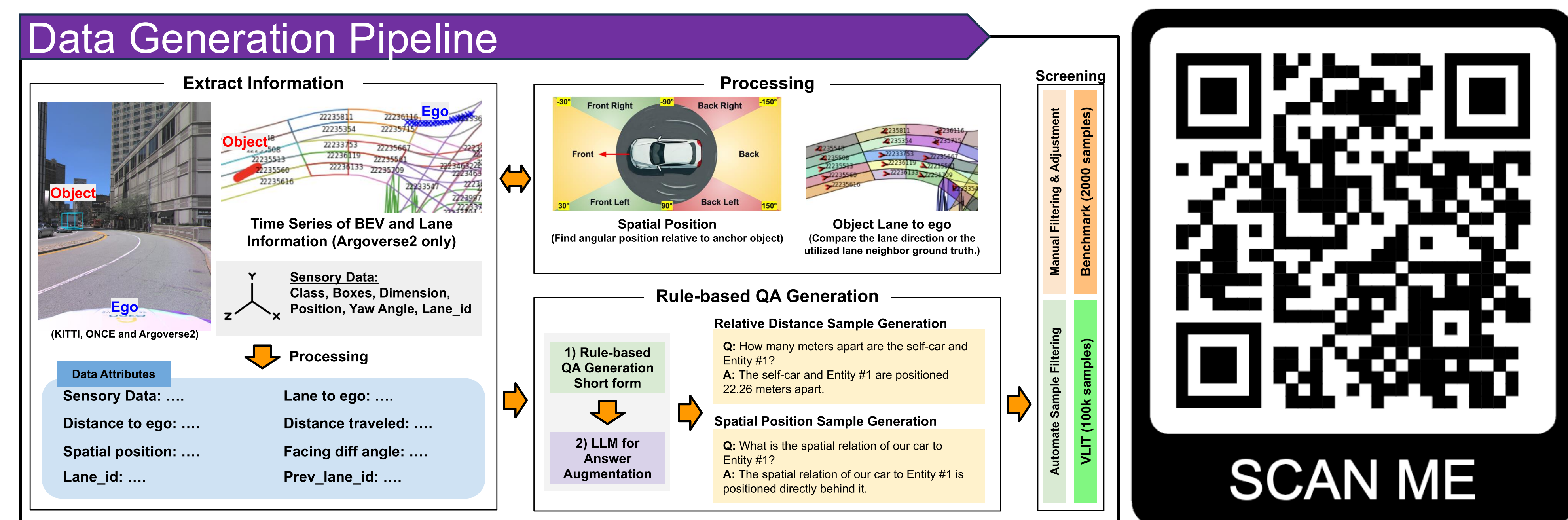


# TB-Bench: Training and Testing Multi-Modal AI for Understanding Spatio-Temporal Traffic Behaviors from Dashcam Images/Videos



