



Project: <https://haofei.vip/VoT/>

Paper: <https://is.gd/fcfZeO>

Code: <https://github.com/scofield7419/Video-of-Thought>

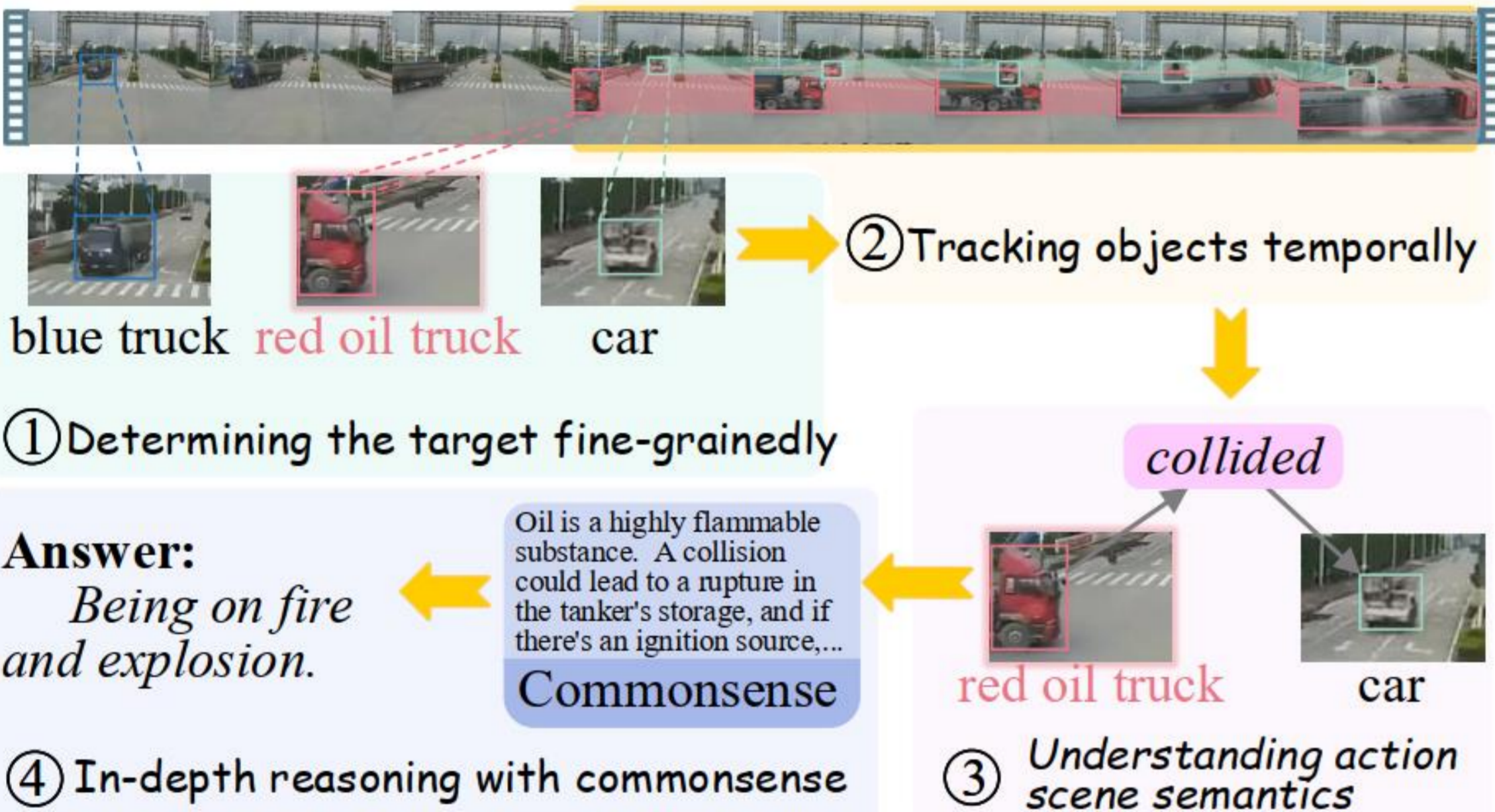
TL;DR

This work for the first time propose the video Chain-of-Thought (CoT) reasoning framework, Video-of-Thought (VoT), for achieving human-level video reasoning. VoT decomposes raw complex problems into a chain of sub-problems, and reasons through multiple steps from low to high levels, enabling not only pixel perceptive recognition but also semantic cognitive understanding of videos.

