



**VIP3D**  
Visual Trajectory Prediction 3D



# ViP3D: End-to-end Visual Trajectory Prediction via 3D Agent Queries



Junru Gu\*



Chenxu Hu\*



Tianyuan Zhang



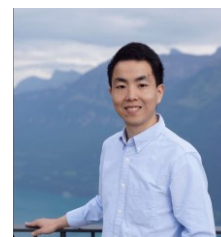
Xuanyao Chen



Yilun Wang



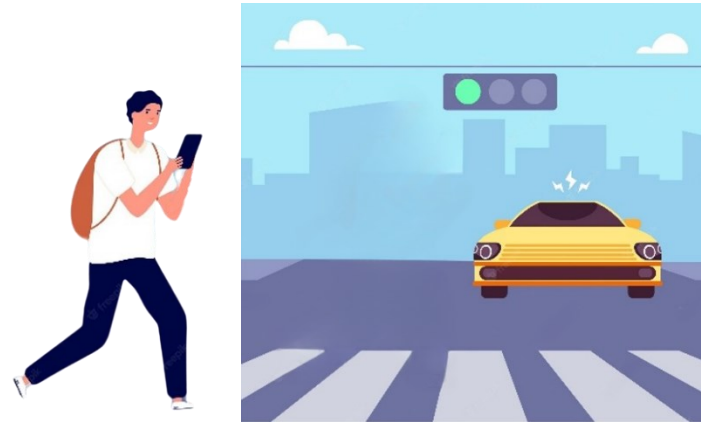
Yue Wang



Hang Zhao

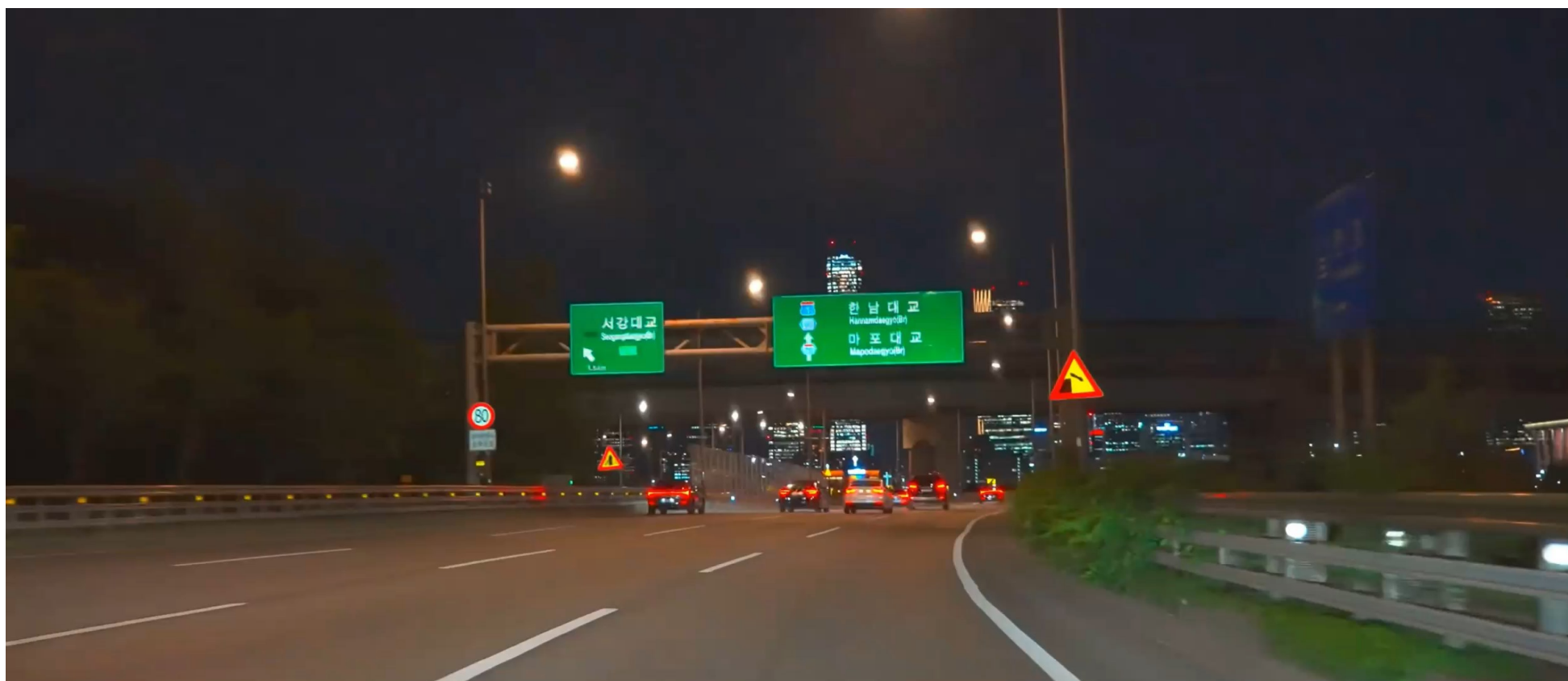
# Motivation: Why Visual?

