

Not All Inconsistency is Equal: Decomposing LVLM Uncertainty into Belief Divergence and Belief Conflict

Jie Shi¹, Xiaodong Yue^{2,3*}, Wei Liu⁴, Yufei Chen⁴, Feifan Dong³

¹ School of Computer Engineering and Science, Shanghai University, Shanghai, China

² Institute of Artificial Intelligence, Shanghai University, Shanghai, China

³ School of Future Technology, Shanghai University, Shanghai, China

⁴ School of Computer Science and Technology, Tongji University, Shanghai, China

jieshi@shu.edu.cn, yswantfly@shu.edu.cn, ldachuan@outlook.com, yufeichen@tongji.edu.cn, dongfeifan@shu.edu.cn

Abstract

Uncertainty Quantification (UQ) is critical for detecting hallucinations in black-box Large Vision-Language Models (LVLMs). However, prevailing methods like Discrete Semantic Entropy (DSE) are unreliable, as their scores are primarily dominated by the number of semantic clusters. This renders them incapable of distinguishing between benign semantic ambiguity (varied but coherent responses) and severe belief conflict (contradictory responses). We address this limitation by proposing a novel framework rooted in Dempster-Shafer theory of evidence, built on the premise that not all inconsistency is equal. Our method decomposes uncertainty into two complementary metrics: Belief Divergence, which quantifies ambiguity by measuring the separation between viewpoints, and Belief Conflict, which captures direct logical contradictions. Extensive experiments demonstrate that our framework provides a more reliable measure of uncertainty.

Introduction

Large Vision-Language Models (LVLMs) have demonstrated powerful capabilities across a spectrum of vision and language tasks, becoming foundational tools in numerous real-world applications (Hu et al. 2024; Peng et al. 2023; Cui et al. 2024). A significant portion of these cutting-edge models are deployed as black-box systems, accessible only through APIs without exposing internal parameters. While this broadens access, it presents a critical challenge: ensuring the reliability of models whose inner workings are opaque (Woo et al. 2025). The phenomenon of *hallucination*, defined as the generation of plausible but factually incorrect information, is particularly concerning in this context (Bai et al. 2024; Liu et al. 2024b). Therefore, developing robust methods for Uncertainty Quantification (UQ) for these black-box models is a crucial necessity for responsible AI deployment.

A prevailing strategy for black-box UQ involves generating a diverse answer set via semantic-preserving perturbations and quantifying its *inconsistency* as a proxy for uncertainty (Zhang, Zhang, and Zheng 2024). This is often realized as Discrete Semantic Entropy (DSE) (Farquhar et al. 2024), which clusters responses by meaning,

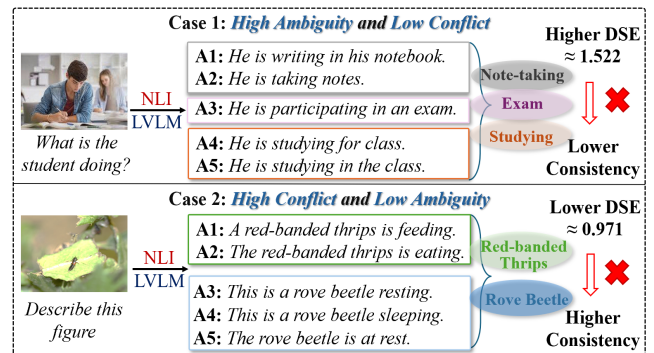


Figure 1: An illustration of the two primary failure modes of Discrete Semantic Entropy (DSE). In the high-ambiguity case (Top), DSE yields a high score (≈ 1.522), incorrectly suggesting high uncertainty. In the high-conflict case (Bottom), it yields a deceptively lower score (≈ 0.971), failing to capture the severe contradiction.

typically using bidirectional entailment predictions from a Natural Language Inference (NLI) model, before calculating entropy from the cluster distribution. However, **as its evaluation is primarily dominated by the number of semantic clusters**, this approach is incapable of distinguishing between two fundamentally different types of inconsistency: a mild form arising from *semantic ambiguity* and a severe form arising from *semantic conflict*.

This dual-limitation is vividly illustrated in Figure 1. The benign ambiguity in Case 1 results in a high DSE score (≈ 1.522), while the severe conflict in Case 2 yields a deceptively lower score (≈ 0.971). These counter-intuitive results reveal a fundamental principle that existing metrics overlook: **not all inconsistency is created equal**. A mild form like ambiguity signals a stable exploration of a topic, whereas a severe form like conflict signals a critical failure in reasoning.

To achieve this differentiation, we propose a novel UQ framework rooted in Dempster-Shafer theory (DST) of evidence (Shafer 1992), chosen for its inherent ability to explicitly model ignorance. Crucially, our framework operates in a fully black-box setting. It moves beyond a single, coarse-grained score by decomposing uncertainty

