## **APPENDIX**

## A FRAMES VISUALIZATION

This section visualizes sample input keyframes from the videos to provide a clearer understanding of the data used in our experiments. Each row in Figure 7 represents a distinct video sequence fed into the LVLMs for analysis. These examples are representative of the scenarios in our dataset, encompassing a variety of everyday actions, objects, and environments. By visualizing the raw inputs, we aim to illustrate the visual complexities, such as changes in viewpoint, object scale, and partial occlusions, that the model must handle to perform accurate semantic reasoning.



Figure 7: Visualization of Input Keyframes. Each row displays a sequence of frames provided to the model as input for a specific video. The red bounding boxes highlight the ground-truth object pertinent to the task's question (e.g., the object being picked up, kicked, or taken). It is important to note that these bounding boxes are included here for clarity and were not provided to the model during inference.

## B THEORETICAL FRAMEWORK FOR VISUAL SEMANTIC CIRCUITS

In this section, we propose the theoretical framework built on three core principles that we hypothesize govern the internal computations of LVLMs: Information Localization, Progressive Semantic Refinement, and a Two-Stage Reasoning flow. We model the LVLM as a probabilistic system to formalize these principles

into specific, falsifiable predictions. This framework provides a principled foundation for understanding and predicting the model's behavior, which we then validate through targeted experiments.

### B.1 PROBABILISTIC MODEL FORMULATION

Let V be a video represented by a sequence of key frames, and Q be a textual question. The LVLM, denoted by  $\mathcal{M}$ , aims to generate an answer A. The process begins by encoding the video V into a set of N visual tokens,  $E_V = \{e_1, e_2, \ldots, e_N\}$ . The question Q is tokenized into M text tokens,  $E_Q = \{q_1, \ldots, q_M\}$ . The model then computes the probability of an answer A:

$$P(A|V,Q) = \mathcal{M}(E_V, E_Q). \tag{3}$$

Central to our investigation is the hypothesis that the set of visual tokens  $E_V$  can be partitioned into the subset  $E_O$  containing primary information about the specific object o, and the complementary subset  $E_C$  containing contextual information, such that  $E_V = E_O \cup E_C$  and  $E_O \cap E_C = \emptyset$ .

# B.2 PRINCIPLE OF INFORMATION LOCALIZATION

We begin with the hypothesis that to answer a specific question, the LVLM does not treat all visual tokens equally. Instead, we propose the *Principle of Information Localization*: task-critical information is spatially concentrated in the subset of tokens corresponding to the object of interest. Let this subset be  $E_O$ , with the remainder being contextual tokens  $E_C$ .

This principle leads to a direct, testable prediction. The informational value of a token set can be quantified by the degradation in model performance upon its ablation. We model this degradation as the KL divergence between the original and ablated posterior distributions:

$$\mathcal{L}_{\text{drop}}(E_S) = D_{KL} \left( P(A|E_V, E_Q) \parallel P(A|E_V \setminus E_S, E_Q) \right). \tag{4}$$

Our principle predicts that the information is concentrated in  $E_O$ . Formally, if we ablate the object tokens  $E_O$ , the resulting information loss should be significantly greater than ablating any other random subset of tokens  $E_R$  of the same size. This leads to the following inequality, which we aim to verify experimentally:

$$\mathbb{E}_{E_R \subset E_V, |E_R| = |E_O|} \left[ \mathcal{L}_{\text{drop}}(E_R) \right] \ll \mathcal{L}_{\text{drop}}(E_O). \tag{5}$$

Furthermore, we hypothesize that the model internally reasons over abstract concepts. This predicts that injecting a clean, symbolic representation of the object,  $e_{w_{\rm correct}}$ , should be even more effective than the noisy visual tokens  $E_O$ . This can be formalized as:

$$P(A^*|e_{w_{\text{correct}}}, E_C, E_Q) > P(A^*|E_O, E_C, E_Q).$$
 (6)

The ablation and injection experiments presented in Table 1 were designed to test these formal predictions.

# **B.3** Hypothesis of Progressive Semantic Refinement

We hypothesize that visual information is not processed into its final semantic form in a single step. Instead, we propose the model of *Progressive Semantic Refinement*, where hidden states associated with visual tokens transition from encoding low-level perceptual features in early layers to abstract, language-aligned concepts in later layers.