Motivation

- Keynotes of human cognition patterns on video understanding/reasoning:

Question: What will happen to *the red oil tanker truck*?



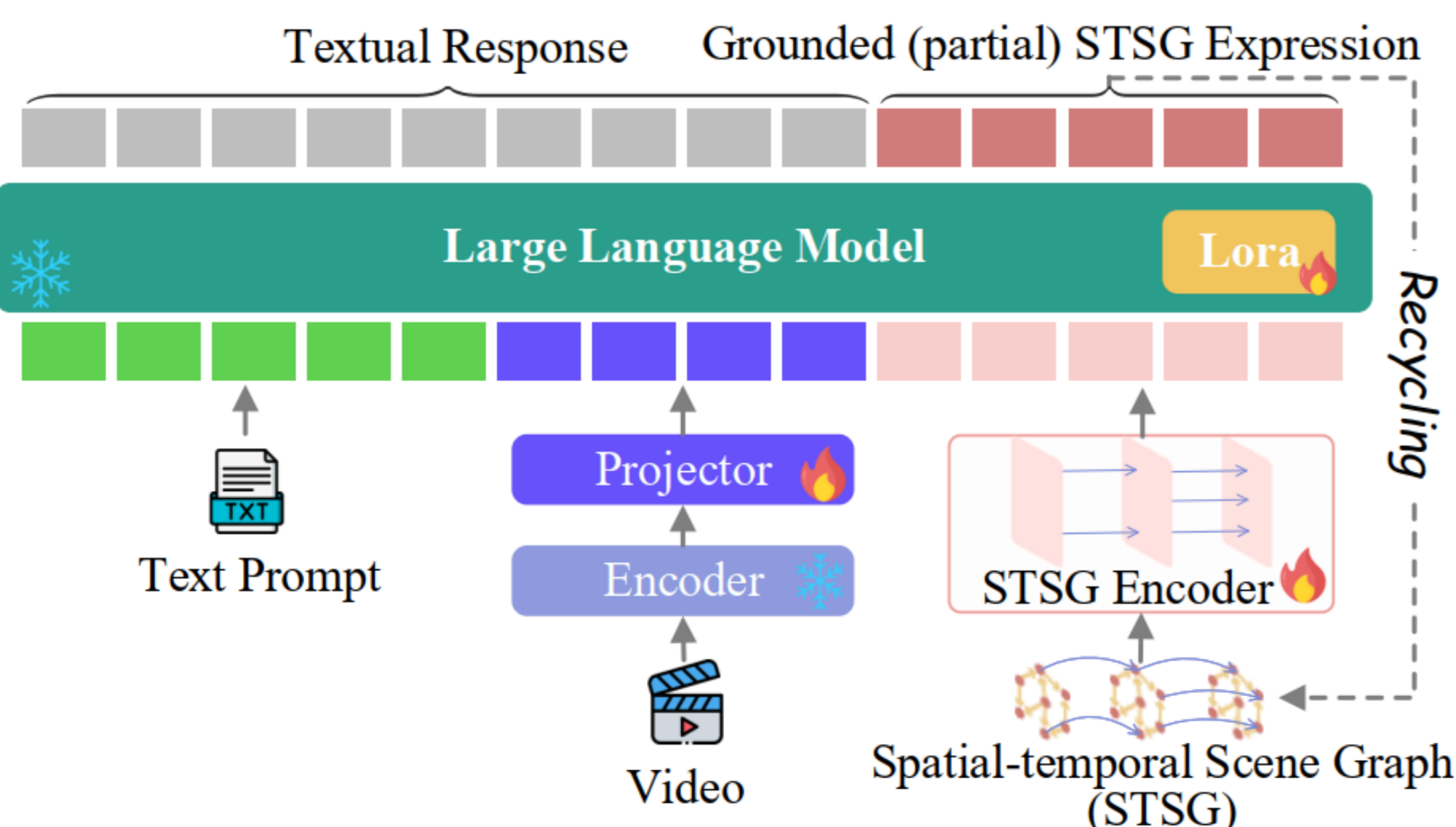
First, to achieve precise content perception, a fine-grained perceptive pixel understanding of the video movement is necessary