*Corresponding author

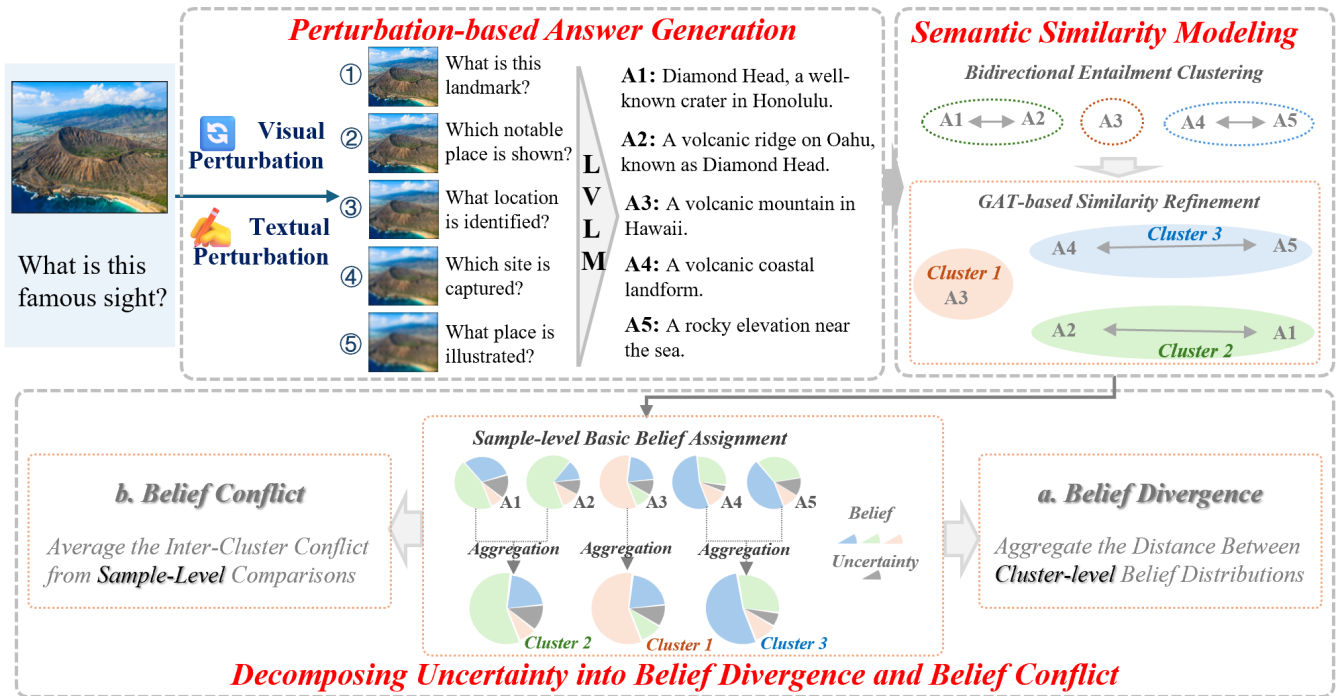


Figure 2: An overview of our proposed framework for decomposing LVLM uncertainty. (1) First, we generate a diverse answer set by applying semantic-preserving perturbations to the initial vision-language input. (2) Next, the answer set’s semantic structure is modeled through NLI-based hard clustering and then refined with a Graph Attention Network (GAT). (3) Finally, within our DST-based module, each answer is converted into a Basic Belief Assignment (BBA), which is used to quantify uncertainty via two complementary metrics: Belief Divergence and Belief Conflict.

into two complementary components: *Belief Divergence*, which quantifies semantic ambiguity, and *Belief Conflict*, which captures direct logical contradictions. By integrating these two dimensions, our framework provides a more robust, fine-grained, and interpretable measure of LVLM uncertainty. In summary, our contributions are as follows:

- We are the first to conceptualize and decompose LVLM inconsistency into two distinct forms: benign semantic ambiguity and critical semantic conflict.
- We develop a novel, black-box framework using Dempster-Shafer evidence theory to quantify uncertainty via two complementary metrics: Belief Divergence and Belief Conflict.
- Empirical analyses validate the effectiveness of the proposed method in providing a more reliable measure of uncertainty.

Related Work

A prevailing strategy in UQ of black-box LVLMs is to generate a diverse answer set via semantic-preserving perturbations and then quantify its inconsistency as a proxy for uncertainty. Existing methods primarily operate on these answer sets (Kadavath et al. 2022; Zhang et al. 2024). A foundational approach is to first perform *Semantic Clustering*, typically using a pre-trained NLI model to group semantically equivalent answers. Methods like DSE (Zhang,

Zhang, and Zheng 2024) and NumSet (Kuhn, Gal, and Farquhar 2023; Lin, Trivedi, and Sun 2024) then quantify uncertainty based on the distribution or simply the number of these clusters. However, the primary limitation of these methods is that their scores are dominated by the cluster count, ignoring the finer-grained relationships between them. To address this, other methods incorporate pairwise similarity. These include lexical-based approaches like LexSim (Fomicheva et al. 2020), nearest-neighbor inspired techniques like SNNE (Nguyen, Payani, and Mirzasoleiman 2025), and a family of graph-based methods that derive uncertainty from the properties of a semantic similarity graph, such as Deg, SumEigv, and Eccen, proposed by (Lin, Trivedi, and Sun 2024). While these approaches provide a more nuanced view than simple cluster counting, they still treat inconsistency as a monolithic concept, lacking a formal mechanism to distinguish between benign semantic ambiguity and severe semantic conflict. In contrast, Evidence Theory, also known as Dempster-Shafer Theory (DST) (Shafer 1992), offers a dedicated framework for reasoning under uncertainty. It has been successfully applied to diverse tasks, including trustworthy multi-view learning (Han et al. 2022; Liu, Chen, and Yue 2025a,b; Liang et al. 2025a,b; Xu et al. 2024; Liu et al. 2023c, 2022), ensemble learning (Fu et al. 2022; Lv et al. 2021), fine-grained image classification (Xu et al. 2023), patent classification (Zhang et al. 2022), and medical diagnosis (Liu et al.

Method

In this section, we introduce a novel framework rooted in DST that decomposes inconsistency into its complementary parts. The framework unfolds in three stages, as illustrated in Figure 2.

Modeling Answer Set Semantics

Perturbation-based Answer Generation: To assess a model’s consistency, we first generate a diverse answer set $\mathcal{A} = \{a_1, a_2, \dots, a_N\}$ by querying the LVLM with N semantically equivalent but perturbed input pairs, a strategy inspired by prior work (Zhang, Zhang, and Zheng 2024).