- Take advantage of **rich visual information** for motion prediction.



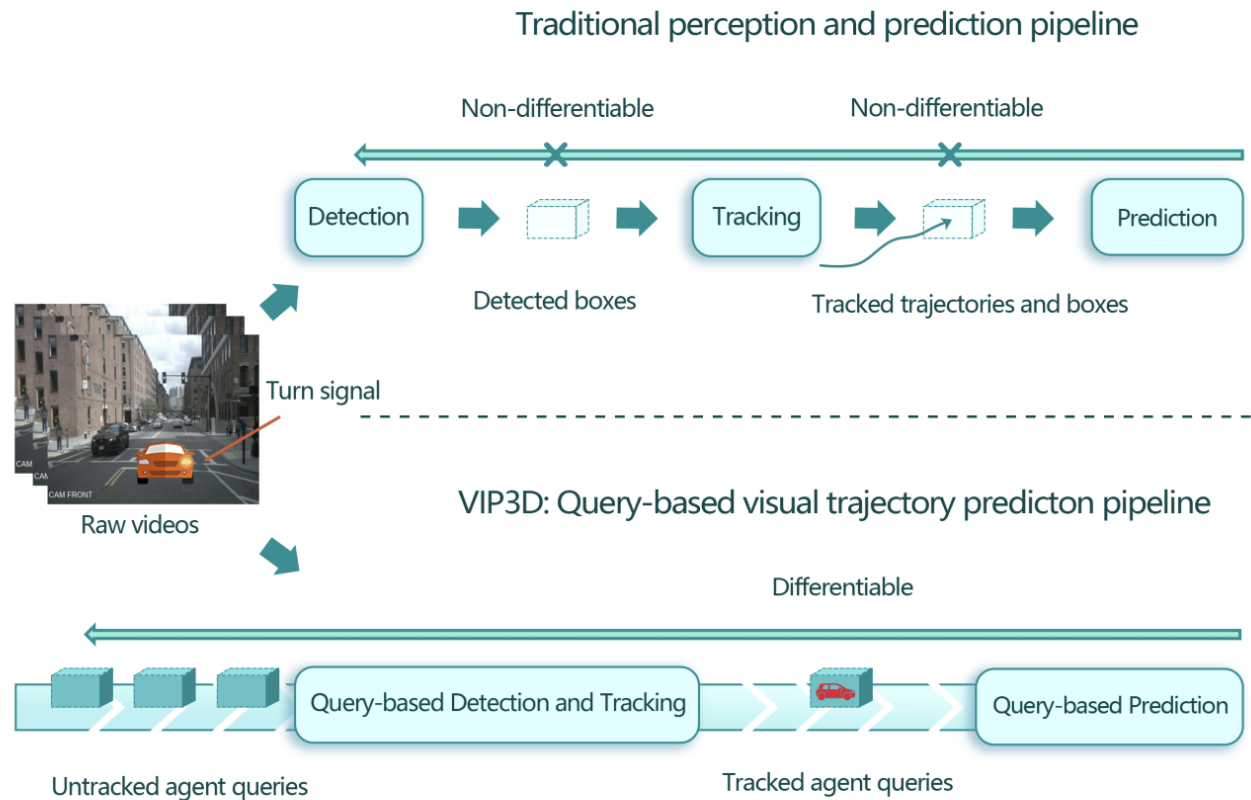
# Motivation: Why Visual?

