

TUMTraf VideoQA: Dataset and Benchmark for Unified Spatio-Temporal Video Understanding in Traffic Scenes

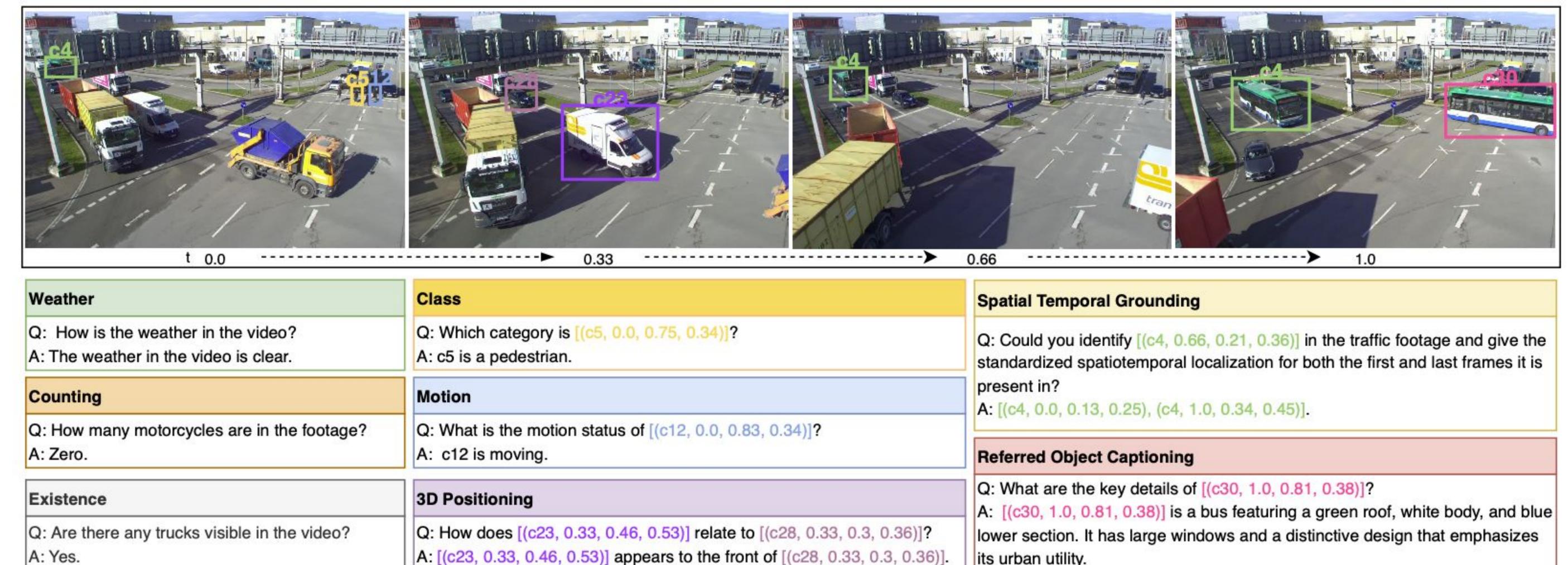
Xingcheng Zhou[†], Konstantinos Larintzakis, Hao Guo, Walter Zimmer, Mingyu Liu, Hu Cao[†], Jiajie Zhang, Venkatnarayanan Lakshminarasimhan, Leah Strand, Alois C. Knoll



