

# Understanding What Vision-Language Models See in Traffic: PixelSHAP for Object-Level Attribution in Autonomous Driving

Anonymous ICCV 2025 DriveX submission

Paper ID XXXX

## Abstract

001 Vision-Language Models (VLMs) are increasingly used in  
002 autonomous driving for scene understanding, hazard de-  
003 tection, and decision-making support. Yet, knowing which  
004 traffic objects these models prioritize is crucial for safety  
005 validation and trust. Existing interpretability methods pro-  
006 vide pixel-level attributions but fail to answer the key ques-  
007 tion: “Which specific objects—vehicles, pedestrians, traffic  
008 signs—influence the model’s driving decisions?”

009 We introduce **PixelSHAP**, a model-agnostic framework  
010 for object-level explainability in Vision-Language Models  
011 applied to traffic scenarios. PixelSHAP extends Shapley-  
012 based attribution to structured visual entities, systemati-  
013 cally quantifying how individual traffic participants influ-  
014 ence a VLM’s reasoning about driving situations. Operat-  
015 ing purely on input-output behavior, our method is compat-  
016 ible with both open-source models (LLaVA, LLaMA-Vision)  
017 and commercial systems (GPT-4V, Gemini) commonly used  
018 in autonomous driving applications.

019 Our approach introduces novel masking strategies in-  
020 cluding Bounding Box with Overlap Avoidance (BBOA) that  
021 address fundamental challenges in traffic scene attribution,  
022 achieving complete object occlusion while minimizing in-  
023 terference with neighboring vehicles or infrastructure. We  
024 evaluate PixelSHAP on traffic scene understanding tasks,  
025 demonstrating its ability to reveal which objects VLMs pri-  
026 oritize for different driving scenarios. Compared to simple  
027 baselines, PixelSHAP provides semantically meaningful at-  
028 tributions that align with human expectations about traffic  
029 safety priorities.

030 Beyond technical contribution, PixelSHAP enables  
031 safety engineers to audit VLM behavior in autonomous  
032 driving contexts, identify potential failure modes, and vali-  
033 date that models focus on safety-critical objects. Our imple-  
034 mentation provides immediate practical value for develop-  
035 ing more transparent and trustworthy autonomous driving  
036 systems.

## 1. Introduction

037 Vision-Language Models (VLMs) are increasingly integral  
038 to autonomous driving systems, supporting scene under-  
039 standing, hazard detection, and driving decision assistance.  
040 As these models move from prototypes to safety-critical de-  
041 ployments, understanding their decision-making is crucial  
042 for ensuring passenger safety and public trust.  
043 Consider a scenario: a VLM analyzes a busy intersec-  
044 tion and outputs: “Pedestrian visible in crosswalk, vehi-  
045 cle should yield.” This triggers braking protocols. Yet,  
046 among multiple pedestrians—on sidewalks, near the cross-  
047 walk, and one crossing—which person influenced the de-  
048 cision? Attribution is essential to validate that the system  
049 responded to the correct participant.  
050

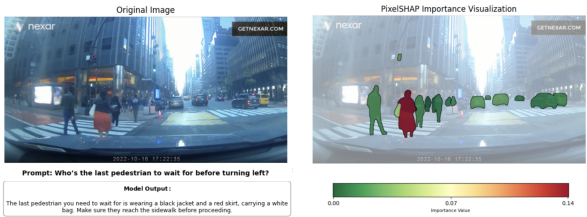


Figure 1. PixelSHAP reveals object-level attribution in traffic scenes. It identifies which pedestrian influenced the VLM’s safety assessment, enabling validation that the model focused on the actual crossing pedestrian.

The core challenge is the semantic gap between how VLMs process visual information and how we interpret their decisions for safety validation. Existing methods fall short: gradient-based approaches like GradCAM require model internals unavailable in commercial VLMs, while pixel-level perturbation methods such as RISE blur distinct traffic participants into indecipherable importance regions.

We propose PixelSHAP, a model-agnostic framework for object-level interpretability in traffic scenarios. Extending Shapley value attribution from tokens to structured visual entities, PixelSHAP quantifies how vehicles, pedestrians, traffic signs, and infrastructure influence VLM assessments.