For the visual modality, we perturb the original image I into a series of images I_i using a 2D Gaussian Blur, controlling the intensity via the blur radius r_i . Concurrently, for the textual modality, we employ a pre-trained Large Language Model (LLM) to generate semantically equivalent paraphrases T_i of the original text T , controlling diversity by the temperature parameter τ_i . These perturbations are synchronized by intensity to form a set of input pairs, $\Lambda = \{\langle I_i, T_i \rangle\}$, which are then fed to the LVLM to produce the answer set \mathcal{A} for our analysis.

Structured Semantic Modeling: To model the internal semantic structure of the answer set \mathcal{A} , we perform an initial hard clustering followed by a graph-based refinement. First, we use a pre-trained NLI model to apply a strict bidirectional entailment test to every answer pair, grouping semantically equivalent answers into a set of N_c clusters, $\{C_k\}_{k=1}^{N_c}$.

While this provides a coarse overview, we refine it using a Graph Attention Network (GAT (Velickovic et al. 2017)) to create a more nuanced similarity measure. We first obtain initial embeddings, \mathbf{X} , for all answers using a standard sentence-embedding model. We then encode the clustering result into a graph where edges connect answers belonging to the same cluster. This graph and the embeddings \mathbf{X} are processed by a pre-trained GAT to produce refined, “structure-aware” embeddings. From these, we compute the final similarity matrix \mathbf{S} using a Gaussian kernel function, where each element S_{ij} captures the deep semantic consistency between answers a_i and a_j .

Belief Representation of Answers

After modeling the semantic structure, we move to the formal framework of DST, chosen for its inherent ability to explicitly represent ignorance. Each answer a_i is represented as a Basic Belief Assignment (BBA), denoted m_i .

Our frame of discernment (FOD), $\Theta = \{C_1, \dots, C_{N_c}\}$, is the set of mutually exclusive semantic clusters. The focal sets for our BBA, m_i , are restricted to the *singleton sets* $\{C_k\}$ and the *universal set* Θ . The belief masses are calculated as follows. First, we calculate a semantic affinity distribution $\mathbf{p}^i = [p_1^i, \dots, p_{N_c}^i]$ for answer a_i across all clusters based on the refined similarity matrix \mathbf{S} :

$$p_k^i = \tilde{S}_i(C_k) / \sum_{l=1}^{N_c} \tilde{S}_i(C_l), \quad (1)$$

where

$$\tilde{S}_i(C_k) = \frac{1}{|C_k|} \sum_{a_j \in C_k} S_{ij} \quad (2)$$

is the support score for each answer a_i with respect to each cluster C_k by averaging its similarity to the cluster’s members. The ignorance term $m_i(\Theta)$ is then quantified using the normalized Shannon entropy of this distribution:

$$m_i(\Theta) = \frac{-\sum_{k=1}^{N_c} p_k^i \log_2 p_k^i}{\log_2 N_c}. \quad (3)$$

The remaining belief mass is then distributed proportionally among the singleton sets:

$$m_i(\{C_k\}) = (1 - m_i(\Theta)) \cdot p_k^i. \quad (4)$$

Quantifying Consistency via Two Metrics

The BBA representation allows us to decompose inconsistency into two complementary components.

Belief Divergence (Measuring Ambiguity): This metric quantifies inconsistency from *semantic ambiguity*. A higher divergence corresponds to a higher degree of ambiguity.

To calculate this, we first establish a centroid belief distribution $\bar{m}_k = \frac{\sum_{i \in C_k} m_i}{|C_k|}$ for each semantic cluster by averaging the BBAs of its members. Our raw inconsistency score for divergence, d , is then the total sum of the pairwise Jous-selme distances between every unique cluster centroid:

$$d = \sum_{k=1}^{N_c-1} \sum_{l=k+1}^{N_c} d_J(\bar{m}_k, \bar{m}_l), \quad (5)$$

where

$$d_J(\bar{m}_k, \bar{m}_l) = \sqrt{\frac{1}{2}(\bar{m}_k - \bar{m}_l)^T \mathbf{D}(\bar{m}_k - \bar{m}_l)}. \quad (6)$$

Here, \mathbf{D} is the matrix of pairwise Jaccard similarities between focal sets. This summation-based design, rather than an average, ensures that, *similar to DSE, the metric is sensitive to the number of distinct semantic clusters, while the Jous-selme distance provides a nuanced measure of their actual separation*. This total distance d is then mapped to a normalized consistency score, $S_{div} = \exp(-d)$.

Belief Conflict (Measuring Contradiction): This metric quantifies inconsistency from *semantic conflict*. To obtain a global measure of conflict, we employ a hierarchical averaging strategy. The process begins at the sample-level by calculating the Dempster conflict coefficient, κ , between any two answers, a_i and a_j , from different clusters. Since the focal sets in our framework are restricted to singletons, we have a simplified form:

$$\kappa(m_i, m_j) = \sum_{k \neq l} m_i(\{C_k\}) m_j(\{C_l\}). \quad (7)$$

Next, these sample-level values are averaged to find the representative conflict between any two clusters, C_k and C_l :

$$\kappa(C_k, C_l) = \frac{1}{|C_k||C_l|} \sum_{a_i \in C_k} \sum_{a_j \in C_l} \kappa(m_i, m_j). \quad (8)$$

Finally, our raw inconsistency score for conflict, $\bar{\kappa}$, is:

$$\bar{\kappa} = \frac{1}{\binom{N_c}{2}} \sum_{k=1}^{N_c-1} \sum_{l=k+1}^{N_c} \kappa(C_k, C_l). \quad (9)$$

This conflict level is then mapped to a normalized consistency score, $S_{con} = \exp(-\bar{\kappa})$.

Final Integrated Consistency Score: While the diagnostic power lies in the individual metrics, a single score is necessary for benchmarking. We combine the two consistency scores into a final score, $S_{overall}$, using a weighted average:

$$S_{overall} = \alpha S_{div} + (1 - \alpha) S_{con}, \quad (10)$$

where $\alpha \in [0, 1]$ is a weighting parameter. This score represents the model’s overall consistency and is inversely proportional to its uncertainty.