Contact: Xingcheng.zhou@tum.de

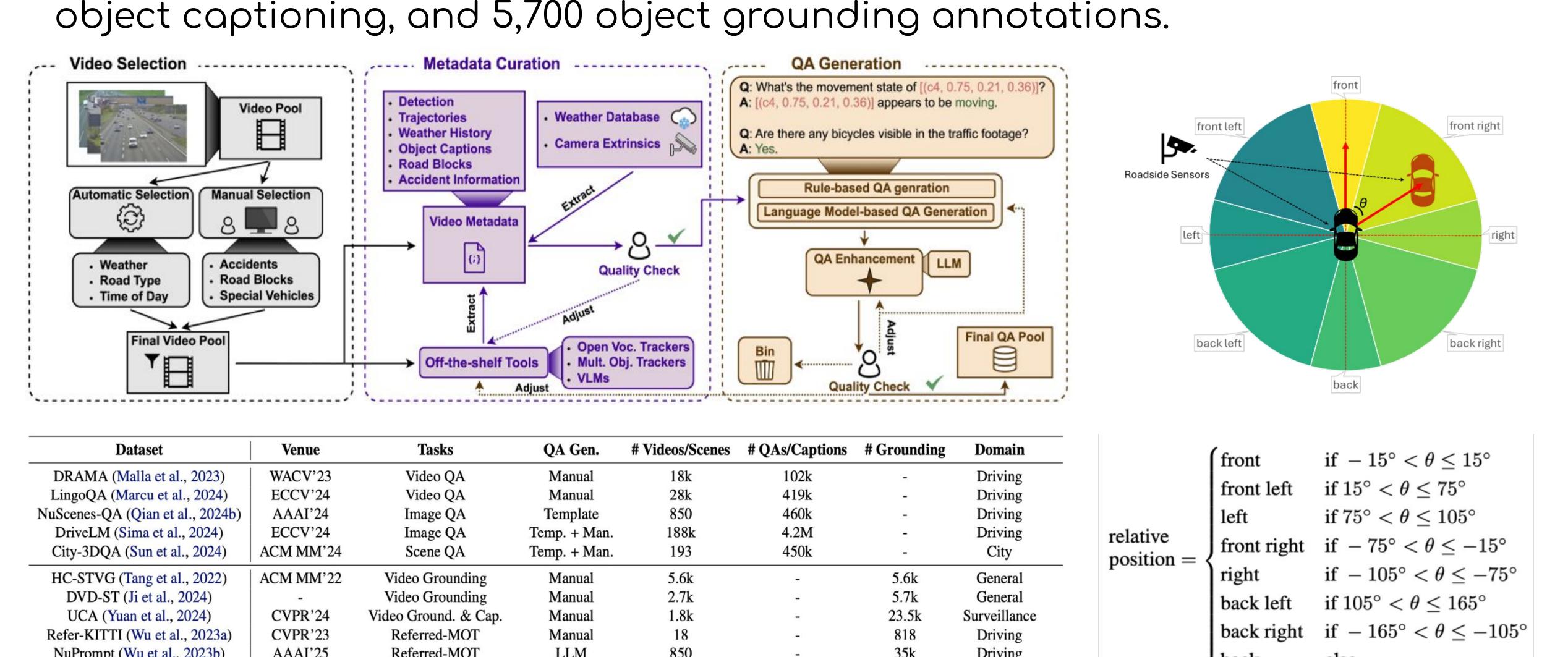
Introduction

- First large-scale video-language dataset for roadside traffic scenes understanding under diverse real-world conditions, including weather variations and accidents.
- Provide unified annotations for multi-choice QA, spatio-temporal object grounding, and referred object captioning, enabling fine-grained traffic video reasoning.



