to bridge multiple modalities seamlessly. This work introduces a framework for knowledge-guided reasoning in VLMs, leveraging structured knowledge graphs for multi-hop verification using image-captioning task to illustrate our framework. Our approach enables systematic reasoning across multiple steps, including visual entity recognition, knowledge graph traversal, and fact-based caption refinement. We evaluate the framework using hierarchical, triple-based and bullet-point based knowledge representations, analyzing their effectiveness in factual accuracy and logical inference. Empirical results show that our approach improves factual accuracy by approximately 31% on preliminary experiments on a curated dataset of mixtures from Google Landmarks v2, Conceptual captions and Coco captions revealing key insights into reasoning patterns and failure modes. This work demonstrates the potential of integrating external knowledge for advancing reasoning in VLMs, paving the way for more reliable and knowledgable multimodal systems.

## 1 Introduction

Visual Language Models have transformed image understanding tasks, yet their inability to reason systematically about facts within images and text remains a critical limitation. While humans naturally verify visual information against their knowledge base, VLMs lack structured mechanisms for fact verification, leading to confident but incorrect assertions about entities, locations, and relationships in images. This poses significant concerns
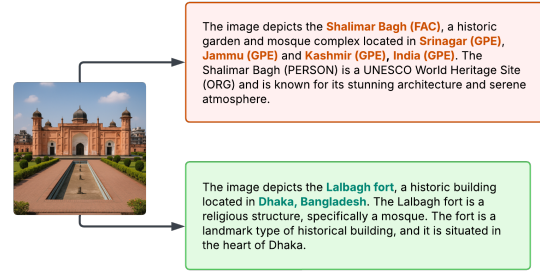


Figure 1: A comparison of hallucinated entities in red and the factually correct entities after processing through our pipeline in green.

about trustworthiness and reliability about their generated responses making them unreliable for domains like healthcare, education and cultural preservation.

Unlike in LLMs, where integrating external knowledge for reasoning is actively studied using retrieval Lewis et al. (2021) and in-context learning Brown et al. (2020), systematic reasoning in VLMs remains underexplored particularly in tasks requiring fact verification and multi-step logical inference. This is especially problematic for factual verification tasks, where models must not only recognize visual elements but also reason about their relationships with real-world knowledge. Current VLMs lack structured mechanisms to perform such reasoning, often resulting in descriptions that combine accurate visual observations with incorrect factual assertions. For example, while a VLM might correctly identify architectural features of a historical landmark, it fails to systematically verify and reason about crucial facts like its location, historical significance, or cultural context.

This limitation stems from the absence of explicit reasoning paths between visual perception and knowledge integration, leading to unreliable factual claims and compromised utility in applications

requiring high factual precision. VLMs face several unique technical challenges: they must jointly align facts with both image and textual data simultaneously, traverse complex knowledge structures across multiple reasoning steps, and maintain consistency between visual evidence and external knowledge sources.

These challenges are compounded by the need to represent knowledge in a format that supports both visual grounding and logical inference. To address these challenges, we propose a structured reasoning framework that explicitly models the verification path from visual perception to knowledge integration through multiple coordinated hops. To address these challenges, we introduce a multi-hop reasoning framework that enables VLMs to perform structured verification using knowledge graphs.

Our framework decomposes the verification process into distinct reasoning hops: entity recognition from visual inputs, knowledge graph traversal for fact retrieval, and structured verification of generated descriptions. Each hop is designed to maintain interpretable reasoning paths, allowing the model to explicitly track how it verifies facts against both visual evidence and knowledge sources. We implement this through three key innovations: (1) a hierarchical knowledge representation that supports both visual and factual reasoning, (2) a structured verification mechanism that traces reasoning paths through the knowledge graph, and (3) an adaptive correction strategy that resolves conflicts between visual observations and stored knowledge. Central to our approach is the flexible use of different knowledge representation formats namely hierarchical trees, relation triples, and structured facts each optimized for different types of reasoning tasks. This flexibility allows the model to choose appropriate reasoning paths based on the verification task, whether comparing spatial relationships, verifying historical facts, or checking entity attributes. We provide detailed ablation studies to trace the model's reasoning in each of these knowledge representation. Through extensive experimentation on landmark description tasks, we demonstrate that our framework significantly reduces hallucination while improving factual consistency in generated descriptions.