- Take advantage of **rich visual information** for motion prediction.



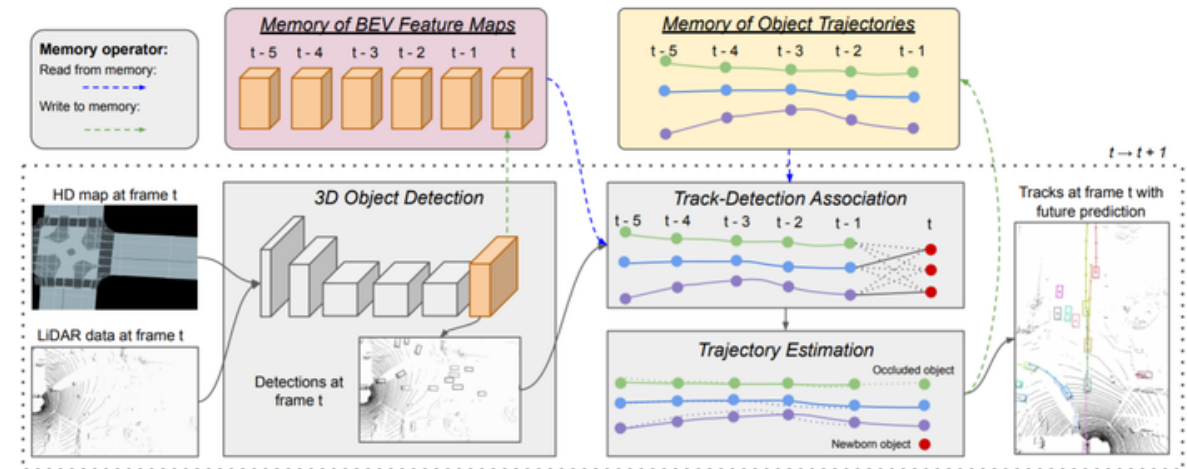
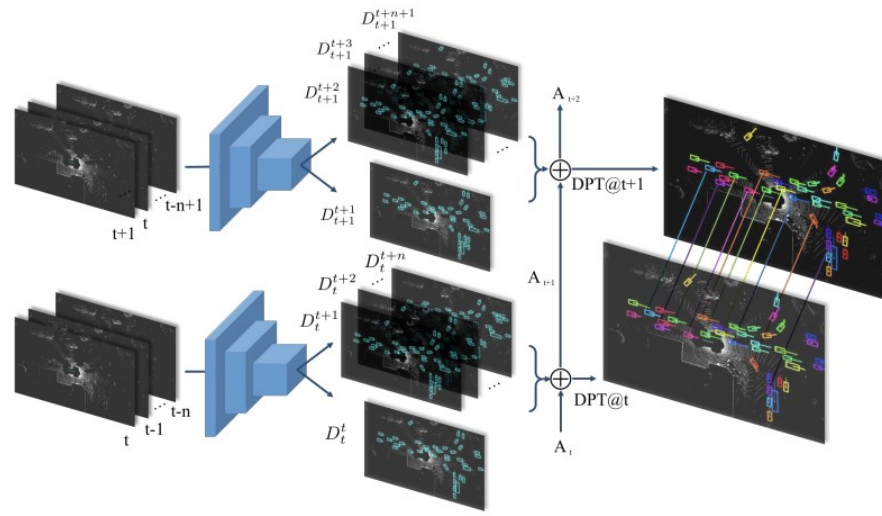
# Motivation: Why End-to-End?

- Avoid **error accumulation** and **data distribution shift**.



# Previous End-to-End Prediction Methods

- FaF, CVPR 2018
  - The first CNN model for **joint detection and prediction**, from **LiDAR** inputs.
- PnPNet, CVPR 2020
  - A **tracking-in-the-loop** model for **LiDAR**-based trajectory prediction.