Complexity Analysis

The overall time complexity of our framework is $O(N^2)$, where N is the number of generated answers. This complexity is comparable to standard methods like DSE-based method ($O(N^2)$) and is primarily dominated by the pairwise semantic clustering step. The cost of each major module is as follows. 1) The initial **Answer Generation** module has a complexity of $O(N)$, as it requires N forward passes through the LVLm. The **Structured Semantic Modeling** Module contains the main computational bottleneck with the NLI-based hard clustering requiring $O(N^2)$ pairwise inferences. The subsequent GAT refinement is an efficient, single forward pass over the constructed graph. Finally, the **Belief Framework Quantification** stage, including the BBA construction and the calculation of our two metrics, also has a complexity of approximately $O(N^2)$, as its operations depend on pairwise comparisons between the N answers and N_c clusters (where $N_c \leq N$).

Experiments

Experimental Setup

Datasets: We evaluate our method on three diverse benchmarks with various LVLms. First, **LLaVA-Bench** (Liu et al. 2023a) provides 60 challenging questions designed to test convention, detail, and complex reasoning. Second, **MM-Vet-V2** (Yu et al. 2024), a comprehensive benchmark with 218 questions, systematically evaluates 6 core and 16 combined capabilities, from OCR to chart understanding. Finally, **VQA-RAD** (Lau et al. 2018), a medical visual question answering dataset with over 3500 questions on 315 radiological images, tests domain-specific abilities like identifying abnormalities.

Evaluated LVLms: We conduct experiments on 13 LVLms from four distinct model families: **Qwen** (Wang et al. 2024; Bai et al. 2025), **InternVL** (Chen et al. 2024; Zhu et al. 2025), **LLaVA** (Liu et al. 2023b, 2024a), and **DeepSeek** (Lu et al. 2024). While many of these models are open-source, we treat all of them as **black-box** in our experiments, as our proposed method only requires access to their final text outputs.

Implementation Details: For each sample, we first generate a reference answer at a low temperature (0.1) for correctness evaluation. To create the answer set for our uncertainty analysis, we generate N perturbed inputs. **Unless otherwise specified, we set $N = 10$ and the weighting parameter $\alpha = 0.5$ for all main experiments.** Visual perturbations are 2D Gaussian Blurs with radii r_i varying from 0.1 to 0.55 (step 0.05), while textual perturbations are paraphrases from the **Qwen2.5-7B-Instruct** model (Team 2024) with temperatures τ_i in the same range and step size. The LVLm under evaluation then generates an answer for each perturbed input. We deliberately use a low temperature (0.1) for this step, as we empirically found that for our VQA tasks, the input perturbations alone were sufficient to elicit a rich and semantically diverse answer set. Moreover, we use the **jina-embeddings-v3** (Sturua et al. 2024) to obtain initial answer embeddings. The same **Qwen2.5-7B-Instruct** model also serves as the NLI model for the hard clustering step. The GAT-based refinement is implemented using the *GATConv* layer from the PyTorch Geometric library with a single-layer network. To ensure a fair comparison, all baseline methods were evaluated under the exact same experimental settings as our proposed framework. All experiments were conducted on a single NVIDIA A100 40GB GPU.

Evaluation metrics: Following prior work (Nguyen, Payani, and Mirzasoleiman 2025), we adopt two standard metrics: **AUROC**, which measures the ability of an uncertainty score to discriminate between correct and incorrect answers, and **AUARC**, which reflects the practical utility of using the score to reject low-quality answers.

Baselines: We compare our method against eight existing black-box UQ baselines. These include methods based on lexical similarity (**LexSim** (Fomicheva et al. 2020)), semantic clustering and entropy (**DSE** (Zhang, Zhang, and Zheng 2024), **NumSet** (Kuhn, Gal, and Farquhar 2023)), three graph-based methods proposed by (Lin, Trivedi, and Sun 2024) (**SumEigv**, **Deg**, **Eccen**), and other recent approaches (**LUQ** (Zhang et al. 2024), **SNNE** (Nguyen, Payani, and Mirzasoleiman 2025)). For SNNE, we use the variant with the ROUGE similarity function (**SNNE-ROUGE**), as it was demonstrated to be the most effective in the original work.

Experimental Results

Comparison with Baselines: We evaluate our proposed framework against a suite of strong black-box UQ baselines, with the AUROC and AUARC results presented in Table 1 and Table 2, respectively. The results demonstrate the robust performance of our framework, which achieves state-of-the-art results on the vast majority of configurations. Our framework’s advantage is particularly pronounced on the more complex and nuanced benchmarks, MM-Vet-V2 and VQA-RAD. On these datasets, our method achieves the highest AUROC and AUARC scores for nearly every model, highlighting its robustness in handling challenging questions that require deep comprehension. On the simpler LLaVA-Bench dataset, our method remains highly competitive, though other specialized methods also perform well, particularly on the AUARC benchmark. We attribute this to