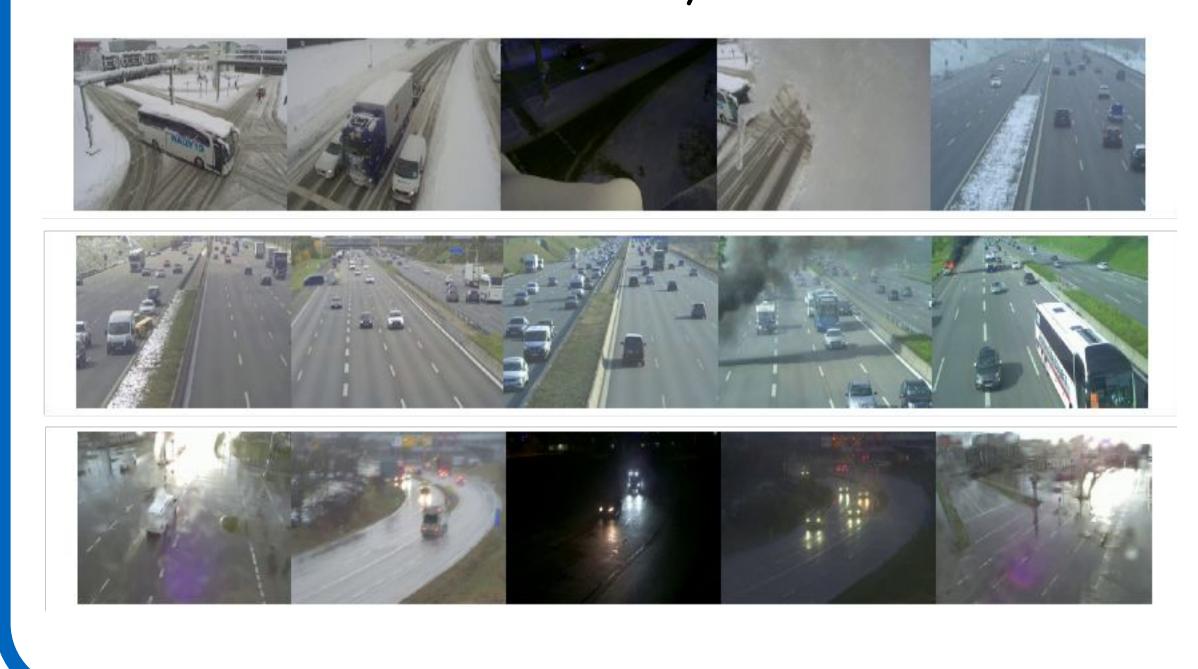
Dataset Curation

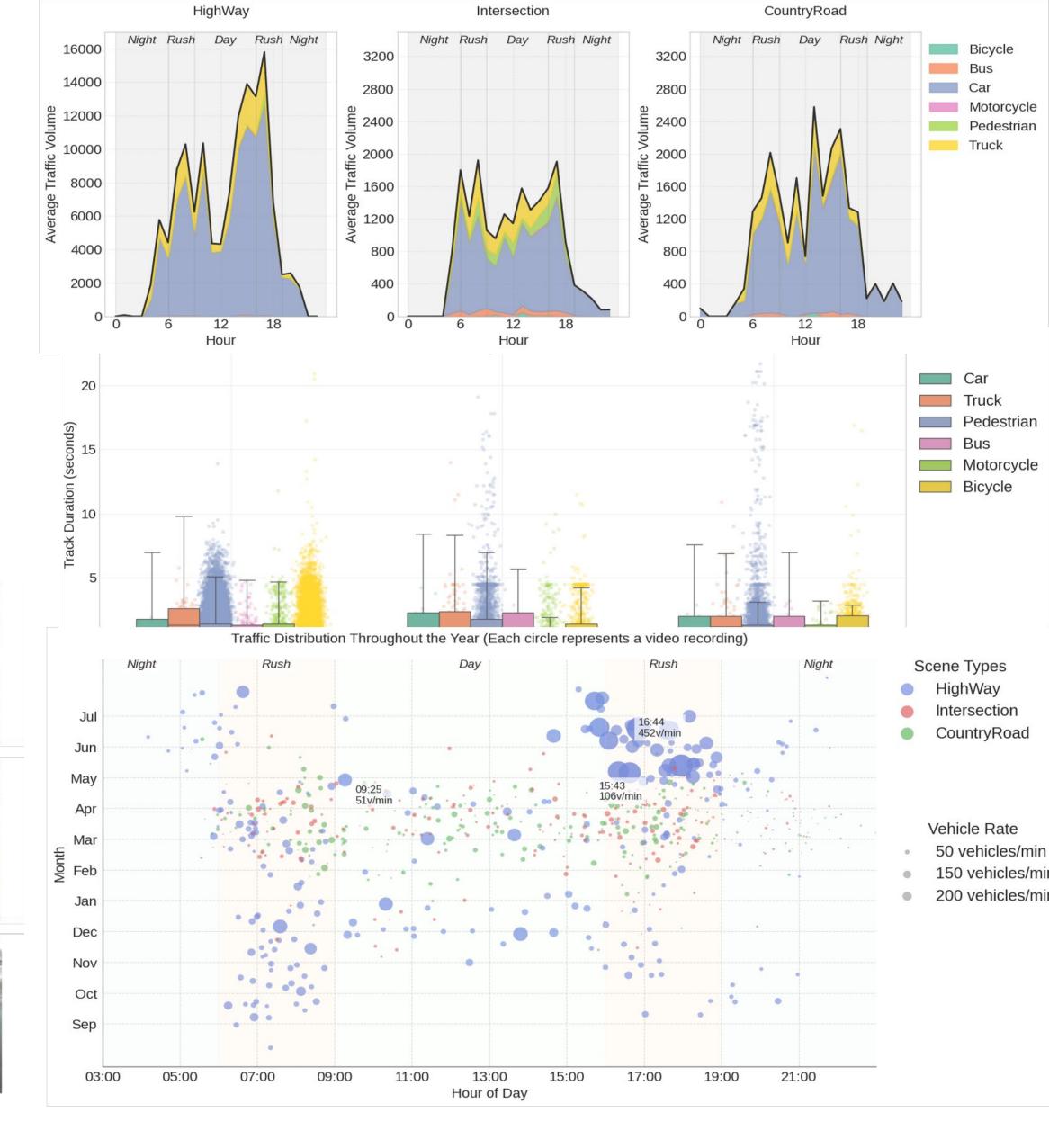
• Semi-automatic data pipeline leveraring LLM- and template-based methods with human verification, featuring 85,000 multiple-choice QA pairs, 2,300 object captioning, and 5,700 object grounding annotations.