# Previous End-to-End Prediction Methods

- 🙄 • Unable able to leverage the abundant **visual information** from cameras
- 🙄 • Use **convolutional feature maps** as their intermediate representations



Suffering from **non-differentiable operations**

- E.g. Non-maximum suppression in object detection
- E.g. Object association in tracking

# Previous End-to-End Prediction Methods

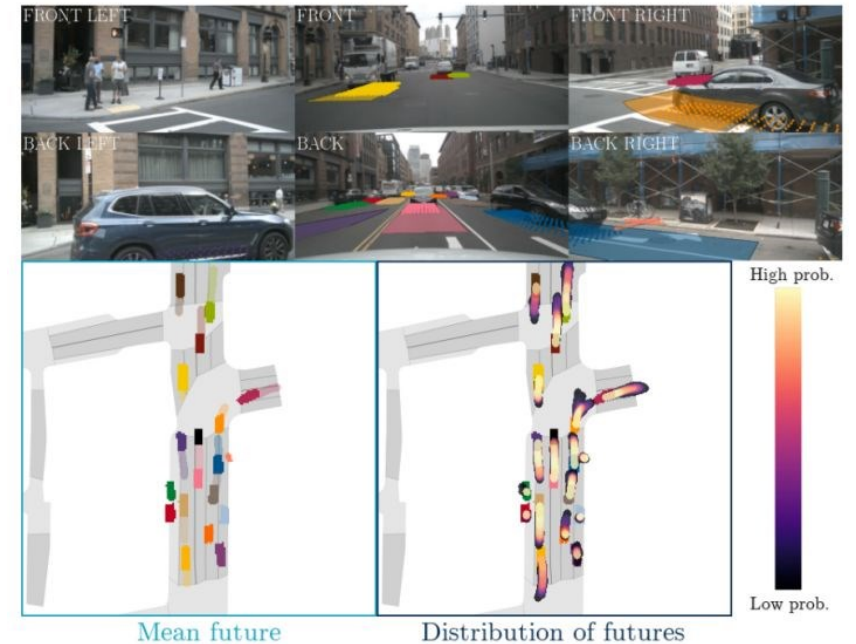
- FIERY, ICCV 2021
  - Predicts future **BEV occupancy heatmaps** from **visual** data.



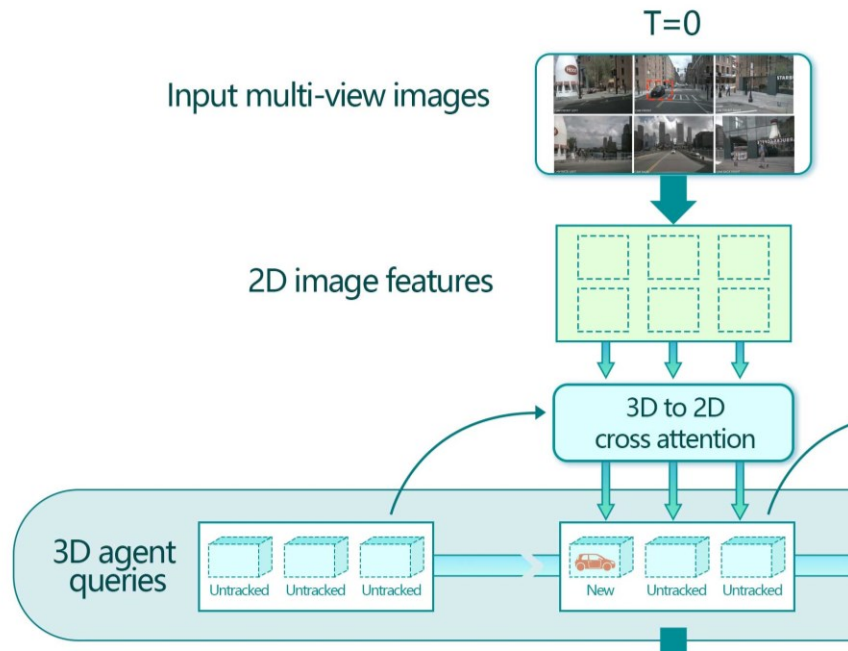
- **Not modularized, non-interpretable**, thus not engineering-wise friendly



- Output representation is **not compatible** with downstream modules like motion planning