GAT	S_{con}	S_{div}	Qwen2-VL 7B	Qwen2.5-VL 7B 32B 72B			LLaVA-1.5 7B 13B		LLaVA-NeXT 7B 13B		InternVL2 8B	InternVL3 8B 9B 14B			DeepSeek-VL 7B
LLaVA-Bench															
	✓		0.6899	0.5989	0.6261	0.6004	0.5839	0.5916	0.6504	0.6875	0.6382	0.6574	0.6358	0.6748	0.6438
		✓	0.6618	0.5897	0.6250	0.6015	0.6557	0.5885	0.6370	0.6704	0.6224	0.6574	0.6358	0.6655	0.6284
✓	✓	✓	0.6716	0.6034	0.6261	0.5993	0.6340	0.5978	0.6370	0.6804	0.6276	0.6609	0.6375	0.6690	0.6322
✓	✓		0.6752	0.6000	0.6217	0.6093	0.5904	0.6294	0.6711	0.6558	0.6422	0.6605	0.6817	0.6840	0.6476
✓	✓	✓	0.7069	0.5989	0.6261	0.6004	0.6840	0.6397	0.6607	0.6804	0.6580	0.6609	0.6473	0.6735	0.6463
✓	✓	✓	0.7106	0.6160	0.6272	0.6049	0.6449	0.6553	0.6622	0.6989	0.6620	0.6620	0.6604	0.6736	0.6605
MM-Vet-V2															
	✓		0.6621	0.6947	0.6798	0.6394	0.7430	0.7149	0.6872	0.6622	0.6859	0.6891	0.6863	0.6669	0.6558
		✓	0.6517	0.6984	0.6805	0.6414	0.7391	0.6961	0.6958	0.6607	0.6816	0.6934	0.6873	0.6743	0.6622
✓	✓	✓	0.6572	0.6956	0.6813	0.6410	0.7412	0.7024	0.6961	0.6623	0.6842	0.6898	0.6876	0.6703	0.6626
✓	✓		0.6397	0.7007	0.6804	0.6426	0.6971	0.6607	0.6591	0.6205	0.6474	0.6879	0.6702	0.6689	0.6261
✓	✓	✓	0.6683	0.6987	0.6842	0.6441	0.7682	0.7365	0.6903	0.6695	0.6942	0.6950	0.6884	0.6731	0.6596
✓	✓	✓	0.6727	0.7019	0.6851	0.6453	0.7703	0.7393	0.7033	0.6753	0.7019	0.6954	0.6940	0.6748	0.6627
VQA-RAD															
	✓		0.6130	0.6469	0.6651	0.6629	0.6403	0.6324	0.5862	0.5974	0.6811	0.6955	0.6910	0.6912	0.5970
		✓	0.5984	0.6424	0.6627	0.6574	0.6267	0.6000	0.5904	0.6118	0.6701	0.6941	0.6888	0.6867	0.5864
✓	✓	✓	0.6047	0.6460	0.6643	0.6613	0.6322	0.6086	0.5886	0.6081	0.6633	0.6960	0.6909	0.6908	0.5912
✓	✓		0.5725	0.6310	0.6461	0.6479	0.6221	0.5922	0.5853	0.5951	0.6286	0.6728	0.6621	0.6797	0.5732
✓	✓	✓	0.6170	0.6505	0.6648	0.6658	0.6446	0.6399	0.5896	0.6160	0.6763	0.6972	0.6916	0.6964	0.5990
✓	✓	✓	0.6171	0.6504	0.6651	0.6661	0.6484	0.6450	0.5943	0.6171	0.6756	0.6986	0.6946	0.6986	0.5991

Table 3: Ablation study of our framework’s components, reported in AUROC across all datasets. Bold text indicates the best-performing configuration for each LVLm. Experimental settings are identical to those in Table 1.

the fact that the answer sets generated from LLaVA-Bench’s simpler questions are less semantically diverse, making the full power of our decomposition less critical. **Most importantly, the primary contribution of our framework is not just the superior performance of the final integrated score, but its diagnostic power.** While the overall score is necessary for benchmarking, the true value lies in the decomposition of uncertainty. Our framework, by providing two distinct metrics for ambiguity and contradiction, offers a level of interpretability that single-score methods cannot.

Ablation Study: We conduct an ablation study to validate the contribution of each component, with AUROC results shown in Table 3. The analysis provides strong empirical evidence for our core thesis that decomposing inconsistency is a more effective approach. The full framework, which integrates S_{div} and S_{con} and utilizes our GAT-based refinement, almost consistently achieves the best performance.

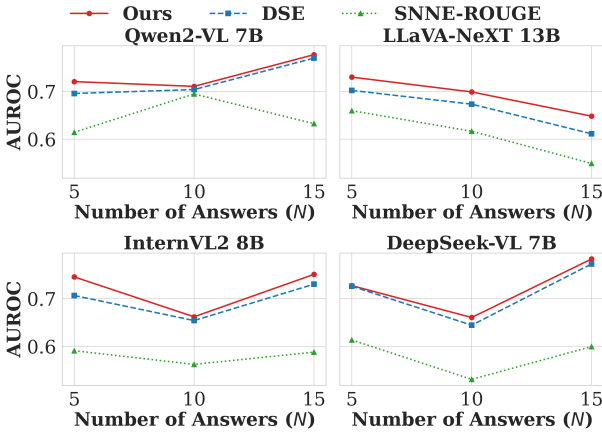


Figure 3: The effect of N on AUROC for our method and key baselines on the LLaVA-Bench dataset.

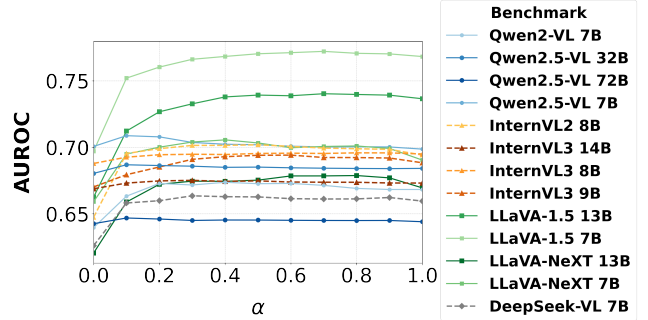


Figure 4: The effect of α on AUROC for all 13 LVLms on the MM-Vet-V2 dataset, using $N = 10$ answers.