Statistics

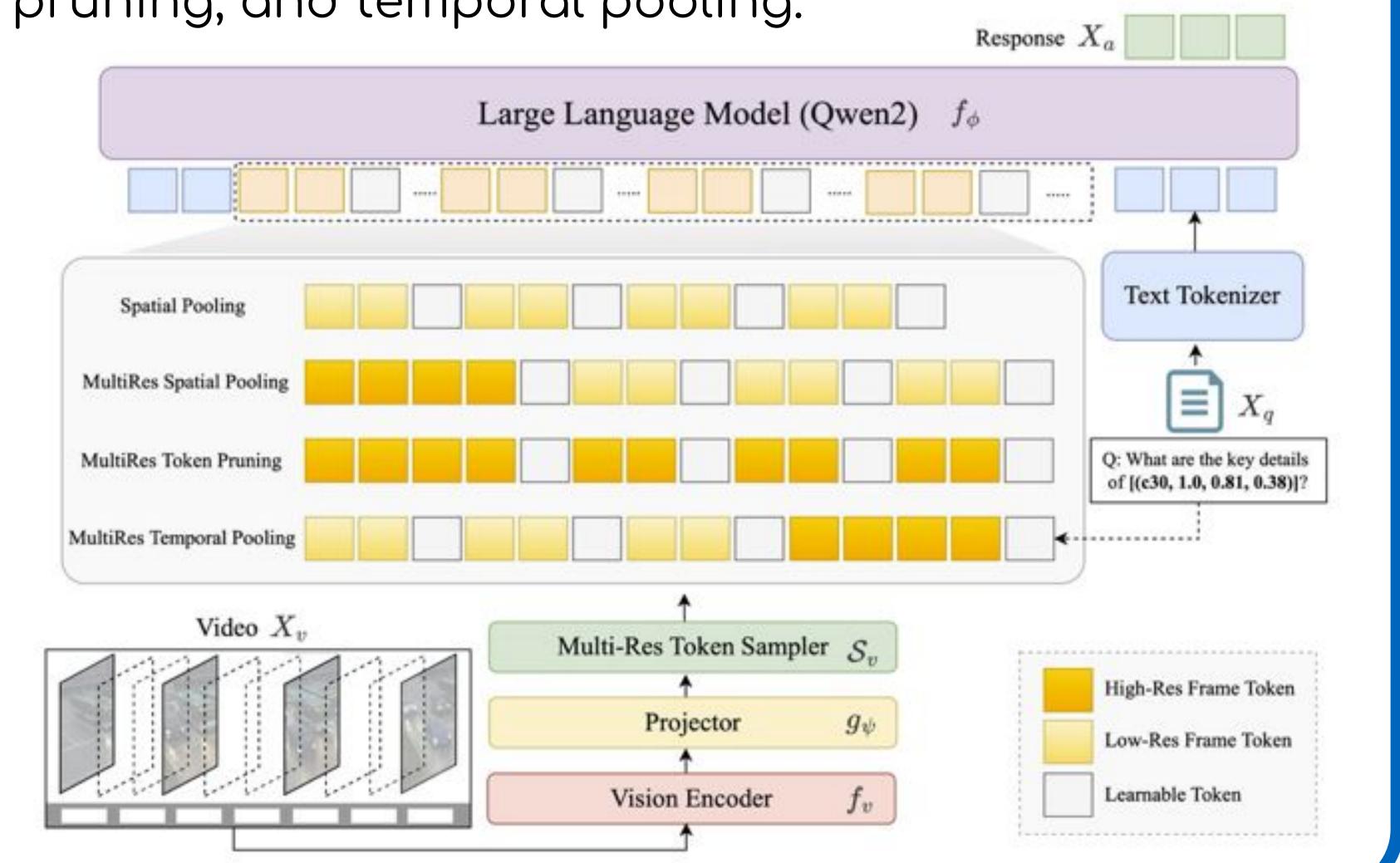
- Covers 1,000 roadside videos with diverse weather, lighting, and incidents across highway, urban, and rural scenes.
- Balanced coverage across times of day and seasons under varied traffic conditions, reflecting real-world traffic dynamics.