063	A key innovation is our object-level perturbation approach.	112
064	Our Bounding Box with Overlap Avoidance (BBOA) achieves complete object occlusion while preserv-	113
065	ing context for neighboring elements, enabling clean attrib-	114
066	ution critical for safety validation.	115
067		116
068	This paper makes four contributions to interpretable AI	117
069	for autonomous driving:	118
070	• <b>Traffic-Focused Object-Level Attribution:</b> A frame-	119
071	work identifying traffic participants influencing VLM	120
072	decisions, compatible with open-source and commercial	121
073	models.	122
074	• <b>BBOA Masking Strategy:</b> A perturbation method that	123
075	occludes target objects while preserving surrounding	124
076	context.	
077	• <b>Multi-Model Validation:</b> Evaluation across four	
078	VLMs, comparing against adapted interpretability base-	
079	lines.	
080	• <b>Traffic Scene Evaluation:</b> Protocols for assessing at-	
081	tribution quality in driving-relevant contexts.	
082	The remainder of this paper describes our methodology,	
083	experimental validation, and implications for interpretable	
084	autonomous driving systems.	
085	<b>2. Related Work</b>	
086	Understanding VLM decisions in traffic scenarios requires	
087	explainability methods that can identify which specific vi-	
088	sual objects influence model outputs. The choice of explain-	
089	ability approach is fundamentally constrained by model ac-	
090	cessibility and the semantic granularity required for safety	
091	validation in autonomous driving applications.	
092	<b>2.1. White-Box vs. Black-Box Explainability</b>	
093	Explainability methods for VLMs divide into white-box ap-	
094	proaches requiring access to model internals and black-box	
095	methods operating solely on input-output behavior. White-	
096	box methods like Grad-CAM [14] analyze internal gradi-	
097	ents and activations to generate attribution maps. LVLM-	
098	Interpret [16] provides attention visualization, relevancy	
099	maps, and causal interpretation for vision-language models	
100	by accessing transformer weights and gradients.	
101	White-box methods offer detailed insights into model	
102	mechanisms but face limitations for practical autonomous	
103	driving applications. Many state-of-the-art VLMs de-	
104	ployed in commercial autonomous systems, including GPT-	
105	4V [7] and Gemini-2.0 [1], do not provide access to inter-	
106	nal weights or gradients. For applications requiring inter-	
107	pretability of production-deployed models, black-box ap-	
108	proaches become essential.	
109	<b>2.2. Black-Box Perturbation-Based Methods</b>	
110	Black-box methods explain model decisions through sys-	
111	tematic input perturbation and output analysis, making	
	them compatible with any VLM regardless of architec-	112
	ture. RISE [10] generates importance maps by randomly	113
	masking image regions and measuring output changes.	114
	LIME [13] learns local linear approximations around input	115
	instances using perturbation-based sampling.	116
	These pixel-level approaches face limitations when ana-	117
	lyzing traffic scenes with multiple objects. When vehicles,	118
	pedestrians, and infrastructure appear in proximity, pixel-	119
	-based attribution creates blended importance maps that can-	120
	not isolate individual traffic participants. For autonomous	121
	driving safety validation, understanding which specific ob-	122
	ject influenced a model’s assessment requires object-level	123
	granularity that pixel-based methods cannot provide.	124
	<b>2.3. Shapley Values for Principled Attribution</b>	125
	Shapley values from cooperative game theory [15] provide	126
	mathematically principled feature attribution with desirable	127
	properties including efficiency, symmetry, and additivity.	128
	TokenSHAP [3] demonstrated their effectiveness for lan-	129
	guage model interpretability by quantifying individual to-	130
	ken contributions. MM-SHAP [8] applied Shapley values	131
	to multimodal models, measuring the relative importance	132
	of visual versus textual modalities using image patches.	133
	While Shapley-based approaches offer theoretical rigor,	134
	existing applications focus on different granularities and	135
	questions than object-level attribution in traffic scenarios.	136
	Extending Shapley principles to semantic object-level anal-	137
	ysis while maintaining black-box compatibility remains an	138
	active area of development.	139
	<b>2.4. Multimodal Interpretability: Related Ap-</b>	140
	<b>proaches and Distinctions</b>	141
	Recent interpretability frameworks for VLMs address com-	142
	plementary aspects of multimodal understanding, though	143
	with different focus areas than object-level attribution:	144
	MM-SHAP [8] provides valuable insights into	145
	modality-level contributions, quantifying whether models	146
	rely more on textual or visual information. However, its	147
	patch-based granularity cannot isolate individual traffic	148
	participants within scenes. While MM-SHAP can reveal	149
	that a model used “60% vision, 40% text,” it cannot	150
	distinguish which specific vehicle or pedestrian drove that	151
	visual contribution—a distinction critical for autonomous	152
	driving safety validation.	153
	LVLM-Interpret [16] offers comprehensive analysis	154
	through attention visualization and causal interpretation,	155
	providing detailed insights into model reasoning processes.	156
	However, its dependency on white-box access to attention	157
	weights and gradients limits applicability to commercial	158
	VLMs commonly deployed in autonomous systems. Ad-	159
	ditionally, its patch-based visualizations operate at spatial	160
	resolutions that may not align with semantic object bound-	161
	aries essential for traffic safety analysis.	162