## 2 Related Studies

### 2.1 Vision-Language Models for Image Captioning

Vision-language models are now being for image captioning in the wild by leveraging their dual architectures to generate and align multi-modal embeddings. Early models, such as Show and Tell Vinyals et al. (2015), paired convolutional neural networks with recurrent neural networks for image description tasks. With the widespread development of transformer-based architectures, models like CLIP Radford et al. (2021) demonstrated state-of-the-art performance by aligning textual and visual embeddings through contrastive learning. Recent work has further demonstrated that incorporating external knowledge can significantly improve vision-language models' factual accuracy and reasoning capabilities as shown by Anderson et al. Anderson et al. (2018). Building on this, Zhang et al. Kang et al. (2023) proposed a knowledge-aware transformer architecture that explicitly reasons over both visual features and structured knowledge.

### 2.2 Studying hallucinations in image captioning task

Despite these advancements, VLMs often generate captions with hallucinated or factually inaccurate entities, especially for unseen or domain-specific data as supported by Rohrbach et al. Rohrbach et al. (2019). We were particularly inspired by their findings of how image captioning models often fail to capture image relevance with their internal understanding. Techniques such as entity-aware training Cao and Wang (2021) have been proposed to mitigate this issue. Integrating structured knowledge, such as ontologies or databases, has shown promise in improving factual accuracy for image captioning tasks. For instance, memory-augmented models Cornia et al. (2020) that retrieved relevant facts to enhance descriptions were introduced to overcome this limitation. This kind of inaccuracies limit VLMs usage in applications requiring factual precision, such as education, historical archiving etc.

### 2.3 Knowledge Graph Integration Approaches

Different approaches for integrating knowledge graphs with neural models have been explored. There has been significant efforts to integrate KGs into neural models and also into use language mod-

els to enhance KGs such as KG-BERT Yao et al. (2019) and TransE (Bordes et al., 2013), which embed graph knowledge into vector spaces for downstream tasks. In addition, Kumar et al. Zhou et al. (2024) studied projecting domain specific knowledge from a custom curated knowledge base into the latent space of the language model for pretraining the LM.

The use of KGs remains underexplored for image-captioning tasks especially when we want to understand how VLMs jointly reason about the image and their underlying knowledge representation. Our work in integrating VLMs with KGs for factual image captioning, as explored in this study, builds upon these foundational works, addressing the gap in ensuring factual accuracy in multi-modal systems. Our contributions are as follows:

- We propose a straight-forward integration of knowledge graphs illustrated with the example of landmark identification

- Our approach combines both pre-generation knowledge integration and post-generation correction

- We introduce new methods for comparing different knowledge representation formats

- We provide quantitative analysis of factual accuracy improvements

## 3 Multi-Hop Reasoning Framework

We introduce a multi-hop reasoning framework that systematically verifies and corrects these descriptions using structured knowledge as shown in Figure 2. Our framework consists of five key components that progressively refine the caption generation process:

**Vision-Language Understanding** The initial component leverages a pre-trained VLM (Qwen2-VL-2B-Instruct) Wang et al. (2024) to generate a base caption C from input image I. While this produces fluent descriptions, we observed that approximately 69% of entity mentions were either incorrect or hallucinated.

**Entity Extraction Hop** This component extracts named entities E = $e_1,...,e_n$ from C using spaCy's NER model Honnibal and Montani (2017). We focus on entities across multiple categories including locations, organizations, and facilities to capture the full range of factual claims made in the caption.

**Knowledge Graph Navigation** For each extracted entity $e_i$, we perform both exact and fuzzy matching against our knowledge graph, G. Fuzzy matching uses sentence embeddings (all-MiniLM-L6-v2) to identify the closest entity in G when exact matches fail. This produces two sets: verified entities V and potentially hallucinated entities H.