TraffiX-Qwen

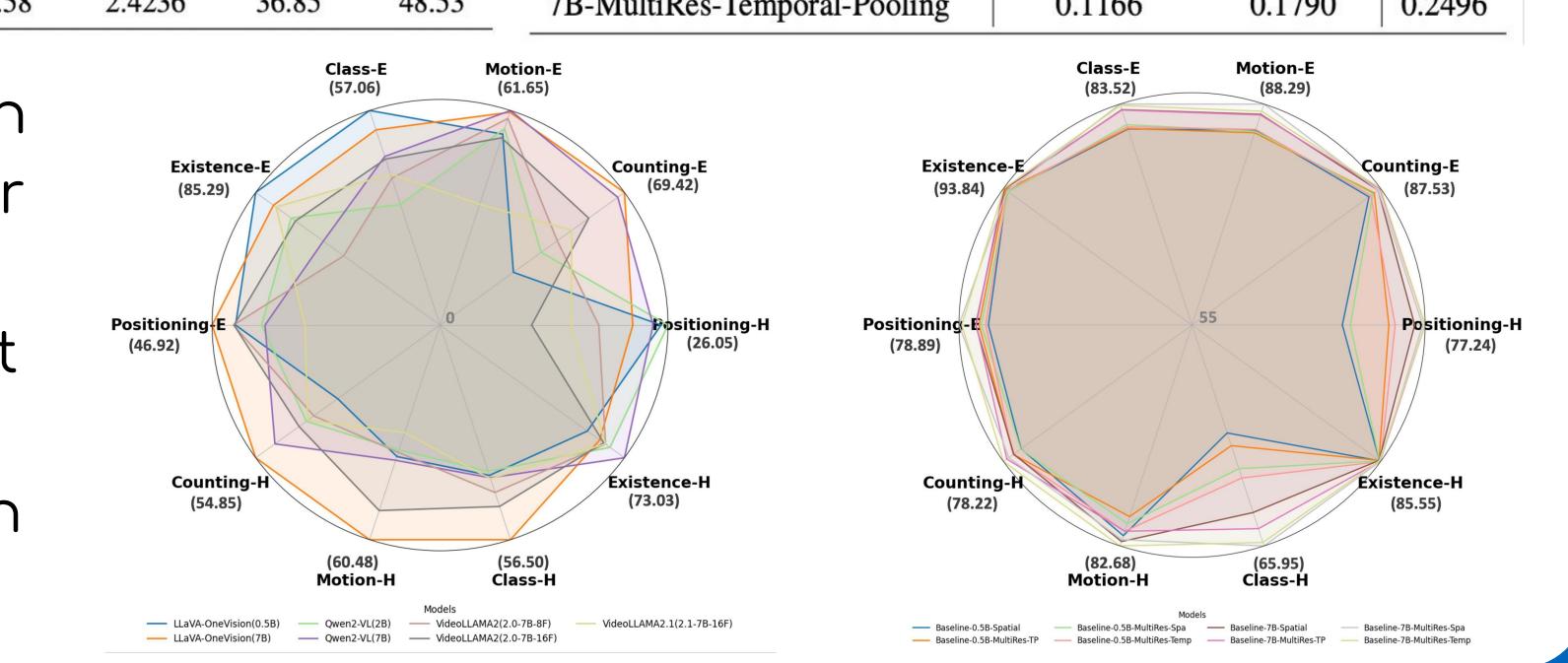
- Video-LLM for unified traffic video understanding, across multi-choice QA, spatio-temporal grounding, and referred object captioning.
- Leveraging multi-resolution visual sampling strategies, i.e., multi-res spatial pooling, token pruning, and temporal pooling.