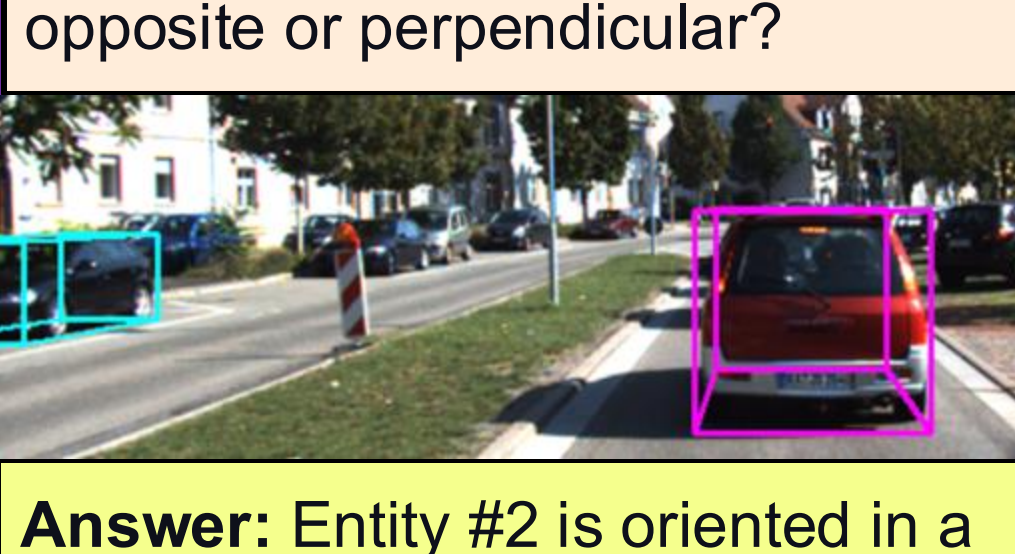
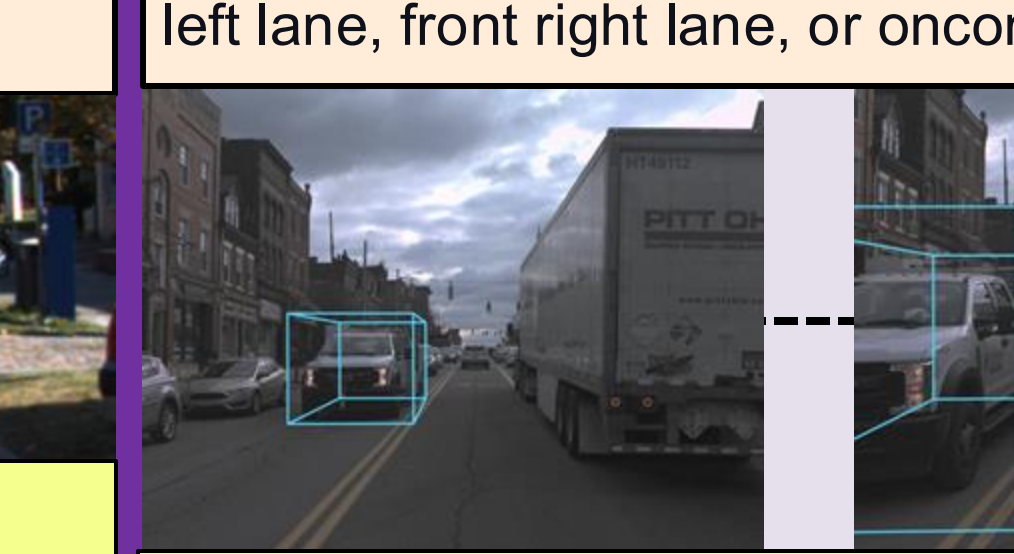
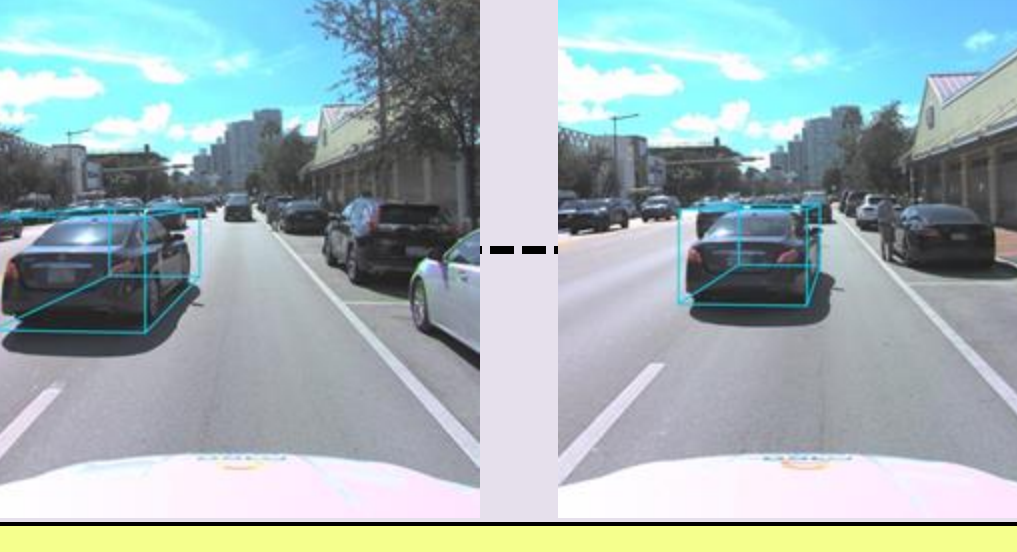
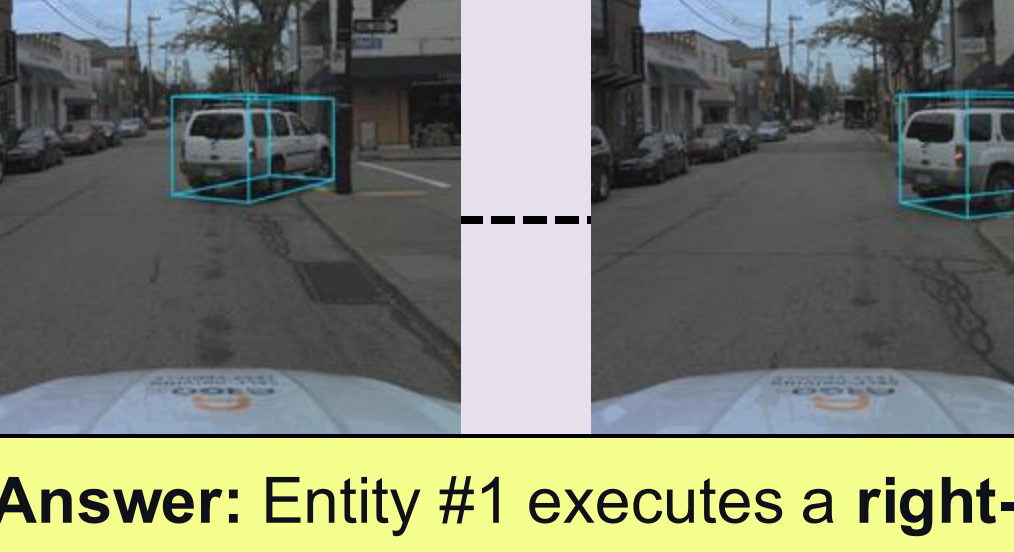