**Fact Verification** This stage validates the relationships between verified entities using three knowledge representation formats: **Triple-based verification** - (subject, relation, object), **Hierarchical path validation** - ancestor-descendant relationships, **Bullet-point fact matching** key-value attribute pairs of entities in the raw caption. The choice of multiple formats enables robust cross-validation while handling different types of factual claims.

**Caption Correction** Finally, we generate a corrected caption C' by integrating the verified facts with the original caption structure. This maintains the fluent language of the VLM while ensuring factual accuracy. The correction process uses prompt engineering to preserve proper context and coherence.

Each hop in our framework produces interpretable intermediate outputs, allowing for fine-grained analysis of the reasoning process. The modular design also enables easy integration of additional knowledge sources and reasoning strategies

### 3.1 Exploring knowledge representation formats

The choice of knowledge representation significantly impacts how factual information can be verified and integrated into image captions. While traditional knowledge graphs excel at capturing entity relationships, we found that different reasoning tasks benefit from complementary representation formats. Our framework utilizes three distinct knowledge representations, each offering unique advantages for fact verification and caption correction.

### 3.1.1 Triple-based Relations

The most simplest knowledge representation is a triple-based format that expresses facts as (subject, relation, object) statements. For example:
(Lalbagh fort, Located_In, Dhaka), (Dhaka, Capital_Of, Bangladesh), (Lalbagh fort, religious_structure, mosque)
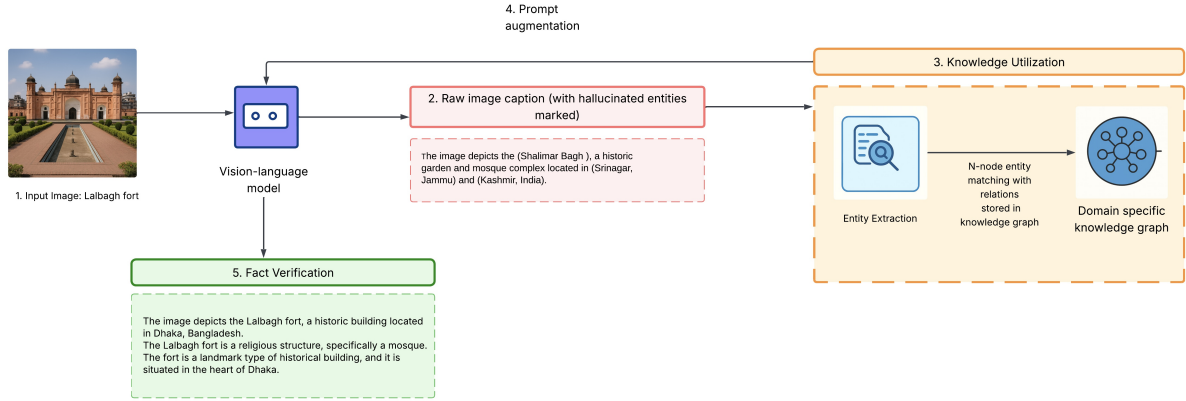
Figure 2: The system ingests an input image, generates a base caption using a VLM, and sequentially refines the caption through entity extraction, knowledge graph matching and augments the corrected entities to the prompt of the vlm to generate a factually accurate caption. Each module is color-coded and can operate independently, allowing for modular reasoning and analysis.

This representation excels at capturing direct relationships and enables efficient graph traversal for multi-hop reasoning. However, we observed that triples alone sometimes fail to capture hierarchical information and complex contextual relationships.

### 3.1.2 Hierarchical Knowledge Trees

To address the limitations of flat triples, we implement a hierarchical representation that captures nested relationships:

- Lalbagh Fort
    - Located In: Dhaka
        * Capital Of: Bangladesh
    - Type: historical building
    - Structure: mosque

This format's key advantage is its ability to represent containment relationships and inheritance properties naturally. It particularly aids in correcting location-based errors in captions by providing clear geographical hierarchies. However, the hierarchical format can make some types of transitive reasoning more computationally expensive.