Benchmark

<u>-</u>																			2 			
Models	Category	Positi	ioning	Cou	nting	9	tion	Cl	ass	Ex	istence	Overall	Model	Bleu_4	ROUGE_L	CIDEr	METEOR	SPICE	Model	Temporal E↓	Spatial E↓	ST E↓
	omegez,	E	Н	E H		E	Н	E	Н	E	Н			Open-	Source Model	S		*2	Open-So	ource Models		
Open-Source Models													LLAVA-OneVision (0.5B)	0.48	10.16	0.0102	227	<u>**</u>	*	P	0.7010	0.0415
LLAVA-OneVision	0.5B	42.10	25.26	27.62	30.45	54.87	37.04	57.06	39.57	85.29	58.35	45.82	LLAVA-One Vision (0.3b) LLAVA-One Vision (7B)	5.77	14.09	0.0102		ī. 2	LLAVA-OneVision (0.5B)	0.7285	0.7212	0.8415
	7B							51.92			63.25		Qwen2-VL (2B)	8.72	17.93	0.2086	_	_	LLAVA-OneVision (7B)	0.7680	0.7750	0.8142
Owen 2 MI								1				1	Qwen2-VL (7B)	10.47	20.14	0.4119	-	_	Qwen2-VL (2B)	0.7729	0.7793	0.8127
Qwen2-VL	2B 7B	36.73	1000 00 000000	38.10		1	35.19	.00000000000000000000000000000000000000				070-90-0 0-900-0-03	VideoLLaMA2 (7B-8F)	6.25	19.94	0.2391	_	-	Qwen2-VL (7B)	0.7615	0.7647	0.8032
	/ Б	30.03	24.35			61.65	7.5-2-22 (2-27-27)		505-50 555		AS SUPPLY DATE OF		VideoLLaMA2 (7B-16F)	6.87	18.69	0.2111	-	-	VideoLLaMA2 (7B-8F)	0.6225	0.6360	0.6896
VideoLLaMA2	2.0-7B-8F	42.54		J		59.37		1-0-3000-3000-3000-3000		24 50000000000	65.56	2007-010-010-010-010-01						VideoLLaMA2 (7B-16F)	0.7218	0.7383	0.7895	
	2.0-7B-16F	42.41	10.47	55.98	41.94	53.80	52.26	44.16	47.75	66.93	64.82	48.05	TraffiX-Qwen Baseline 0.5B						TraffiX-Qwen Baseline			
TraffiX-Qwen Baseline													Spatial-Pooling	34.99	50.44	2.5195	35.24	46.35			0.1003	0.000
	Spatial Pooling	75.54	68.47	85.31	75.82	83 92	81.26	79.95	59 73	93.06	5 85.37	78.84	MultiRes Spatial-Pooling	34.91	50.26	2.4306	35.20	45.75	0.5B-Spatial-Pooling	0.1220	0.1892	0.2600
Baseline-0.5B (Ours)	MultiRes Spatial-Pooling	(A. Sarkovaris C. 18)		86.10		5,500 -000 -000 -000	79.59	03/07/03/03/03/03		WARRENCE CONTROL CONTR	***************************************	MultiRes Token-Pruning	35.07	50.79	2.5730	35.30	46.48	0.5B-MultiRes-Spatial-Pooling	0.1211	0.1894	0.2607	
	MultiRes Token-Pruning	76.61	1000014000 1100000	Seed to the second state of the Seeding			78.60		0 2 0		85.27		MultiRes Temporal-Pooling	35.63	51.00	2.5464	35.77	47.17	0.5B-MultiRes-Token-Pruning	0.1230	0.1934	0.2650
	MultiRes Temporal-Pooling	F112-F212-F212-F212-F212-F212-F212-F212	230000-00 250000	The state of the state of the state of		100000000000000000000000000000000000000		100000000000000000000000000000000000000		10.000.000.000.000.000.000		FOUNDAME 201 40		TraffiX-Qwen Baseline 7B					0.5B-MultiRes-Temporal-Pooling	0.1228	0.1912	0.2629
	Spatial Pooling					72					2 85.27		Spatial-Pooling	36.74	52.04	2.5613	36.42	47.32	7B-Spatial-Pooling	0.1083	0.1737	0.2382
	MultiRes Spatial-Pooling	17-15-15-15-15-15-15-15-15-15-15-15-15-15-		I The Section		and the same of the same		la Constitution of the con		A Charles Show	85.51	0 2000 10 00000000	MultiRes Spatial-Pooling	37.60	53.26	2.6113	37.31	49.16	7B-MultiRes-Spatial-Pooling	0.1136	0.1822	0.2493
Baseline-7B (Ours)	MultiRes Token-Pruning									100000000000000000000000000000000000000	85.48	7419/2004 - 4 1400 2409 S	MultiRes Token-Pruning	37.83	52.31	2.6162	36.56	47.80	7B-MultiRes-Token-Pruning	0.1152	0.1748	0.2454
	MultiRes Temporal-Pooling										8 85.34		MultiRes Temporal-Pooling	0.7000 F (0.701 0.007 0.004	52.58	2.4236	36.85	48.53	7B-MultiRes-Temporal-Pooling	0.1166	0.1790	0.2496