163 These methods address important questions about multi-  
164 modal reasoning and provide valuable debugging capabilities.  
165 Our work complements these approaches by focusing  
166 specifically on the object-level attribution question that existing  
167 methods cannot directly address due to granularity  
168 and accessibility constraints.

## 169 2.5. Object-Level Attribution: Addressing the 170 Granularity Gap

171 Current black-box methods cannot directly answer ques-  
172 tions critical for traffic scene understanding: “Which spe-  
173 cific vehicle influenced the model’s safety assessment?” or  
174 “Did the model focus on the crossing pedestrian or back-  
175 ground elements?” This limitation stems from the granular-  
176 ity mismatch between available attribution methods (pixels  
177 or patches) and the semantic units relevant for autonomous  
178 driving validation (objects representing traffic participants).

179 The autonomous driving context amplifies these chal-  
180 lenges because safety validation requires understand-  
181 ing attribution at the semantic level of traffic partic-  
182 ipants—vehicles, pedestrians, cyclists, and infrastruc-  
183 ture—rather than abstract visual regions. Existing pixel-  
184 level methods cannot distinguish between a pedestrian ac-  
185 tively crossing versus one standing on a sidewalk when both  
186 appear in the same image region, yet this distinction is crit-  
187 ical for validating autonomous driving decisions.

## 188 2.6. Our Approach

189 We introduce PixelSHAP to address the object-level attri-  
190 bution gap by extending Shapley-based attribution to indi-  
191 vidual traffic objects while maintaining black-box compati-  
192 bility with commercial VLMs. Our approach builds on the  
193 theoretical foundation of Shapley values while adapting the  
194 methodology to operate on semantic objects rather than pix-  
195 els or patches.

196 PixelSHAP complements existing interpretability meth-  
197 ods by focusing on the specific granularity and accessibil-  
198 ity requirements of autonomous driving applications. We  
199 demonstrate improvements over adapted versions of exist-  
200 ing methods (RISE-Objects) and simple heuristics, showing  
201 that principled Shapley attribution can provide more accu-  
202 rate object-level explanations for traffic safety validation.  
203 Our evaluation includes comparison with gradient-based  
204 methods where applicable, providing insight into the rela-  
205 tive performance of black-box versus white-box approaches  
206 for object-level attribution tasks.

## 207 3. Problem Statement

208 We formalize object-level attribution in Vision-Language  
209 Models (VLMs) as a black-box interpretability challenge:  
210 quantifying how individual visual objects contribute to a  
211 model’s textual output.

## 3.1. Problem Formulation

212 Given a VLM  $f$  mapping an image  $I$  and optional text  
213 prompt  $p$  to a response  $y = f(I, p)$ , our goal is to as-  
214 sign an attribution score  $\phi_i$  to each object  $o_i$  in  $O =$   
215  $\{o_1, o_2, \dots, o_n\}$ , representing its influence on  $y$ . Attribution  
216 scores must satisfy:  
217

- 218 1. **Efficiency:**  $\sum_{i=1}^n \phi_i = f(I, p) - f(\emptyset, p)$ , where  $\emptyset$  is the  
219 scene with all objects removed.
- 220 2. **Symmetry:** Identical contributors receive equal scores.
- 221 3. **Additivity:** Scores combine consistently across object  
222 subsets.