Second, profound reasoning demands cognitive capabilities allowing reasonable explanation and even causal imagination, i.e., with a reservoir of commonsense knowledge to link video pixels to the factual world.

Third, for humans, video reasoning is not an instantaneous process but follows a multi-hop procedure from lower level to higher level.

Methodology

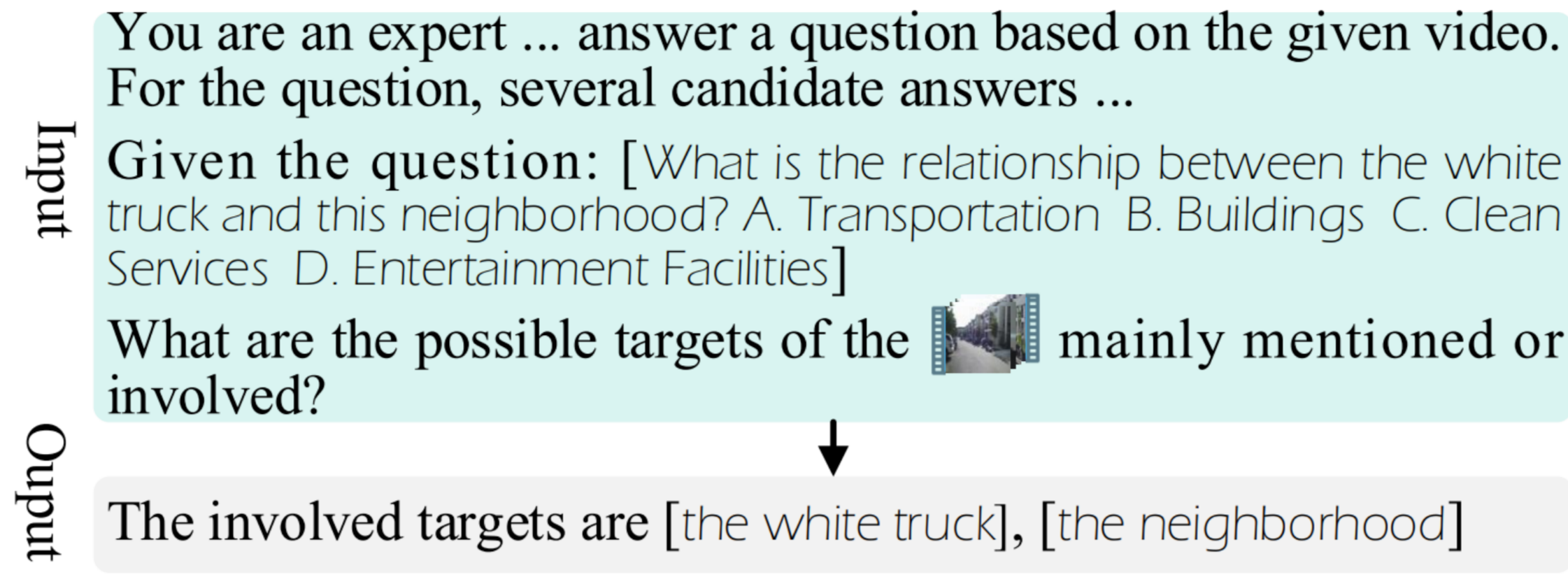
- MotionEpic**: Fine-grained Spatial-temporal Grounded Video MLLM