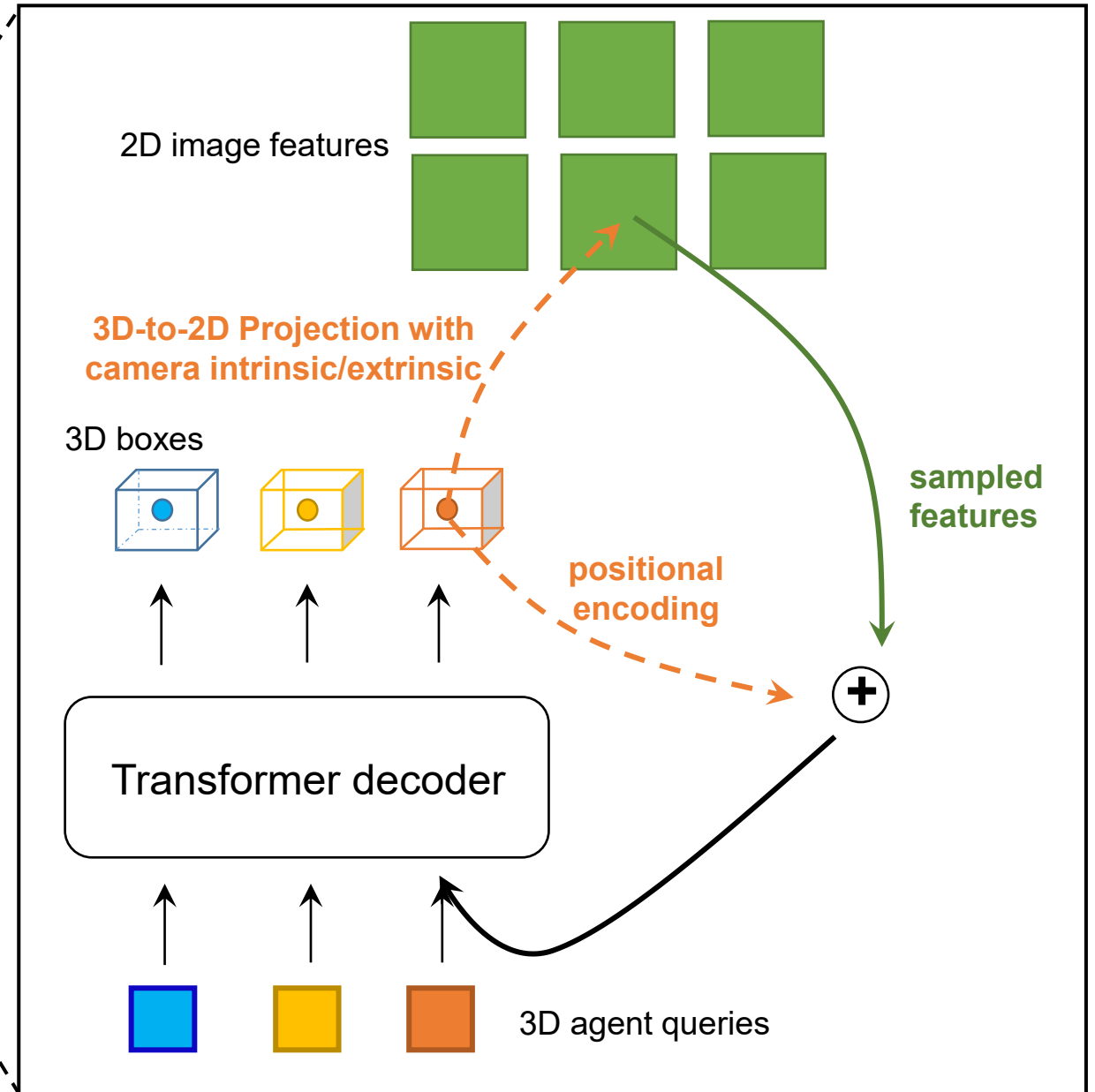
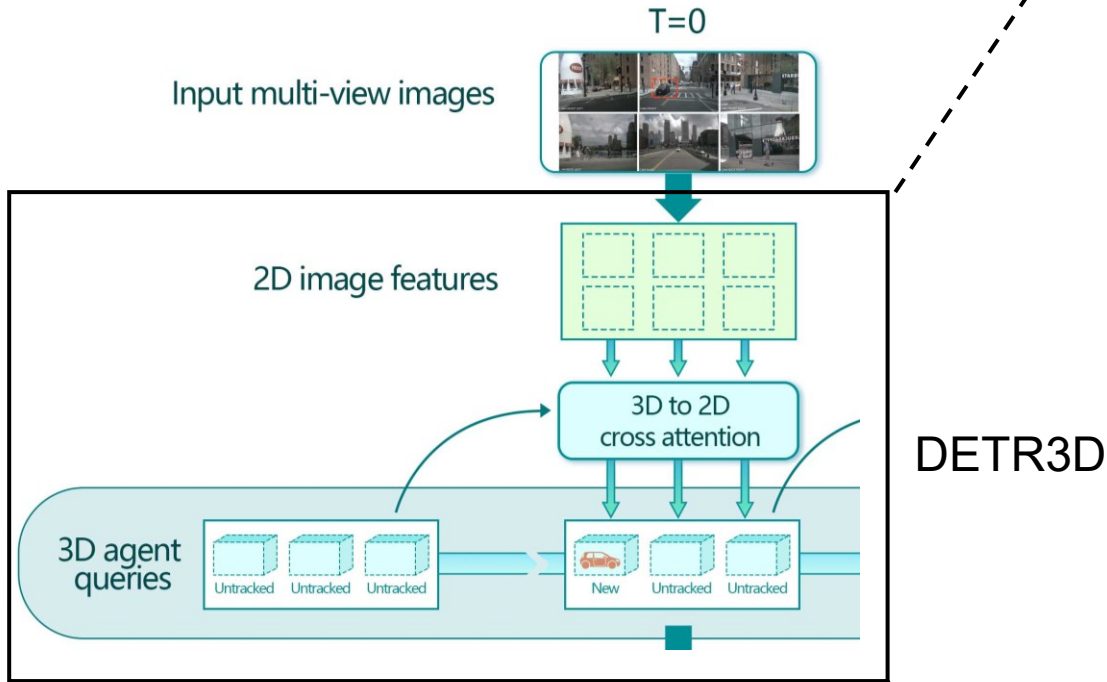