Let  $h_i^{(l)}$  be the hidden state for a token i at layer l. Let  $\mathcal{S}(w_{\text{correct}})$  be the semantic space associated with the correct object concept, represented by its text embedding  $e_{w_{\text{correct}}}$ . Our hypothesis predicts that for an object token  $i \in E_O$ , its representation  $h_i^{(l)}$  will become progressively more aligned with this semantic space as it passes through the network. We can formalize this predicted monotonic increase in alignment for layers l beyond a critical depth  $l_{\text{crit}}$  using a similarity metric:

$$sim(h_i^{(l)}, e_{w_{correct}})$$
 is a monotonically increasing function of  $l$  for  $l > l_{crit}$ . (7)

# 709 710 711

716 717

722

723

728

733 734 735

736

746

747

741

748 749 750 To test this prediction, we employ the logit lens technique, which projects intermediate hidden states into the vocabulary space. We measure the Correspondence Rate  $(C_R^{(l)})$  and Answer Probability  $(A_P^{(l)})$  to track this alignment across layers. The experimental results in Figure 4 are used to validate the existence and location of the predicted critical layer depth  $l_{crit}$ .

## B.4 THE TWO-STAGE REASONING HYPOTHESIS

Building on the previous principles, we hypothesize that the model's reasoning is not monolithic but follows an efficient, cognitively plausible two-stage process.

- 1. Stage 1 (Contextual Grounding): In the early layers ( $\mathcal{L}_{early}$ ), the model first processes contextual tokens  $(E_C)$  to establish a general understanding of the scene and the query.
- 2. Stage 2 (Focal-Point Refinement): In the late layers ( $\mathcal{L}_{late}$ ), after the context is established, the model focuses its attention on the specific object tokens  $(E_O)$  to extract fine-grained details necessary for a precise answer.

This hypothesis can be formalized by considering the sensitivity of the final prediction to attention weights at different layers. Let  $\nabla_{\alpha_n^{(l)}} \log P(a_1^*)$  be the gradient of the log-probability of the correct answer with respect to the attention weights from a set of tokens S at layer l. Our two-stage hypothesis predicts a shift in sensitivity:

$$\sum_{l \in \mathcal{L}_{\text{early}}} \left\| \nabla_{\alpha_C^{(l)}} \log P(a_1^*) \right\| > \sum_{l \in \mathcal{L}_{\text{early}}} \left\| \nabla_{\alpha_O^{(l)}} \log P(a_1^*) \right\|, \tag{8}$$

$$\sum_{l \in \mathcal{L}_{\text{late}}} \left\| \nabla_{\alpha_O^{(l)}} \log P(a_1^*) \right\| > \sum_{l \in \mathcal{L}_{\text{late}}} \left\| \nabla_{\alpha_C^{(l)}} \log P(a_1^*) \right\|. \tag{9}$$

These inequalities formalize the "context-first, detail-later" strategy as a testable prediction. We designed the attention masking experiments in Table 2 to directly probe these sensitivities and validate our two-stage reasoning hypothesis.

# SCALING EXPERIMENTS

We conducted supplementary experiments to verify that our conclusions generalize to larger-scale models. We replicated our core analyses on the LLaVA-NeXT-34B model variants, with results that closely mirror those presented in the main body of the paper. The visual token ablation study on the 34B models reaffirms the principle of spatial localization for semantic information. Ablating object-specific tokens incurs significantly more substantial performance degradation than removing larger quantity of random tokens (in Table 3).

Furthermore, our semantic tracing analysis on the 34B architecture, depicted in Figure 8, reveals a conceptual emergence pattern consistent with our earlier observations. Both the Correspondence Rate and Answer Probability remain negligible through the initial layers before exhibiting a sharp, concurrent rise beginning around layer 40. This trend indicates that abstract, language-aligned concepts are consolidated in the deeper layers of the network, irrespective of model scale. These scaling experiments provide robust evidence that the mechanisms of semantic localization and late-stage conceptual formation are fundamental properties of the tested LVLM architectures.