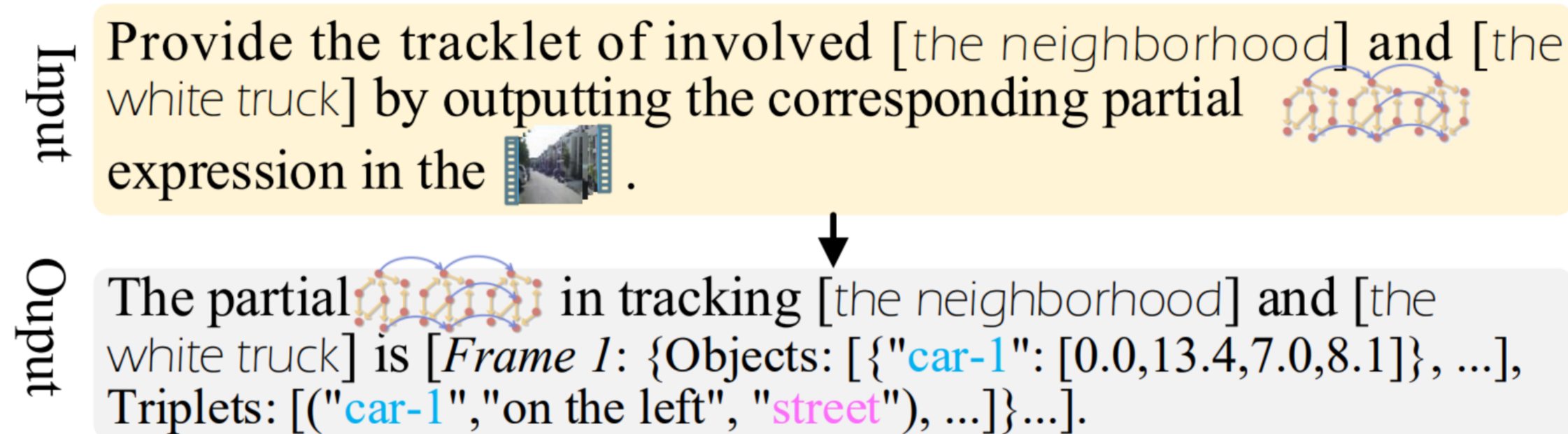
Methodology

- Video-of-Thought Reasoning Framework**

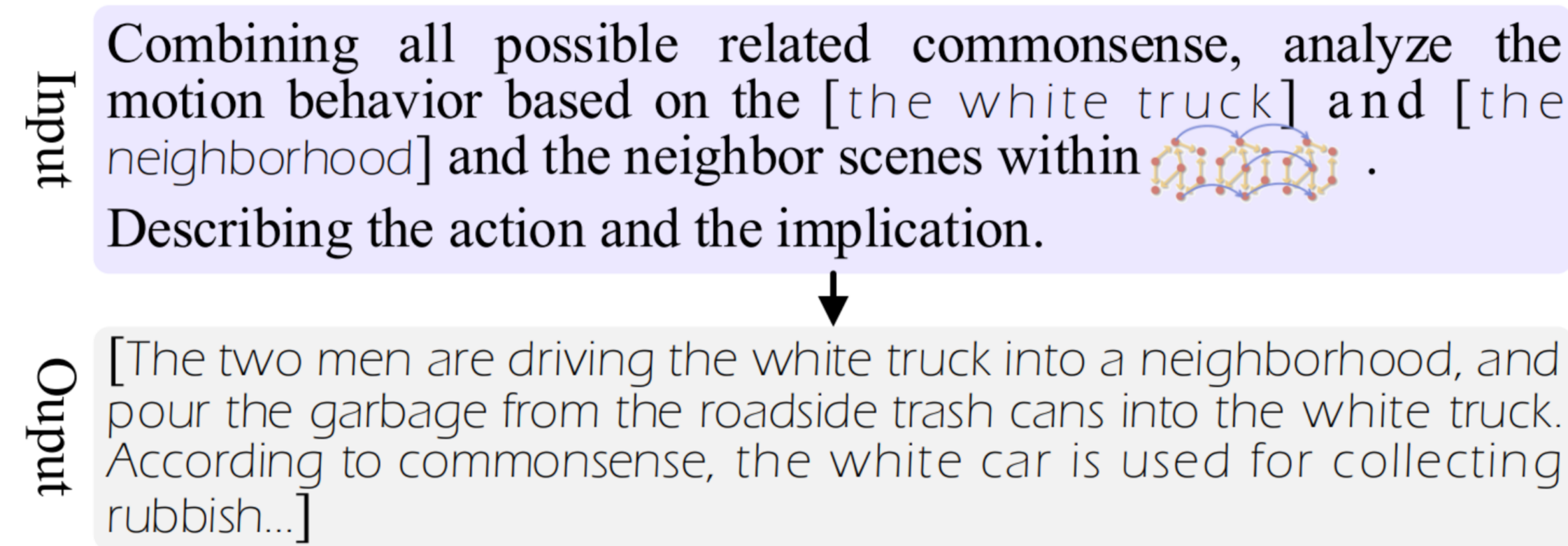
Step-1: Task Definition and Target Identification