Parameter Analysis of N : Here, we analyze the sensitivity of our framework to the number of generated answers on the LLaVA-Bench dataset across $N \in \{5, 10, 15\}$. We select one model from each of the four families under evaluation (Qwen2-VL, LLaVA-NeXT, InternVL2, and DeepSeek-VL) and compare our method against two key baselines, DSE and SNNE-ROUGE.

The results, shown in Figure 3, reveal a complex and non-monotonic relationship between performance and sample size. Our method’s AUROC on several models shows a notable dip at $N = 10$ before surging to its highest point at $N = 15$. This suggests an intricate interplay between the diversity of the answer set and our decomposition-based approach. Despite the strong performance at $N = 15$ on this dataset, we selected $N = 10$ as the default setting for our main experiments to maintain a pragmatic balance between performance and the significant computational cost required across all models and datasets.

Parameter Analysis of α : To analyze the relative contributions of two metrics, we vary the weighting parameter α from 0 (purely Conflict-based) to 1 (purely Divergence-

Case 1: Semantic Conflict		Case 1	Case 2
A1: Someone is slicing an apple over the cutting board. G1		NumSet	2 = 2
A2: Someone is cutting an apple over the cutting board.		LexSim	-0.7746 < -0.6781
A3: Someone is chopping an apple over the cutting board.		SumEigv	2.0662 < 2.2422
A4: Someone is rinsing a banana. G2		Deg	2.4569 < 2.5428
A5: Someone is cleaning a banana.		Eccen	1.0000 < 1.4142
Case 2: Semantic Ambiguity		LUQ	0.0009 < 0.0030
A1: Someone is chopping an apple into chunk. G1		DSE	0.9710 = 0.9710
A2: Someone is slicing an apple into bits.		SNNE-ROUGE	-2.2721 < -2.1975
A3: Someone is cutting up an apple into pieces.		Ours ($1 - S_{overall}$)	0.4847 > 0.4042
A4: Someone is processing a piece of red fruit into chunks on a cutting slab. G2		Uncertainty: Case 1 > Case 2 ✓	
A5: Someone is handling a red fruit, transforming it into tiny bits on a clean cutting board.		G1: Cut apple G1: Cut apple G2: Wash banana G2: Process fruit	
(a) Semantic Ambiguity vs Semantic Conflict			
Case 1: Semantic Low Ambiguity		Case 1	Case 2
A1: A person is running through the park. G1		NumSet	2 = 2
A2: A person is running in a park.		LexSim	-0.6212 > -0.6727
A3: A person is on a run within the park.		SumEigv	1.3827 > 1.3526
A4: A form of outdoor exercise that this figure is doing is running. G2		Deg	1.3323 > 1.2813
A5: A person is engaged in an outdoor run.		Eccen	1.0000 = 1.0000
Case 2: Semantic High Ambiguity		LUQ	0.0003 > 0.0002
A1: A person is running through the park. G1		DSE	0.9710 = 0.9710
A2: A person is running in a park.		SNNE-ROUGE	-2.2000 > -2.2350
A3: A person is on a run in the park.		Ours ($1 - S_{div}$)	0.3850 < 0.4810
A4: A person is getting some outdoor exercise. G2		Uncertainty: Case 1 < Case 2 ✓	
A5: A person is doing a physical activity outdoors.		G1: Running in Park G2: Exercising G2: Running Outdoors Outdoors	
(b) Semantic Low Ambiguity vs Semantic High Ambiguity			
Case 1: Semantic Low Conflict		Case 1	Case 2
A1: A person is folding clothes into a suitcase. G1		NumSet	2 = 2
A2: The individual is arranging clothes into a wardrobe.		LexSim	-0.6735 > -0.6916
A3: Someone is packing their attire into a suitcase.		SumEigv	4.1608 > 3.3368
A4: The figure is choosing an outfit from the wardrobe. G2		Deg	3.7537 > 3.4040
A5: The person is selecting clothes from the suitcase.		Eccen	2.0000 > 1.7320
Case 2: Semantic High Conflict		LUQ	0.5650 > 0.1986
A1: A person is folding clothes into a suitcase. G1		DSE	0.9710 = 0.9710
A2: The individual is arranging clothes into a wardrobe.		SNNE-ROUGE	-0.2126 > -2.1615
A3: Someone is packing their attire into a suitcase.		Ours ($1 - S_{con}$)	0.3470 < 0.4390
A4: The figure is making the bed. G2		Uncertainty: Case 1 < Case 2 ✓	
A5: The person is making the bed.		G1: Pack Clothes G1: Pack Clothes G2: Select Clothes G2: Make the Bed	
(c) Semantic Low Conflict vs Semantic High Conflict			

Figure 5: Illustrative experiments validating our decomposition-based approach. (a) A targeted experiment demonstrating that baseline methods fail to distinguish between semantic ambiguity and conflict, proving the necessity of decomposition. (b) Validation of our Belief Divergence metric ($1 - S_{div}$) on a controlled ambiguity gradient. (c) Validation of our Belief Conflict metric ($1 - S_{con}$) on a controlled conflict gradient.

based) on MM-Vet-V2 benchmark. The result, shown in Figure 4, provides strong empirical evidence for our core thesis that a decomposed approach is superior. For the vast majority of models, performance peaks when $0 < \alpha < 1$, demonstrating a clear synergistic effect where the combination of both metrics consistently outperforms using either in isolation. Interestingly, we also observe that Belief Divergence ($\alpha = 1$) often serves as a stronger individual baseline than Belief Conflict ($\alpha = 0$), suggesting our perturbation strat-

egy primarily elicits semantic ambiguity. Given the robust performance of a balanced combination, we select $\alpha = 0.5$ as a general-purpose setting for all our main experiments.