## 3.2. Key Constraints and Requirements

223 **Black-Box Compatibility:** The method must function  
224 without access to model internals, gradients, or attention  
225 weights, ensuring compatibility with commercial VLMs.  
226 **Object-Level Granularity:** Beyond pixel-level maps, we  
227 require semantic object attribution to answer, e.g., “Which  
228 specific vehicle influenced the decision?” **Semantic Preser-**  
229 **vation:** Perturbations must fully remove an object’s contri-  
230 bution while maintaining scene context.  
231

## 3.3. Applications and Use Cases

232 This formulation supports critical interpretability needs: In  
233 *autonomous systems*, identifying which traffic participants  
234 (vehicles, pedestrians, signs) influenced a VLM’s assess-  
235 ment validates correct prioritization of safety-critical ob-  
236 jects. In *content moderation*, it clarifies which visual ele-  
237 ments trigger policy violations, improving automated sys-  
238 tems. In *medical imaging*, object-level attribution aids in  
239 validating diagnostic outputs and building clinician trust. In  
240 *general scene understanding*, it verifies that VLMs attend  
241 to relevant elements rather than spurious correlations.  
242

## 3.4. Technical Challenges

243 **Object Segmentation Dependency:** Reliable attribution  
244 depends on accurate detection and segmentation of objects.  
245 **Occlusion Strategy:** Removing an object cleanly while  
246 preserving scene context requires sophisticated masking to  
247 avoid artifacts or distortion of neighboring elements. **Com-**  
248 **putational Efficiency:** Exact Shapley value computation is  
249 infeasible; efficient approximations are essential. **Evalu-**  
250 **ation Methodology:** Assessing attribution quality requires  
251 ground truth aligned with human judgments of object im-  
252 portance.  
253

254 The following sections describe how PixelSHAP ad-  
255 dresses these challenges.

## 4. Methodology

256 PixelSHAP extends Shapley value attribution from textual  
257 tokens to visual objects, enabling principled object-level in-  
258 terpretability for Vision-Language Models. Our approach  
259



operates through three stages: object identification with segmentation, systematic perturbation, and attribution computation.

#### 4.1. Framework Design

The framework requires both object detection (bounding boxes) and segmentation masks for each object. Users can integrate results from any detection system suited to their application domain, including category-specific models like YOLO [11] variants or open-vocabulary systems like GroundingDINO [6]. When detection systems provide only bounding boxes, we automatically generate segmentation masks using SAM2 [4] within the detected regions to ensure complete object-level analysis.

We formulate object attribution as a cooperative game where detected objects serve as players and the VLM’s response represents the outcome. For objects  $O = \{o_1, o_2, \dots, o_n\}$ , each object’s Shapley value  $\phi_i$  quantifies its contribution:

$$\phi_i = \sum_{S \subseteq O \setminus \{o_i\}} \frac{|S|!(|O| - |S| - 1)!}{|O|!} [v(S \cup \{o_i\}) - v(S)]$$

where  $v(S)$  measures the VLM’s response when only objects in subset  $S$  remain visible.

#### 4.2. Object Perturbation Strategy

The central challenge lies in removing target objects while preserving scene context for accurate attribution. We propose Bounding Box with Overlap Avoidance (BBOA) and evaluate it against two established baselines.

Precise masking applies exact segmentation boundaries but creates irregular occlusions that may introduce visual artifacts. Bounding box masking uses rectangular regions but risks occluding adjacent objects in dense scenes.

BBOA combines the advantages of both approaches through a three-step process: first masking the target object’s bounding box region, then identifying other objects whose segmentation masks intersect this region, and finally restoring those overlapping objects by unmasking their precise boundaries. This strategy ensures complete target removal while preserving neighboring objects regardless of scene density.

#### 4.3. Computational Implementation

Exact Shapley computation requires evaluating  $2^n$  object subsets, which becomes computationally prohibitive for scenes with many objects. We employ sampling-based approximation that reduces VLM queries from exponential to linear scaling, typically requiring 100-300 evaluations for scenes with 10-15 objects and completing analysis within 30-60 seconds.