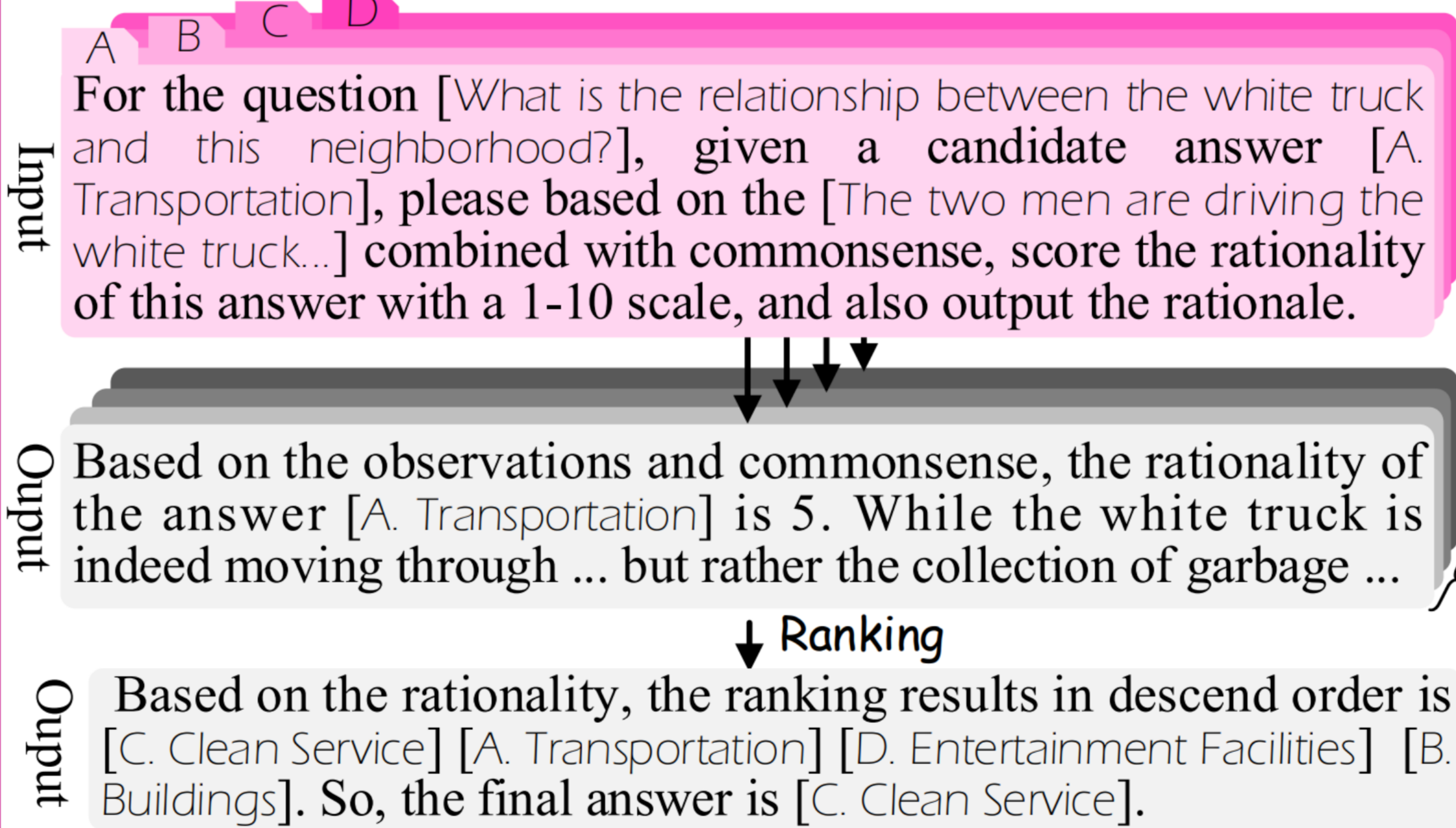
Step-2: Object Tracking



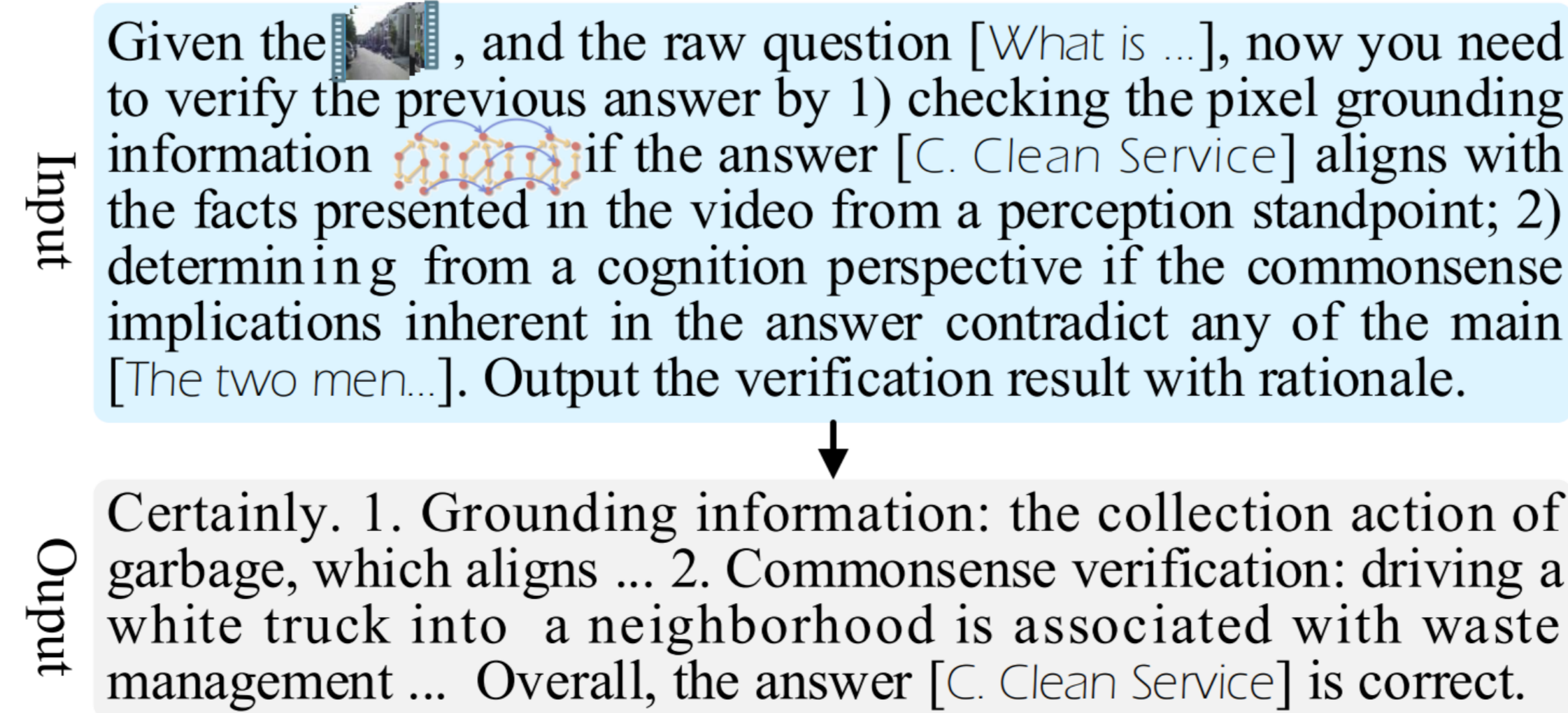
Step-3: Action Analyzing



Step-4: Question Answering via Ranking



Step-5: Answer Verification



Experiment

- Quantitative Results**

Table 1: Results on four VideoQA datasets. STAR data includes four subsets: Interaction (Int.), Sequence (Seq.), Prediction (Pre.), Feasibility (Fea.). The best scores of baselines are underlined, and the new best results are **bold**.

Model	VLEP		STAR		IntentQA		Social-IQ	
	Int.	Seq.	Pre.	Fea.	2-Way	4-Way	2-Way	4-Way