- Multi-resolution visual sampling strategies yield moderate improvements in multi-choice QA and referred object captioning, but introduce challenges for spatio-temporal grounding due to inter-frame association ambiguities.
- Tasks requiring fine-grained 3D spatial and temporal reasoning, including object positioning and motion analysis, remain challenging for current Visual LLMs.
- Accurate spatio-temporal object localization and cross-frame object association remain difficult for visual LLMs in dynamic traffic environments.



Qualitative



Referred Object Captioning Q: What are the main characteristics of [(c1, 0.01, 0.42, 0.87), (c1, 0.97, 0.58, 0.17)]? GT: [(c1, 0.01, 0.42, 0.87), (c1, 0.97, 0.58, 0.17)] is a dark purple coupe with a sleek design and distinctive rounded roofline. It features sporty alloy wheels and a modern rear light design.

purple coupe with a sleek design and distinctive round roofline. It features sporty alloy wheels and a modern light design.

Baseline-0.5B: [(c1, 0.01, 0.42, 0.87), (c1, 0.97, 0.58, is a black hatchback vehicle with a sleek design. It feat distinctive rear lights and a modern tailgate. The vehic compact size, making it suitable for urban environmenth BLEU-4: 0.33 ROUGE-L: 0.61 METEOR: 0.38 SPICE

Baseline-7B: [(c1, 0.01, 0.42, 0.87), (c1, 0.97, 0.58, 0 a black SUV with a sleek design. It features distinctive wheels and has a modern aesthetic, characterized by lines and a spacious appearance.

BLEU-4: 0.39 ROUGE-L: 0.69 METEOR: 0.42 SPICE

Q: What are the main features of [(c4, 0.01, 0.21, 0.21, 0.22]]