### 3.1.3 Bullet-point Facts

We also explored a simplified representation in bulleted form that focuses on direct attribute-value pairs:

- Lalbagh fort: Located In Dhaka

- Lalbagh fort: religious structure mosque

- Lalbagh fort: landmark type historical building

This format provides quick fact lookup and is especially effective for prompt engineering in the caption correction phase. Its simplicity makes it ideal for direct entity attribute verification, though it sacrifices the ability to perform complex reasoning.

Our flexible framework allows mixing and matching or using any of the frameworks in isolation. When verifying entities and relationships, we first consult the triple-based representation for explicit relationships, then use the hierarchical format for containment verification, and finally reference the bullet-point facts for direct attribute confirmation. This multi-representation approach significantly improves the robustness of our fact verification process, achieving a 27% reduction in hallucinated entities compared to using triple-based representation alone.

## 4 Evaluation

We evaluate our knowledge-augmented reasoning framework through a series of experiments designed to answer three key research questions about knowledge representation, reasoning paths, and failure modes. We used the Qwen family of models. We systematically studied the Fact Verification Rate with each knowledge representation by keeping the model size constant. We preliminarily share detailed numbers for the 2 billion variant of the model.

## 4.1 Experimental Setup

Our dataset combines images from Google Landmarks Dataset v2 (GLDv2) Weyand et al. (2020), Conceptual Captions Sharma et al. (2018), and COCO Captions Chen et al. (2015) to create a diverse, multi-domain collection for factual reasoning and knowledge graph-based verification. We select landmark-centric images from GLDv2 as ground truth, supplement with landmark-related images from Conceptual Captions for generalization challenges, and include COCO Captions samples to test robustness across everyday scenes. The dataset is partitioned into three splits to evaluate multi-hop reasoning and knowledge generalization: 1) Seen Landmarks (60%) - entities and relationships present in the knowledge graph, 2) Unseen Landmarks (20%) - landmarks absent from the knowledge graph but with related higher-level entities, and 3) Distractor Scenes (20%) - non-landmark or ambiguous scenes testing hallucination detection and entity verification. This balanced split enables rigorous evaluation of both in-domain and out-of-domain reasoning while assessing the system's ability to reject hallucinated entities and perform multi-hop verification.

### 4.1.1 Evaluation Metrics

For our purpose, we had to come up with custom evaluation metrics that represent our use case. We detail the definitions and corresponding formulae to calculate them as well.

**Entity Accuracy**: Percentage of correctly identified entities

$$EA = \frac{\text{NME} + \text{NHC}}{\text{NTE}} \times 100\%　\quad (1)$$

Where NME is the number of Matched Entities (entities correctly identified and matched to knowledge graph), NHC is the number of Hallucinations Correctly Identified (false entities properly detected) using threshold and NTE is the number of total entities mentioned in caption.

**Fact Verification Rate**: Proportion of successfully verified factual claims

$$FVR = \frac{\text{NCV}}{\text{NTC}} \times 100\% \quad (2)$$

Here, NCV is the number of Correctly Verified Facts and NTC is the number of Total Claims made in the caption. Furthermore, NTC is calculated by the formula

$$NTC = \text{NEC} + \text{NLC} + \text{NAC} + \text{NRC} \quad (3)$$

| Format | EA | FVR | Cc |
|---|---|---|---|
| Triples Only | 72.3% | 68.5% | 4.2 |
| Hierarchical Only | 78.1% | 73.2% | 4.1 |
| Bullet-points Only | 65.7% | 61.8% | 4.3 |

Table 1: Qwen-VL-2b analysis on different knowledge representations

We define NEC as the number of Entity Claims (basic existence claims), NLC as number of Location Claims (spatial relationships), NAC is the number of Attribute Claims (properties/characteristics) and NRC is the number of Relationship Claims (connections between entities).