• SoTA baselines								
InternVideo	63.9	62.7	65.6	54.9	51.9	-	-	-
LLaMA-VQA	71.0	66.2	67.9	57.2	52.7	-	-	-
VLAP	69.6	<u>70.0</u>	<u>70.4</u>	<u>65.9</u>	<u>62.2</u>	-	-	-
SeViLA	68.9	63.7	70.4	63.1	62.4	-	-	-
VideoChat	62.0	63.2	66.8	54.1	49.6	59.3	67.7	37.8
Video-LLaVA	65.8	64.3	67.0	56.5	50.1	62.5	68.9	39.2
• CoT								
Video-LLaVA	65.7	65.0	67.7	57.8	52.0	63.2	69.5	40.4
Video-LLaVA+stsg	67.0	65.9	68.9	58.7	53.7	64.9	70.4	41.7
MotionEpic	68.2	66.8	69.6	60.6	57.4	<u>66.1</u>	<u>71.7</u>	<u>43.0</u>
• VoT								
MotionEpic	73.4	71.5	72.6	66.6	62.7	70.8	72.8	45.0

Table 3: Results on NEXT-QA data.

Model	Acc@All	Acc@C	Acc@T	Acc@D
• SoTA baselines				
InternVideo	63.2	62.5	58.5	75.8
HiTeA	63.1	62.4	58.3	75.6
LLaMA-VQA	72.0	72.7	69.2	75.8
SeViLA	73.8	73.8	67.0	81.8
VLAP	<u>75.5</u>	<u>74.9</u>	<u>72.3</u>	<u>82.1</u>
Video-LLaMA	60.6	59.2	57.4	72.3
VideoChat	61.8	63.5	61.5	74.6
Video-ChatGPT	64.4	66.9	64.1	75.7
Video-LLaVA	66.3	67.7	63.8	75.9
• CoT				
Video-LLaVA	67.7	69.0	65.9	76.5
Video-LLaVA+stsg	68.0	71.6	67.6	78.9
MotionEpic	72.2	73.4	69.1	80.7
• VoT				
MotionEpic	76.0	75.8	74.6	83.3