7	5	2
7	5	3
	5	
7	5	5
	5	
	5	
7	5	8
	5	
7	6	0
7	6	1
	6	
7	6	3
7	6	4
7	6	5
7	6	6
7	6	7
7	6	8
7	6	9
7	7	0
7	7	1
7	7	2
7	7	3
7	7	4
7	7	5
	7	6
7	7	7
7	7	8
7	7	9
7	8	0
7	8	1
7	8	2
7	8	3

Ablation	Tokens	LLaVA-N	eXT-34B-I	LLaVA-NeXT-34B-V				
Type	Number	Open	Close	Open	Close			
Control Groups (Low Tokens)								
Baseline	0	$100.0_{\downarrow 0}$	$100.0_{\downarrow 0}$					
Register	13	$55.3_{44.68}$	$96.9_{\downarrow 3.14}$	$55.3_{44.68}$	$1.8_{\downarrow 98.2}$			
Object-based Ablation								
	304	$9.3_{190.75}$	$29.8_{\downarrow 70.16}$	$14.6_{\pm 85.37}$	$11.7_{188.35}$			
Object	413	$5.8_{194.24}$	$21.1_{\downarrow 78.88}$	$10.7_{\substack{+89.32}}$	$11.3_{188.72}$			
	573	$5.8_{ extstyle 94.24}$	$18.0_{\textcolor{red}{\downarrow}82.02}$	$10.3_{\textcolor{red}{\downarrow 89.74}}$	$11.8_{\downarrow 88.21}$			
Control Groups (High Tokens)								
	100	$54.3_{\downarrow 45.72}$	$96.9_{13.14}$	$93.1_{46.92}$	29.7470.26			
Random	350	$55.5_{144.5}$	$96.5_{13.49}$	$88.5_{\downarrow 11.54}$	31.0 168.9			
	500	$54.5_{\downarrow 45.55}$	$96.3_{\downarrow 3.66}$	$83.1_{\downarrow 16.92}$	30.8 169.23			
	900	$59.9_{40.14}$	$96.8_{\downarrow 3.17}$	$78.3_{1.21.67}$	$29.7_{\downarrow 70.33}$			

Table 3: Accuracy (%) from visual token ablation on question-answering performance across 34B models. The  $\downarrow$  symbol indicates the magnitude of this performance drop.

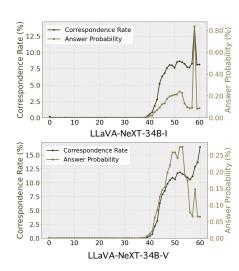


Figure 8: Quantitative analysis of semantic tracing on 34B model size.

Layer 3 Layer 4 Layer 5

Token/Layer Layer 1 Layer 2

1000	20
100	

ŀ	rompt	

<u>USER</u>: The input consists of a sequence of key frames from a video. Question: Which object was taken by the person?

ASSISTANT: The object is the

Answer: Towel

Response: Towel

Tottor is Edy or		Luyo. L	,	Luyo. 4	Luye. c
<s></s>	Institution	Архив	sierp	sierp	sierp
US	ona	Архив	pert	tap	pert
ER	Session	пута	oded	庄	pel
:	Portail	пута	<s></s>	<s></s>	Sav
		пута	Carter	•	•
<img001></img001>	sak	sak	sak	sak	sak
<img002></img002>	gresql	gresql	gresql	gresql	gresql
<img003></img003>	olas	olas	olas	olas	olas
<img004></img004>	<<	<<	<<	<<	<<
<img005></img005>	gresql	gresql	gresql	gresql	gresql
<img006></img006>	sak	sak	sak	sak	sak
<img007></img007>	yter	yter	yter	yter	yter
<img008></img008>	nahm	nahm	nahm	nahm	nahm
<img009></img009>	loop	loop	loop	loop	sak
<img010></img010>	loop	loop	loop	loop	loop
<img011></img011>	gresql	gresql	gresql	gresql gresql	
<img012></img012>	yter	yter	yter yter		yter
<img013></img013>	izi	izi	Emp Emp		conf
<img014></img014>	alle	alle	alle alle		conf
<img015></img015>	yter	yter	yter yter		yter
<img016></img016>	alle	alle	alle	arab	anno
<img017></img017>	yter	yter	yter	yter	yter

Figure 9: Qualitative example of the model correctly identifying an object. The user asks which object was taken by the person. The model correctly identifies the "Towel". The accompanying table shows the layer-by-layer semantic tracing for visual and text tokens.

### D **QUALITATIVE EXAMPLES**

We provide additional qualitative examples to visually illustrate the findings from our circuit-based analysis. These examples showcase the model's process of interpreting video frames to answer specific questions about objects and actions. Each figure includes the input video frames, the posed question, the model's response, and a table showing the semantic evolution of key tokens across different layers, as analyzed through our semantic tracing circuit (Circuit 2). More examples in the anonymous interactive demo website.