Response similarity is measured using semantic embedding approaches through sentence transformers [12] or lexical similarity metrics depending on application requirements. The framework operates entirely through VLM input-output interfaces, maintaining compatibility with both open-source and commercial models without requiring access to internal representations.

### 5. Experimental Evaluation

We evaluate PixelSHAP’s effectiveness for object-level attribution in vision-language models through systematic comparison with existing black-box methods on carefully constructed human-annotated datasets.

#### 5.1. Dataset Construction

The absence of suitable benchmarks for object-level VLM attribution necessitated creating specialized evaluation datasets. We developed two complementary datasets with human annotation protocols designed to assess object-level interpretability across different visual domains.

**BDD10K Traffic Dataset:** We selected 250 representative images from the Berkeley DeepDrive dataset (BDD10K) [18], focusing on driving scenarios containing multiple traffic participants (vehicles, pedestrians, cyclists, traffic signs). Each scene was chosen to represent common driving situations where understanding object-level attention becomes safety-critical: intersections with multiple vehicles, crosswalks with pedestrians, and complex urban environments with mixed traffic.

Three experienced annotators followed a structured protocol: first, they randomly selected one object from each traffic scene, then formulated driving-relevant questions that would require focusing on that specific object to answer correctly (e.g., "Which vehicle poses the greatest safety concern?" or "What traffic element should influence the driver’s next action?"). This approach ensures unbiased ground truth while maintaining realistic question formulation. We measured inter-annotator agreement using Fleiss’ kappa [2] and retained only scenes achieving substantial consensus ( $\kappa > 0.7$ ). Figure 3 illustrates representative examples from this dataset, showing the diversity of objects and question types used in evaluation.

**COCO General Dataset:** To demonstrate broader applicability beyond traffic scenarios, we created a complementary dataset using 250 images selected from COCO [5] validation set. Following identical annotation protocols, annotators first randomly selected objects from general visual scenes, then generated focused questions requiring attention to those specific objects. This dataset enables assessment of PixelSHAP’s effectiveness across diverse visual contexts while maintaining the same evaluation framework.

Both datasets are publicly available on Hugging

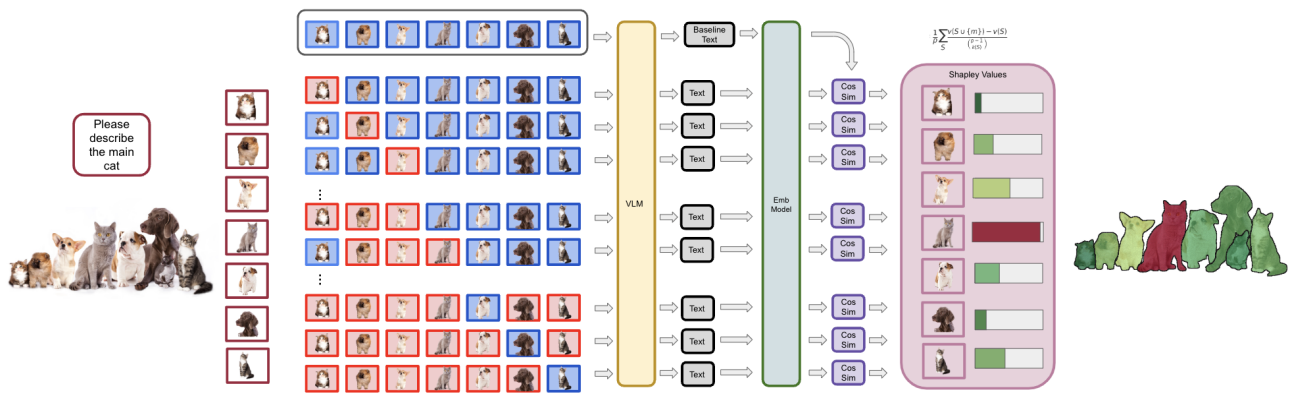


Figure 2. Overview of the PixelSHAP framework. The method systematically perturbs object groups, queries a vision-language model (VLM), and computes Shapley values to quantify object importance.