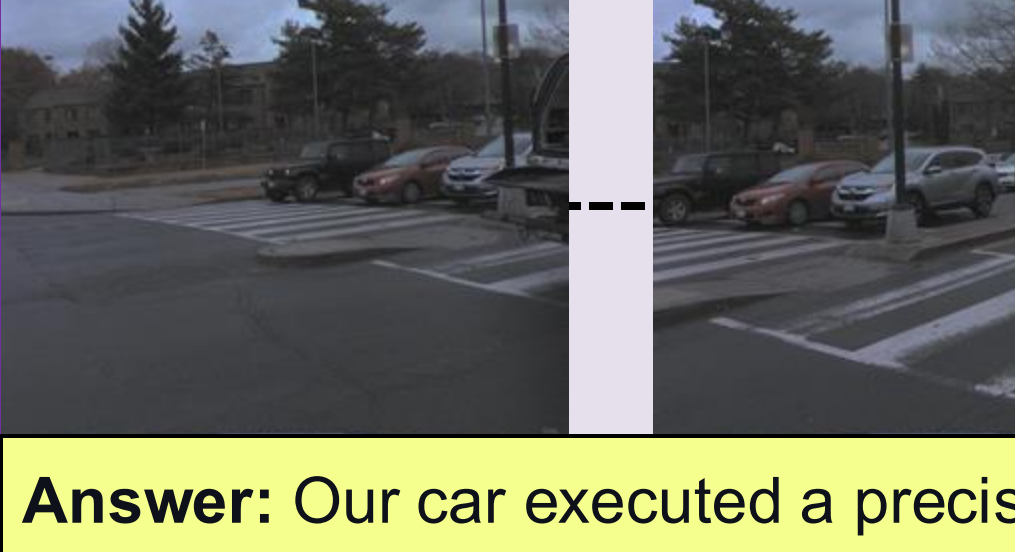
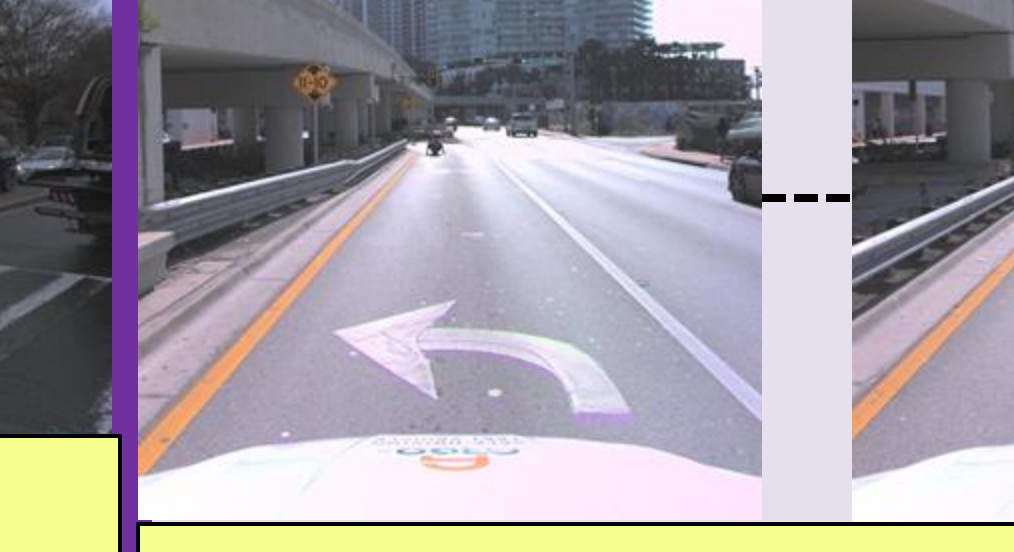
Korawat Charoenpitaks<sup>1</sup>, Van-Quang Nguyen<sup>2</sup>, Masanori Suganuma<sup>1</sup>, Kentaro Arai<sup>3</sup>, Seiji Totsuka<sup>3</sup>, Hiroshi Ino<sup>3</sup>, Takayuki Okatani<sup>1,2</sup>  
<sup>1</sup>Tohoku University <sup>2</sup>RIKEN AIP <sup>3</sup>DENSO Corporation