# Model Architecture



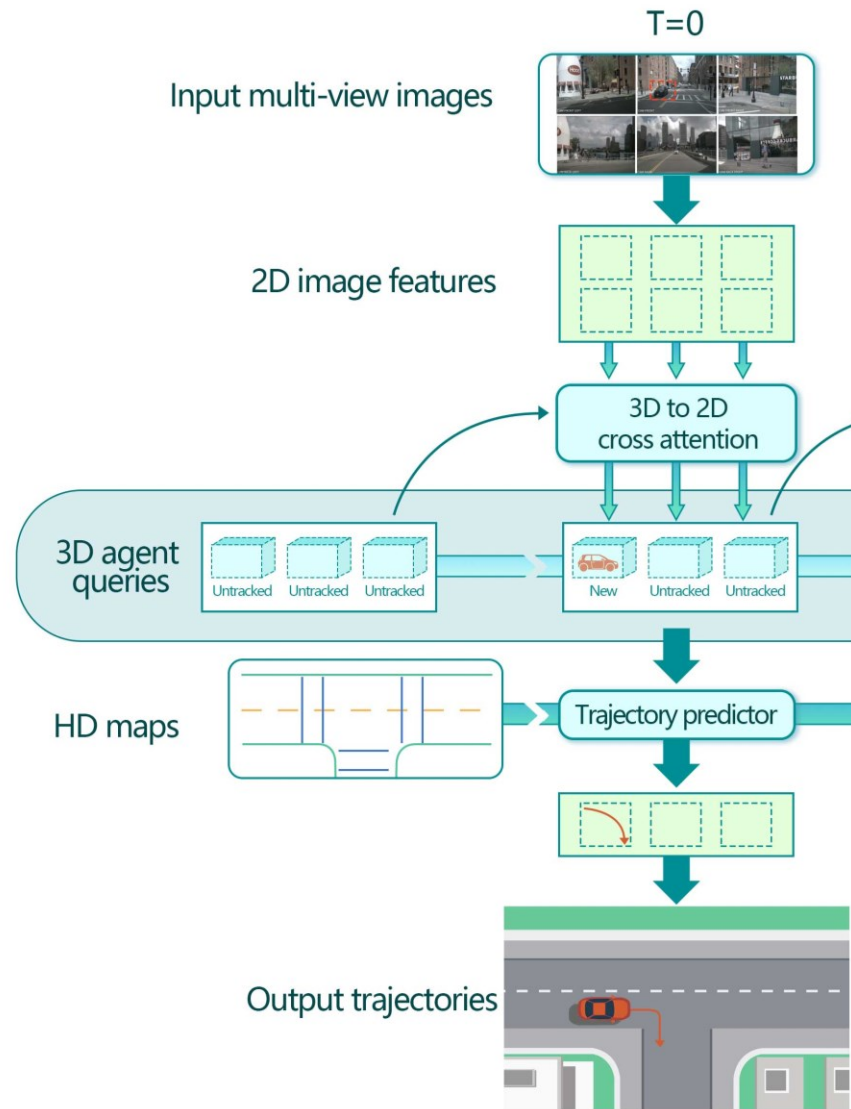


# Model Architecture

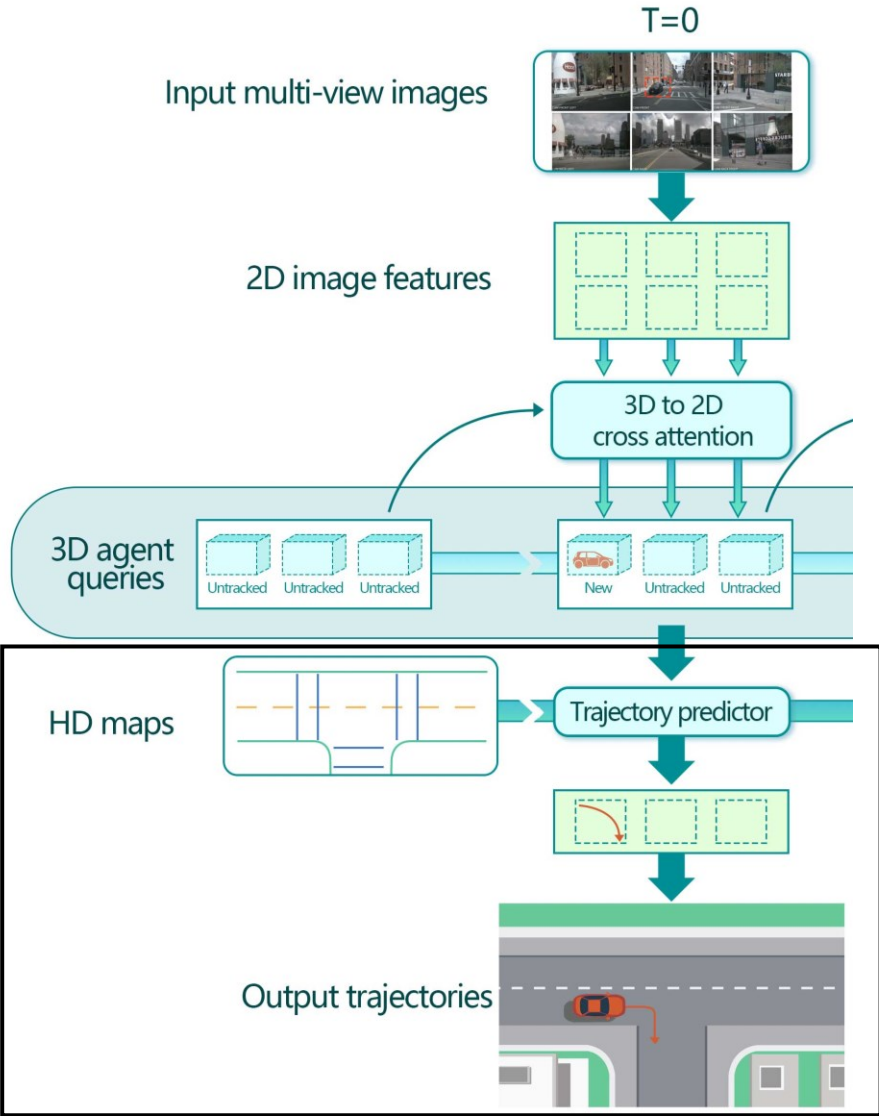


DETR3D: 3D Object Detection from Multi-view Images via 3D-to-2D Queries, Yue Wang, Vitor Guizilini, Tianyuan Zhang, Yilun Wang, Hang Zhao, Justin Solomon, CoRL 2021