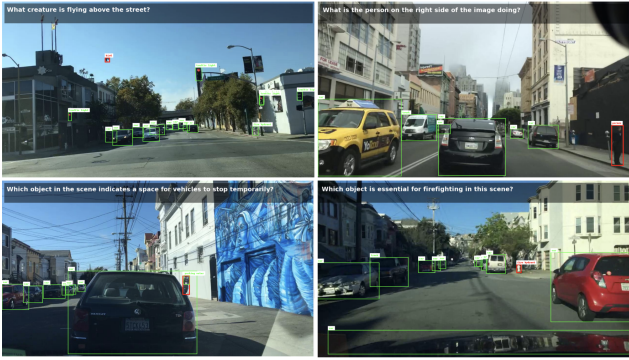


Figure 3. Sample annotations from BDD10K Traffic Dataset showing diverse object types and corresponding questions. Each example demonstrates how human annotators formulated questions requiring attention to specific objects for accurate answering.

Face[17], providing standardized benchmarks for future research in object-level VLM interpretability.

## 5.2. Evaluation Protocol

Our evaluation protocol measures how well attribution methods identify the same objects that human experts consider most relevant for answering given questions. For each image-question pair, we provide the question and corresponding answer to the VLM, then apply different attribution methods to identify which objects the model should focus on. We compare these attributions against human annotations to assess attribution quality.

This framework enables direct comparison of different attribution approaches while maintaining consistency with human reasoning patterns about object relevance in visual question answering tasks.

## 5.3. Masking Strategies

PixelSHAP’s effectiveness depends critically on the masking strategy used during object occlusion. We investigate two primary approaches for handling object removal during attribution computation:

**Precise Masking:** Objects are masked exactly according to their segmentation boundaries, replacing object pixels with neutral background or inpainting [9]. This approach maintains precise object boundaries but may create artificial visual artifacts at object edges.

**Bounding Box Occlusion with Adjustment (BBOA):** Objects are occluded using expanded bounding boxes that fully contain the object while minimizing overlap with neighboring objects. This strategy avoids edge artifacts and prevents unintended masking of adjacent objects that might confound attribution computation.

Figure 4 illustrates these different masking approaches and their impact on attribution quality. The BBOA strategy demonstrates superior performance by ensuring complete object occlusion while preserving the integrity of surrounding visual context.

## 5.4. Baseline Comparison Framework

We establish PixelSHAP’s effectiveness through comparison with existing black-box interpretability methods adapted for object-level analysis. Since direct comparison requires operating at the same semantic granularity, we adapt pixel-level methods to produce object-level attributions.

**RISE-Objects:** The original RISE method [10] generates pixel-level importance maps through random masking. We adapt RISE to operate at object-level granularity by randomly masking subsets of detected objects and measuring resulting changes in model output similarity. This preserves RISE’s core perturbation methodology while enabling fair



Figure 4. Comparison of masking strategies for object occlusion in PixelSHAP. (a) Precise masking follows exact segmentation boundaries. (b) Bounding Box Occlusion with Adjustment (BBOA) uses expanded boxes to ensure complete occlusion while minimizing interference with neighboring objects. BBOA consistently achieves better attribution performance across different scenarios.

comparison at the semantic object level.

**Simple Heuristic Baselines:** We include largest object (by bounding box area) and central object (closest to image center) as basic attribution methods. These baselines test whether sophisticated attribution approaches provide meaningful improvements over simple assumptions about visual attention.

**Random Baseline:** Random object selection establishes the performance floor and validates that our evaluation metrics capture meaningful attribution quality differences.

5.5. Evaluation Metrics

We assess attribution quality using metrics aligned with human annotation protocols and practical interpretability needs:

**Recall@1:** Percentage of test cases where the highest-attributed object matches human expert annotation. This metric directly measures whether attribution methods identify the same object that human experts consider most relevant.

**Recall@3:** Percentage where the human-annotated target object appears among the top-3 attributed objects, providing insight into attribution ranking quality.

**Mean Reciprocal Rank (MRR):** Average inverse rank of the ground-truth object across all test cases, offering a nuanced view of attribution accuracy that accounts for ranking position.

5.6. Results and Analysis

Table 1 presents comprehensive performance comparison across four representative VLMs on both datasets. The results demonstrate that PixelSHAP with BBOA masking achieves the best performance in nearly all scenarios,

though some competitive cases reveal interesting model-specific characteristics.