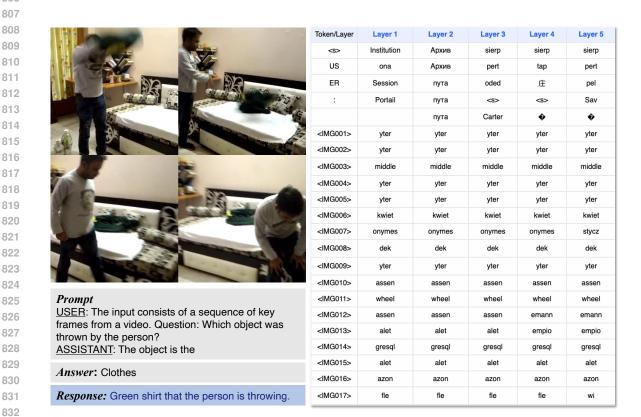


Figure 10: Qualitative example where the model is prompted to identify a thrown object. The model successfully responds that a "Green shirt" was thrown, correctly identifying both the object and its color. The table illustrates the semantic trace, showing how the model processes the visual information through its layers.

#### Ε ETHICS STATEMENT

This research adheres to the ICLR Code of Ethics and aims to enhance the understanding of the internal spatiotemporal reasoning mechanisms within LVLMs through circuit-based analysis, in order to drive the development of more robust and interpretable models. The experiments are based entirely on public academic datasets (e.g., the STAR benchmark), and we acknowledge that these contain videos of human activities, which were used solely for their intended academic analytical purposes. To ensure research integrity, we explicitly state that LLMs were used only for language polishing after the manuscript was written, and that all core scientific contributions originate from the human authors.

# F REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of this research. All experiments are based on publicly available models (the LLaVA-NeXT and LLaVA-OV series) and datasets (the STAR benchmark). We detail our methods for filtering and processing data from the STAR benchmark in the "Dataset Curation" part of Section 3, and provide visualization samples of keyframes in Appendix A. The core methodology of our research, including the specific settings, intervention methods, and evaluation metrics for the three analytical circuits (Visual Information Auditing, Semantic Tracing, and Attention Flow), is thoroughly elaborated in Section 3 (Subsections 3.1, 3.2, and 3.3), which includes key mathematical formulas and parameter definitions. The theoretical framework supporting our experimental design is fully formalized in Appendix B. Furthermore, an anonymous interactive demo website is provided in Appendix D for reviewers to explore additional qualitative results. We believe these detailed descriptions are sufficient to support the reproduction of this work. All code will be made available upon acceptance of the manuscript.

	Token/Layer	Layer 1	Layer 2	Layer 3	Layer 4	Layer 5
	<s></s>	Institution	Архив	sierp	sierp	sierp
	US	ona	пута	pert	tap	pure
	ER	Session	пута	Session	庄	Ori
	:	Portail	пута	<s></s>	<s></s>	Sav
		ctl	пута	Carter	•	IR
	<img001></img001>	МП	МΠ	МП	МΠ	aks
	<img002></img002>	amen	amen	amen	alberga	amen
	<img003></img003>	ката	ката	ката	ката	ката
	<img004></img004>	jes	jes	jes	jes	jes
	<img005></img005>	anno	МΠ	МП	МП	anno
	<img006></img006>	隆	隆	隆	Sito	隆
	<img007></img007>	gresql	gresql	gresql	gresql	gresql
	<img008></img008>	ummy	ummy	ummy	ummy	ummy
	<img009></img009>	osz	osz	osz	Sito	anim
	<img010></img010>	opf	opf	opf	opf	opf
Prompt	<img011></img011>	pur	pur	pur	pur	pur
USER: The input consists of a sequence of key	<img012></img012>	ode	ode	ode	ode	ode
frames from a video. Question: Which object was picked up by the person?	<img013></img013>	wheel	wheel	wheel	ava	置
ASSISTANT: The object is the	<img014></img014>	arab	arab	ode	ode	arab
Answer: Box	<img015></img015>	gresql	gresql	gresql	end	t
Answer. Box		кал	кал	CURLOPT	CURLOPT	cub
Response: Box of shoes	<img017></img017>	way	way	way	way	way

Figure 11: Qualitative example demonstrating the model's ability to recognize an object being picked up. The model correctly identifies the object as a "Box of shoes". The semantic tracing table displays the evolution of token representations across five layers that contribute to this accurate identification.