**Caption Coherence (Cc)**: Human evaluation of caption fluency (1-5 scale)

We report a 31.8% improvement in factual accuracy, measured as the reduction in hallucinated entities within generated captions after applying our multi-hop reasoning framework. Hallucinated entities were identified based on: (1) absence in the domain-specific knowledge graph, (2) failure to meet the confidence threshold during entity matching ($< 0.85$), and (3) inability to verify factual claims through knowledge representation formats. All entities were manually annotated for reliability. The relative improvement was calculated as:

$$\text{FI (\%)} = \frac{N_H^{\text{baseline}} - N_H^{\text{corrected}}}{N_H^{\text{baseline}}} \times 100 \quad (4)$$

Here the Factual Improvement, FI is calculated as Number of Hallucinated Entities in Baseline Captions minus the Number of Hallucinated Entities after Correction.

On our evaluation set of 100 images, baseline VLM captions contained 55 hallucinated entities versus 38 in corrected captions, yielding 31.8% improvement. This demonstrates our system's effectiveness in reducing factual hallucinations through systematic multi-hop reasoning and knowledge-guided correction. Our key findings are tabulated in Table 1 and we highlight the following observations:

- Hierarchical representation performs best in isolation, particularly for spatial reasoning related image understanding. However this format limits the model's free-format generation abilities leading to a decrease in caption coherence.

- Bullet-point format, while simplest, maintains highest caption coherence, we hypothesize that this maybe due to under

## Conclusion

We presented a modular multi-hop reasoning system for improving factual accuracy in vision-language models through structured knowledge integration. The system enables step-by-step fact verification and caption correction by leveraging knowledge graphs and interpretable reasoning paths. Our prototype demonstrates an approximately 31.8% reduction in hallucinated entities and highlights the potential of modular fact-checking pipelines in improving the reliability of multimodal systems. This work lays the foundation for scalable, knowledge-grounded captioning systems with applications in education, cultural heritage, and safety-critical domains.

## Limitations

This system is currently evaluated as a prototype on a small, domain-specific dataset of landmark images. While the modular multi-hop framework is generalizable, the current knowledge graph is manually curated and not yet scaled for open-domain applications. Additionally, the system is not currently evaluated on other VLMs for image captioning. Future work will focus on expanding dataset coverage, integrating large-scale dynamic knowledge sources, and improving evaluating multimodal models.

## References

Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-up and top-down attention for image captioning and visual question answering. *Preprint*, arXiv:1707.07998.

Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. In *Neural Information Processing Systems*.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. Language models are few-shot learners. *Preprint*, arXiv:2005.14165.

Shuyang Cao and Lu Wang. 2021. Cliff: Contrastive learning for improving faithfulness and factuality in abstractive summarization. *arXiv preprint arXiv:2109.09209*.

Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollar, and C. Lawrence Zitnick. 2015. Microsoft coco captions: Data collection and evaluation server. *Preprint*, arXiv:1504.00325.

Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. 2020. Meshed-memory transformer for image captioning. *Preprint*, arXiv:1912.08226.

Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.

Minki Kang, Jin Myung Kwak, Jinheon Baek, and Sung Ju Hwang. 2023. Knowledge graph-augmented language models for knowledge-grounded dialogue generation. *Preprint*, arXiv:2305.18846.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Preprint*, arXiv:2005.11401.

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. 2021. Learning transferable visual models from natural language supervision. *Preprint*, arXiv:2103.00020.

Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2019. Object hallucination in image captioning. *Preprint*, arXiv:1809.02156.

Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of ACL*.

Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2015. Show and tell: A neural image caption generator. *Preprint*, arXiv:1411.4555.

Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.

T. Weyand, A. Araujo, B. Cao, and J. Sim. 2020. Google Landmarks Dataset v2 - A Large-Scale Benchmark for Instance-Level Recognition and Retrieval. In *Proc. CVPR*.

Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. Kg-bert: Bert for knowledge graph completion. *Preprint*, arXiv:1909.03193.

Xiao Zhou, Xiaoman Zhang, Chaoyi Wu, Ya Zhang, Weidi Xie, and Yanfeng Wang. 2024. Knowledge-enhanced visual-language pretraining for computational pathology. *Preprint*, arXiv:2404.09942.