**Model-Specific Performance Patterns:** Gemini-2.0-flash achieves the highest overall performance across both datasets, with particularly strong results on COCO general scenes (67.48% Recall@1) and leading performance on traffic scenarios (64.7% Recall@1). GPT-4o demonstrates competitive performance on traffic scenarios (63.8% Recall@1), while both LLaVA-v1.5-7B and LLaMA-3.2-11B-Vision show more modest but consistent results across datasets.

**Masking Strategy Analysis:** BBOA achieves the best performance in the vast majority of cases, though some notable exceptions highlight the complexity of optimal masking strategies. LLaMA-3.2-11B-Vision shows a rare case where precise masking slightly outperforms BBOA on traffic Recall@1 (56.1% vs 55.8%), while LLaVA-v1.5-7B demonstrates competitive performance where precise masking achieves higher Recall@3 and MRR scores on traffic scenarios. These close margins suggest that masking strategy optimization may be model-dependent in specific contexts.

**Attribution Method Robustness:** The consistently strong performance of BBOA across different models and datasets validates our approach, with typical improvements of 3-8 percentage points over precise masking and 10-20 percentage points over baseline methods. The few competitive cases where precise masking approaches BBOA performance (difference 0.3 percentage points) demonstrate that while BBOA is generally superior, the optimal strategy may require fine-tuning for specific model architectures.

**Baseline Comparison:** PixelSHAP variants substantially outperform simple heuristics and RISE-Objects across all conditions. RISE-Objects achieves moderate performance but consistently lags behind PixelSHAP approaches by 5-15 percentage points in Recall@1, confirming the benefits of principled Shapley-based attribution for object-level interpretability.

**Dataset-Specific Insights:** Performance patterns show interesting domain dependencies. Gemini-2.0-flash maintains strong advantages on both datasets, suggesting robust generalization capabilities. The traffic scenarios generally yield slightly higher absolute performance across models, potentially reflecting the more structured nature of driving scenes compared to diverse COCO imagery.

5.7. Computational Efficiency

PixelSHAP processing requires 35-65 seconds per image depending on object count and VLM inference speed, using approximately 2-3× the number of detected objects in API calls rather than the exponential scaling that naive Shapley computation would require. This represents practical computational requirements suitable for offline analysis and



Model	Method	BDD10K Traffic Dataset			COCO General Dataset		
		R@1 (%)	R@3 (%)	MRR	R@1 (%)	R@3 (%)	MRR
GPT-4o	PixelSHAP (BBOA)	<b>63.8</b>	<b>86.2</b>	<b>0.75</b>	<b>60.56</b>	<b>87.66</b>	<b>0.73</b>
	PixelSHAP (Precise)	59.2	82.1	0.71	57.61	84.71	0.69
	PixelSHAP (BBox)	55.7	78.4	0.67	53.18	85.20	0.68
	RISE-Objects	43.1	67.8	0.57	42.3	69.2	0.56
Gemini-2.0-flash	PixelSHAP (BBOA)	<b>64.7</b>	<b>85.9</b>	<b>0.76</b>	<b>67.48</b>	<b>89.17</b>	<b>0.77</b>
	PixelSHAP (Precise)	62.1	83.2	0.73	59.62	84.73	0.71
	PixelSHAP (BBox)	59.4	80.8	0.71	58.10	88.68	0.72
	RISE-Objects	47.6	72.1	0.61	45.7	71.6	0.59
LLaVA-v1.5-7B	PixelSHAP (BBOA)	<b>48.9</b>	71.4	0.61	<b>49.78</b>	<b>83.28</b>	<b>0.65</b>
	PixelSHAP (Precise)	48.2	<b>72.1</b>	<b>0.62</b>	49.27	75.38	0.61
	PixelSHAP (BBox)	45.6	68.9	0.59	43.88	76.32	0.59
	RISE-Objects	41.3	65.7	0.55	37.2	64.5	0.52
LLaMA-3.2-11B-Vision	PixelSHAP (BBOA)	55.8	<b>78.3</b>	<b>0.68</b>	<b>52.71</b>	<b>86.72</b>	<b>0.68</b>
	PixelSHAP (Precise)	<b>56.1</b>	77.9	0.68	49.76	80.27	0.64
	PixelSHAP (BBox)	53.4	76.2	0.66	50.76	80.32	0.65
	RISE-Objects	44.8	68.5	0.58	38.9	66.4	0.53
Largest Object		38.4	62.5	0.51	23.14	60.85	0.43
Central Object		31.7	58.1	0.46	36.92	70.62	0.52

Table 1. Object-level attribution performance comparison across VLMs and datasets. Results show mean performance over test sets. Bold indicates best performance for each model-method combination.

safety validation applications in autonomous driving systems.