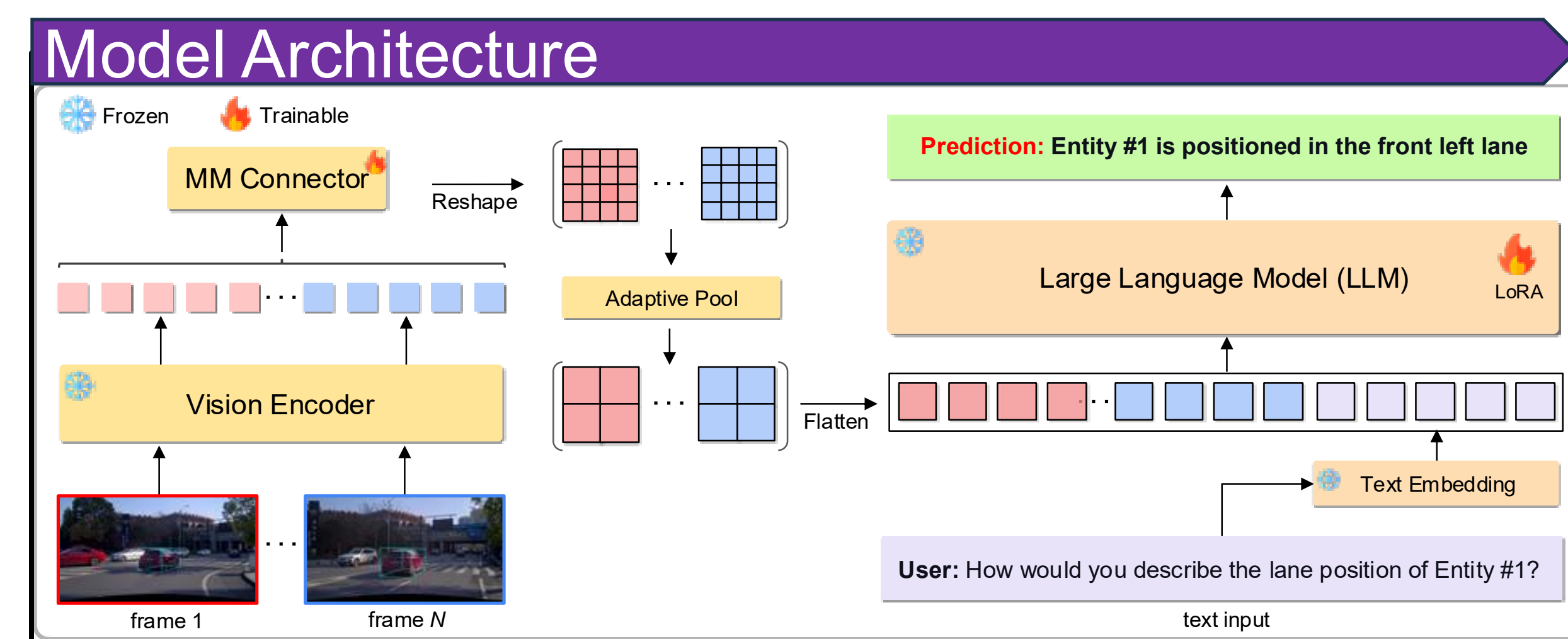
## Motivation & Tasks

- MLLMs in the autonomous driving domain struggle with **spatiotemporal understanding**, lacking traffic-specific data and dedicated benchmarks.
- Existing benchmarks focus on **static spatial tasks** (i.e., object detection) or single-image evaluations, and thus fail to capture **temporal aspects** of traffic behavior.
- Contributions:**
  - Closing the dataset gap:** Create eight spatiotemporal tasks covering critical traffic behaviors\* and provide two training sets (TB-100k and TB-250k) and the TB-Bench benchmark.
  - Findings:** Off-the-shelf proprietary models exhibit poor spatiotemporal performance; fine-tuning substantially improves results.
  - Cross-dataset transfer:** Co-training with our training set boosts other performance.