- Analysis**

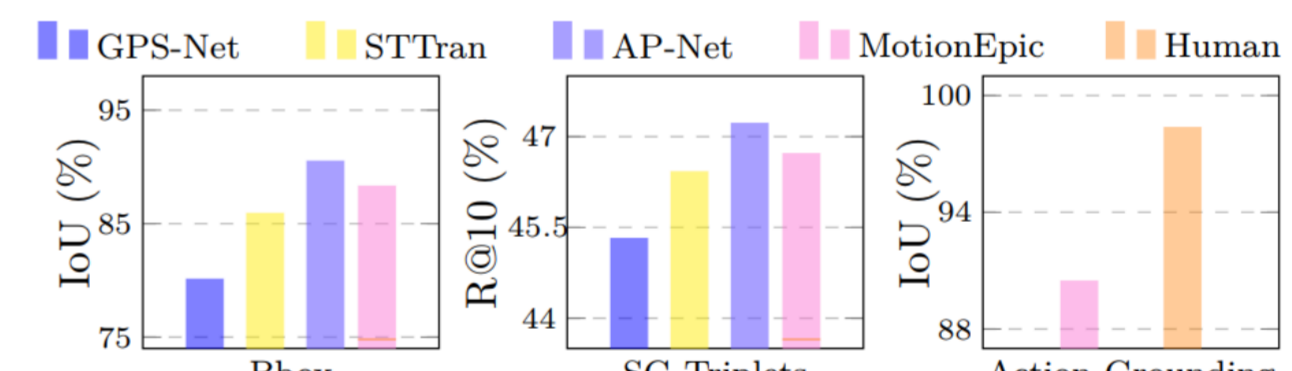


Figure 5: MotionEpic performance on object grounding, scene graph triplet classification, and action grounding.

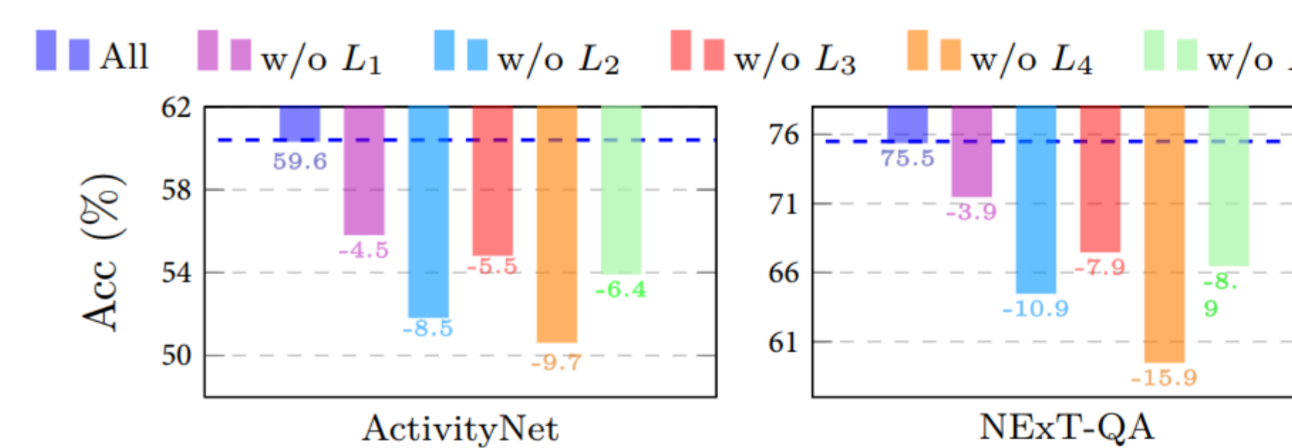


Figure 6: Performance drop (zero-shot) of MotionEpic after ablating different grounding-aware tuning item.

- Video Reasoning Visualization**