Discussion

The Necessity of Decomposing Uncertainty

To demonstrate why decomposing inconsistency is essential, we constructed two targeted scenarios shown in Figure 5(a): a Semantic Conflict case containing a direct factual contradiction (e.g., “Cut apple” vs. “Wash banana”), and a Semantic Ambiguity case containing a benign difference in descriptive granularity (e.g., “Cut apple” vs. “Process fruit”). The results of this illustrative experiment are stark: **every baseline method incorrectly assigns a higher uncertainty score to the benign ambiguity case than to the severe conflict case.** This proves their fundamental inability to distinguish between different types of inconsistency and validates our core thesis that *not all inconsistency is equal*. In stark contrast, our framework’s uncertainty score ($1 - S_{overall}$), even with a balanced $\alpha = 0.5$, correctly identifies the semantic conflict as the more severe and uncertain scenario. This is possible only because our Belief Conflict and Belief Divergence metrics can assess these two phenomena independently.

Independent Validation of Decomposed Metrics

To validate that our two metrics independently measure the phenomena they are designed for, we constructed two pairs of targeted, illustrative cases. First, to test the ability of our Belief Divergence metric to measure semantic ambiguity, we created a Low Ambiguity case and a High Ambiguity case. As shown in Figure 5(b), our uncertainty score ($1 - S_{div}$) correctly and intuitively assigns a lower score to the Low Ambiguity case (0.3850) than to the High Ambiguity case (0.4810), demonstrating its sensitivity to the degree of semantic separation. Second, to test the ability of our Belief Conflict metric to measure semantic conflict, we created a Low Conflict case and a High Conflict case with mutually exclusive answers. As shown in Figure 5(c), our uncertainty score ($1 - S_{con}$) behaves as expected, assigning a significantly lower score to the Low Conflict case (0.3470) compared to the High Conflict case (0.4390), confirming its ability to effectively detect direct logical contradictions.

Conclusion

In this work, we address a fundamental limitation of existing black-box UQ methods for LLMs: their inability to distinguish between benign semantic ambiguity and severe semantic conflict. We propose a novel framework, rooted in DST, that operationalizes the thesis that *not all inconsistency is equal*. Our approach successfully decomposes uncertainty into two distinct and complementary metrics: *Belief Divergence*, to quantify ambiguity, and *Belief Conflict*, to capture direct logical contradictions. Extensive experiments, including targeted illustrative cases and evaluations on 13 LLMs across three diverse benchmarks, demonstrate that our framework provides a more robust, fine-grained, and reliable measure of uncertainty than existing methods.

Acknowledgments

This work was supported by the National Natural Science Foundation of China (Nos. 62476165, 62472315, 62406182), and Fundamental Research Program of Shanxi Province (Serial No. 202403021212176), and Science and Technology Innovation Project for Higher Education Institutions of Shanxi Province (Serial No. 2024L004).

References

- Bai, S.; Chen, K.; Liu, X.; Wang, J.; Ge, W.; Song, S.; Dang, K.; Wang, P.; Wang, S.; Tang, J.; et al. 2025. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*.
- Bai, Z.; Wang, P.; Xiao, T.; He, T.; Han, Z.; Zhang, Z.; and Shou, M. Z. 2024. Hallucination of multimodal large language models: A survey. *arXiv preprint arXiv:2404.18930*.
- Chen, Z.; Wu, J.; Wang, W.; Su, W.; Chen, G.; Xing, S.; Zhong, M.; Zhang, Q.; Zhu, X.; Lu, L.; et al. 2024. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 24185–24198.
- Cui, C.; Ma, Y.; Cao, X.; Ye, W.; Zhou, Y.; Liang, K.; Chen, J.; Lu, J.; Yang, Z.; Liao, K.-D.; et al. 2024. A survey on multimodal large language models for autonomous driving. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 958–979.
- Farquhar, S.; Kossen, J.; Kuhn, L.; and Gal, Y. 2024. Detecting hallucinations in large language models using semantic entropy. *Nature*, 630(8017): 625–630.
- Fomicheva, M.; Sun, S.; Yankovskaya, L.; Blain, F.; Guzmán, F.; Fishel, M.; Aletras, N.; Chaudhary, V.; and Specia, L. 2020. Unsupervised quality estimation for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8: 539–555.
- Fu, H.; Yue, X.; Liu, W.; and Denooux, T. 2022. Stable clustering ensemble based on evidence theory. In *2022 IEEE International Conference on Image Processing (ICIP)*, 2046–2050. IEEE.
- Fu, W.; Chen, Y.; Liu, W.; Yue, X.; and Ma, C. 2023. Evidence reconciled neural network for out-of-distribution detection in medical images. In *International Conference on Medical Image Computing and Computer-assisted Intervention*, 305–315. Springer.
- Han, Z.; Zhang, C.; Fu, H.; and Zhou, J. T. 2022. Trusted multi-view classification with dynamic evidential fusion. *IEEE transactions on pattern analysis and machine intelligence*, 45(2): 2551–2566.
- Hu, Y.; Li, T.; Lu, Q.; Shao, W.; He, J.; Qiao, Y.; and Luo, P. 2024. Omnimedvqa: A new large-scale comprehensive evaluation benchmark for medical lvlm. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 22170–22183.
- Kadavath, S.; Conerly, T.; Askell, A.; Henighan, T.; Drain, D.; Perez, E.; Schiefer, N.; Hatfield-Dodds, Z.; DasSarma, N.; Tran-Johnson, E.; et al. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*.
- Kuhn, L.; Gal, Y.; and Farquhar, S. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. *arXiv preprint arXiv:2302.09664*.
- Lau, J. J.; Gayen, S.; Ben Abacha, A.; and Demner-Fushman, D. 2018. A dataset of clinically generated visual questions and answers about radiology images. *Scientific data*, 5(1): 1–10.
- Liang, X.; Fu, P.; Qian, Y.; Guo, Q.; and Liu, G. 2025a. Trusted multi-view classification via evolutionary multi-view fusion. In *Proceedings of the International Conference on Learning Representations*, 1–14.
- Liang, X.; Wang, S.; Qian, Y.; Guo, Q.; Du, L.; Jiang, B.; Luo, T.; and Li, F. 2025b. Trusted multi-view classification with expert knowledge constraints. In *Proceedings of the International Conference on Machine Learning*, volume 267, 37409–37426.
- Lin, Z.; Trivedi, S.; and Sun, J. 2024. Generating with confidence: Uncertainty quantification for black-box large language models. *Transactions on Machine Learning Research*.
- Liu, H.; Li, C.; Li, Y.; and Lee, Y. J. 2024a. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 26296–26306.
- Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2023a. Visual instruction tuning. In *Advances in Neural Information Processing Systems*.
- Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2023b. Visual instruction tuning. *Advances in Neural Information Processing Systems*, 36: 34892–34916.
- Liu, H.; Xue, W.; Chen, Y.; Chen, D.; Zhao, X.; Wang, K.; Hou, L.; Li, R.; and Peng, W. 2024b. A survey on hallucination in large vision-language models. *arXiv preprint arXiv:2402.00253*.
- Liu, W.; Chen, Y.; and Yue, X. 2024. Building trust in decision with conformalized multi-view deep classification. In *Proceedings of the ACM International Conference on Multimedia*, 7278–7287.
- Liu, W.; Chen, Y.; and Yue, X. 2025a. Enhancing multi-view classification reliability with adaptive rejection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 18969–18977.
- Liu, W.; Chen, Y.; and Yue, X. 2025b. Enhancing testing-time robustness for trusted multi-view classification in the wild. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, 15508–15517.
- Liu, W.; Chen, Y.; Yue, X.; Zhang, C.; and Xie, S. 2023c. Safe multi-view deep classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 8870–8878.
- Liu, W.; Chen, Y.; Yue, X.; Zhang, C.; and Xie, S. 2025. Enhancing reliability in medical image classification of imperfect views. *IEEE Transactions on Circuits and Systems for Video Technology*.