## Do MLLMs Have Spatio-Temporal Capabilities?

Relative Distance (RD)	Spatial Reasoning (SR)	Orientation Reasoning (SR)	Lane relative to ego-car (EGO-LANE)
<b>Question:</b> How many meters away is Entity #2 from Entity #1?  <b>Answer:</b> Entity #2 is situated <b>15.53 meters</b> away from Entity #1.	<b>Question:</b> What is the spatial position of Entity #1 relative to Entity #2?  <b>Answer:</b> Entity #1 is positioned at the <b>back right</b> relative to Entity #2.	<b>Question:</b> What is the orientation of Entity #2 relative to Entity #1, similar, opposite or perpendicular?  <b>Answer:</b> Entity #2 is oriented in a direction that is <b>diametrically opposed</b> to the orientation of Entity #1.	<b>Question:</b> How would you describe the lane position of Entity #1? Options: front lane, front left lane, front right lane, or oncoming traffic lane.  <b>Answer:</b> Entity #1 is positioned in the lane designated for <b>oncoming traffic</b> .
<b>Other lane change (OBJ-LANE)</b> <b>Question:</b> How would you describe the driving scene involving Entity #1? Please explain, focusing on the vehicle's lane change maneuver.  <b>Answer:</b> Entity #1 smoothly performs a <b>right lane change</b> .	<b>Other object turning (OBJ-TURN)</b> <b>Question:</b> How would you describe the driving scene involving Entity #1? Please explain, focusing on the vehicle's turning maneuver.  <b>Answer:</b> Entity #1 executes a <b>right-turn</b> maneuver, smoothly navigating the vehicle in a clockwise direction.	<b>Ego-vehicle turning (EGO-TURN)</b> <b>Question:</b> How would you describe the driving scene involving our car? Please explain, focusing on our car's turning maneuver.  <b>Answer:</b> Our car executed a precise <b>right-turn</b> maneuver, smoothly navigating the corner with adept control.	<b>Ego-vehicle traverse distance (EGO-TRA)</b> <b>Question:</b> How far has our car driven and what kind of steering maneuver did it perform in the current scene?  <b>Answer:</b> The car has driven <b>11.97 meters</b> with a straight steering maneuver.



## Data Statistics

Statistics of TB-Bench, TB-100k, and TB-250K. Source datasets: K (KITTI), O (ONCE), Arv2 (Argoverse2).

Task Type	Sources/ Frames	TB-Bench	TB-250k	TB-100k
<b>Spatial Information:</b>				
Relative Distance	[K, O]/ 1	250	35k	10k
Spatial Reasoning	[K, O]/ 1	250	70k	30k
Orientation Reasoning	[K, O]/ 1	250	70k	30k
<b>Object Behavior:</b>				
Other Lane to Ego	[Arv2]/ 8	250	50k	20k
Other Lane Changing	[Arv2]/ 8	250	1.5k	1.5k
Other Turning	[Arv2]/ 8	250	1.5k	1.5k
<b>Ego Behavior:</b>				
Ego Turning	[Arv2]/ 8	250	1.5k	1.5k
Ego Traverse Distance	[Arv2]/ 8	250	25k	15.5k
<b>Total</b>		2000	254k	110k

## Results on TB-Bench

Model	RD	SR	OR	EGO-LANE	OBJ-LANE	OBJ-TURN	EGO-TURN	EGO-TRA	Avg
<b>Zero-shot:</b> Bunny-v1.1-4B	24.4	20.4	19.6	28.4	16.0	20.0	34.4	0.0	<b>20.4</b>
<b>Zero-shot:</b> Mini-InternVL2-1B-DriveLM	0.0	31.2	20.0	28.4	24.8	47.2	41.6	0.0	<b>24.2</b>
<b>Zero-shot:</b> GPT-4o-2024-08-06	8.4	32.0	40.8	54.4	39.6	43.2	40.4	16.0	<b>34.4</b>
<b>Fine-tuned:</b> Ours with TB-100k	80.4	74.8	88.8	93.6	65.2	76.4	80.0	60.4	<b>77.5</b>
<b>Fine-tuned:</b> Ours with TB-250k	91.2	83.2	94.8	99.6	69.6	80.4	82.8	78.8	<b>85.1</b>

More than 50%

## Cross-dataset results

Metric	Standard training (only BDD-X)	With co-training (Sampling: BDD-X=20, TB-100k=1)
Speed RMSE↓	1.40	1.38
Speed A <sub>0.1</sub> ↑	26.1	26.3
Speed A <sub>0.5</sub> ↑	55.7	57.6
Turning RMSE↓	11.2	11.3
Turning A <sub>0.1</sub> ↑	44.2	44.5
Turning A <sub>0.5</sub> ↑	62.2	63.7

\*Based on Pre-crash Scenarios typology from the National Automotive Sampling System (NASS)