5.8. Limitations and Future Work

Our evaluation reveals several limitations that inform future research directions. Performance degrades in extremely cluttered scenes (> 15 objects) where occlusion becomes pervasive. Attribution quality also depends on segmentation accuracy, creating dependency on upstream computer vision components.

**Segmentation Quality Impact:** We evaluate attribution degradation under noisy segmentation by introducing controlled errors (10-30% mask boundary deviation) to ground-truth objects. Performance drops 8-15% with moderate noise, confirming segmentation dependency while demonstrating reasonable robustness to realistic segmentation errors.

The varying performance patterns across models suggest that future work should explore model-specific attribution strategies that account for architectural differences in visual reasoning capabilities.

5.9. Qualitative Examples

Figure 5 demonstrates PixelSHAP’s ability to provide intuitive, context-sensitive attributions across different query types and scenarios.

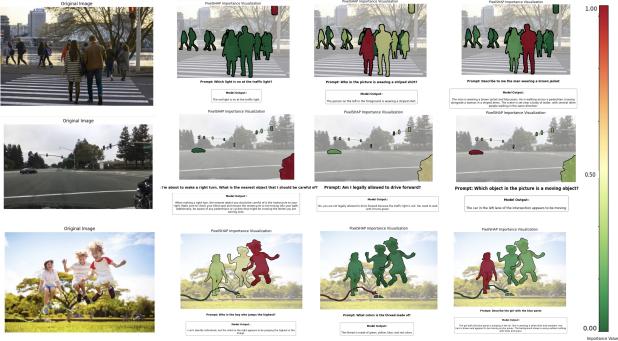


Figure 5. PixelSHAP attribution examples across traffic and general scenarios. Each row shows the same scene analyzed with different questions, demonstrating context-sensitive attribution. Red intensity indicates object importance scores.

**Traffic Scene Analysis:** In driving scenarios, PixelSHAP correctly prioritizes safety-critical objects based on query context. When asked “Which vehicle should the driver monitor?”, the method emphasizes the approaching car rather than parked vehicles. For pedestrian-focused queries like “Is it safe to proceed?”, attribution shifts to highlight the person near the crosswalk while deemphasizing background traffic.





[4] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Yuxin Rolland, Linus Gustafson, Trevor Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *arXiv preprint*, 2023. 4

[5] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. *arXiv preprint*, abs/1405.0312, 2014. 4

[6] Shilong Liu, Feng Li, Hao Zhang, Xiao Zhang, Lei Zhu, Hang Wang, Jianlong Shi, Hongyang Li, and Hao Dong. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint*, 2023. 4

[7] OpenAI. Gpt-4o model card, 2024. Accessed via OpenAI API. 2

[8] L. Parcalabescu and A. Frank. Mm-shap: Multimodal shapley values for model interpretation. *arXiv preprint*, 2022. 2

[9] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context encoders: Feature learning by inpainting. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2536–2544, 2016. 5

[10] V. Petsiuk, A. Das, and K. Saenko. Rise: Randomized input sampling for explanation of black-box models. In *British Machine Vision Conference (BMVC)*, 2018. 2, 5

[11] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 779–788, 2016. 4

[12] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, 2019. 4

[13] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why should i trust you? explaining the predictions of any classifier. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, page 1135–1144, 2016. 2

[14] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Int. Conf. Comput. Vis.*, pages 618–626, 2017. 2

[15] Lloyd S. Shapley. A value for n-person games. *Contributions to the Theory of Games*, 2:307–317, 1953. 2

[16] G. B. M. Stan, R. Saunshi, N. Viswanathan, G. Sundaramoorthi, and A. Menon. Lvlm-interpret: An interpretability tool for large vision-language models. In *CVPR Workshop on Explainable AI for Computer Vision*, 2024. 2

[17] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Perric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, 2020. 5

[18] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 4