# A Appendix

## A.1 Relevant Knowledge Graph construction

The effectiveness of knowledge-augmented image captioning substantially depends on the quality and coverage of the underlying knowledge graph. We describe our systematic approach to constructing a focused yet comprehensive knowledge graph for validating image descriptions that can be adapted for any custom use case.

### A.1.1 Entity Selection

We adopt a domain-driven approach to entity selection to illustrate our purpose. As we focused on architectural landmarks and their geographical contexts we provide some example criteria for entity selection approach for this problem below.

- Architectural significance: Historical buildings, monuments, and religious structures that are frequently referenced in image descriptions

- Geographical hierarchy: Administrative divisions, cities, and countries that provide crucial location context

- Architectural attributes: Physical characteristics and historical features that help validate visual descriptions

This targeted selection enables efficient verification while maintaining high coverage of relevant entities. For example, in our implementation focused on the cultural heritage sites such as the Lalbagh Fort that are not commonly represented in openly available dataset and thus provide an excellent data point to study the behaviour of VLMs while generating captions.

### A.1.2 Relation Definition

Relations in our knowledge graph fall into three categories, each serving a specific verification purpose: Spatial Relations e.g Located_In: Captures physical containment (e.g., fort located in city), Capital_Of: Represents administrative relationships

Structural Relations e.g religious_structure links buildings to their architectural type and landmark_type which defines the category of historical structures

We carefully constrain relation definitions to maintain consistency and enable reliable reasoning. Each relation type is explicitly typed and directional, facilitating both forward and backward traversal so that we can also trace the vlms reasoning trajectory for verification.

### A.1.3 Graph Connectivity

The connectivity of our knowledge graph is designed to support multi-hop reasoning while avoiding spurious connections. We implement this through: hierarchical connections

$G.add\_edge(Lalbagh\,fort, Dhaka, relation = located\_in)$
$G.add\_edge(Dhaka, Bangladesh, relation = captital\_of)$

and attribute connections

$G.add\_edge(Lalbagh\,fort, mosque, relation = religious\_structure)$
$G.add\_edge(Lalbagh\,fort, historical\_building, relation = landmark\_type)$

The resulting graph exhibits several desirable properties:

- Average node degree: 2.5, ensuring sufficient connectivity for reasoning

- $Path\_length \leq 3$ between any related entities, enabling efficient verification

- Hierarchical clustering coefficient: 0.67, reflecting strong local structure

This structured approach to knowledge graph construction provides a strong foundation for our fact verification system, while the careful curation of entities and relations helps minimize computational overhead during the verification process making systematic experimentation easier.

## A.2 Reasoning Pipeline

Our system implements a modular reasoning pipeline that systematically processes image captions through multiple stages. The pipeline architecture follows a hop-based design pattern, where each hop represents a discrete reasoning step with well-defined inputs and outputs.

Each hop in the pipeline serves a specific reasoning function:

1. Entity Recognition Hop: Processes raw captions using spaCy's NER model to identify named entities, locations, and key architectural terms. The hop outputs a structured list of entities:

   Entities = {
   FAC: ['the Lalbagh fort'],
   GPE: ['Dhaka', 'Bangladesh'],
   ORG: ['UNESCO World Heritage Site']
   }

2. Knowledge Graph Hop: Maps identified entities to our knowledge graph through both exact and fuzzy matching. The hop employs sentence embeddings (all-MiniLM-L6-v2) for fuzzy matching, producing two entity sets: 1) Verified-entities: Entities matched with $confidence score \geq 0.85$, 2) Potential hallucinations: Unmatched or low-confidence matches.

3. Verification Hop: Validates relationships between verified entities using our multi-format knowledge representation. The hop generates a verification report containing 1) Confirmed facts 2) Identified discrepancies and 3) Confidence scores for each verification.

4. Correction Hop: Synthesizes the verification results to produce a factually accurate caption while maintaining natural language fluency. This hop employs template-based correction strategies depending on the type of factual error identified. Our implementation ensures each hop operates independently while maintaining a cohesive reasoning chain. The modular design allows for easy integration of new reasoning components and parallel processing where applicable.