- Liu, W.; Yue, X.; Chen, Y.; and Denooux, T. 2022. Trusted multi-view deep learning with opinion aggregation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, 7585–7593.
- Lu, H.; Liu, W.; Zhang, B.; Wang, B.; Dong, K.; Liu, B.; Sun, J.; Ren, T.; Li, Z.; Yang, H.; et al. 2024. Deepseek-vl: Towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*.
- Lv, Y.; Zhang, B.; Yue, X.; Xu, Z.; and Liu, W. 2021. Ensemble of adapters for transfer learning based on evidence theory. In *International Conference on Belief Functions*, 66–75. Springer.
- Nguyen, D.; Payani, A.; and Mirzasoleiman, B. 2025. Beyond semantic entropy: Boosting LLM uncertainty quantification with pairwise semantic similarity. *arXiv preprint arXiv:2506.00245*.
- Peng, Z.; Wang, W.; Dong, L.; Hao, Y.; Huang, S.; Ma, S.; and Wei, F. 2023. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*.
- Shafer, G. 1992. Dempster-shafer theory. *Encyclopedia of Artificial Intelligence*, 1: 330–331.
- Sturua, S.; Mohr, I.; Akram, M. K.; Günther, M.; Wang, B.; Krimmel, M.; Wang, F.; Mastrapas, G.; Koukounas, A.; Koukounas, A.; Wang, N.; and Xiao, H. 2024. jina-embeddings-v3: Multilingual embeddings with task LoRA. *arXiv:2409.10173*.
- Team, Q. 2024. Qwen2.5: A Party of Foundation Models.
- Velickovic, P.; Cucurull, G.; Casanova, A.; Romero, A.; Lio, P.; Bengio, Y.; et al. 2017. Graph attention networks. *stat*, 1050(20): 10–48550.
- Wang, P.; Bai, S.; Tan, S.; Wang, S.; Fan, Z.; Bai, J.; Chen, K.; Liu, X.; Wang, J.; Ge, W.; et al. 2024. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Woo, S.; Zhou, K.; Zhou, Y.; Wang, S.; Guan, S.; Ding, H.; and Cheong, L. L. 2025. Black-Box visual prompt engineering for mitigating object hallucination in large vision language models. In *Proceedings of the Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, 529–538.
- Xu, C.; Si, J.; Guan, Z.; Zhao, W.; Wu, Y.; and Gao, X. 2024. Reliable conflictive multi-view learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, 16129–16137.
- Xu, Z.; Yue, X.; Lv, Y.; Liu, W.; and Li, Z. 2023. Trusted fine-grained image classification through hierarchical evidence fusion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 10657–10665.
- Yu, W.; Yang, Z.; Ren, L.; Li, L.; Wang, J.; Lin, K.; Lin, C.-C.; Liu, Z.; Wang, L.; and Wang, X. 2024. Mm-vet v2: A challenging benchmark to evaluate large multimodal models for integrated capabilities. *arXiv preprint arXiv:2408.00765*.
- Zhang, C.; Liu, F.; Basaldella, M.; and Collier, N. 2024. Luq: Long-text uncertainty quantification for llms. *arXiv preprint arXiv:2403.20279*.
- Zhang, L.; Liu, W.; Chen, Y.; and Yue, X. 2022. Reliable multi-view deep patent classification. *Mathematics*, 10(23): 4545.
- Zhang, R.; Zhang, H.; and Zheng, Z. 2024. VI-uncertainty: Detecting hallucination in large vision-language model via uncertainty estimation. *arXiv preprint arXiv:2411.11919*.
- Zhu, J.; Wang, W.; Chen, Z.; Liu, Z.; Ye, S.; Gu, L.; Tian, H.; Duan, Y.; Su, W.; Shao, J.; et al. 2025. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*.