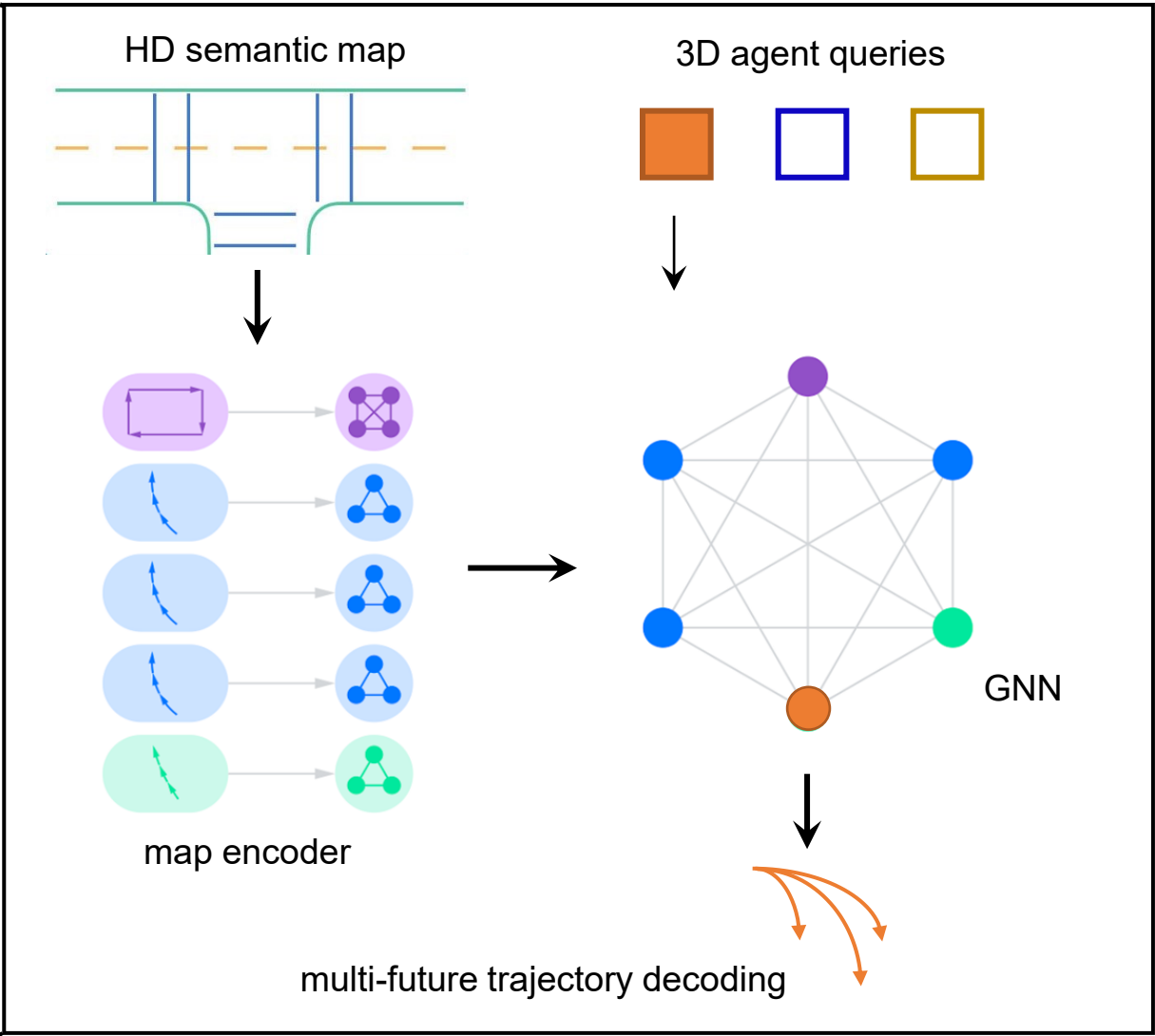
# Model Architecture



# Model Architecture

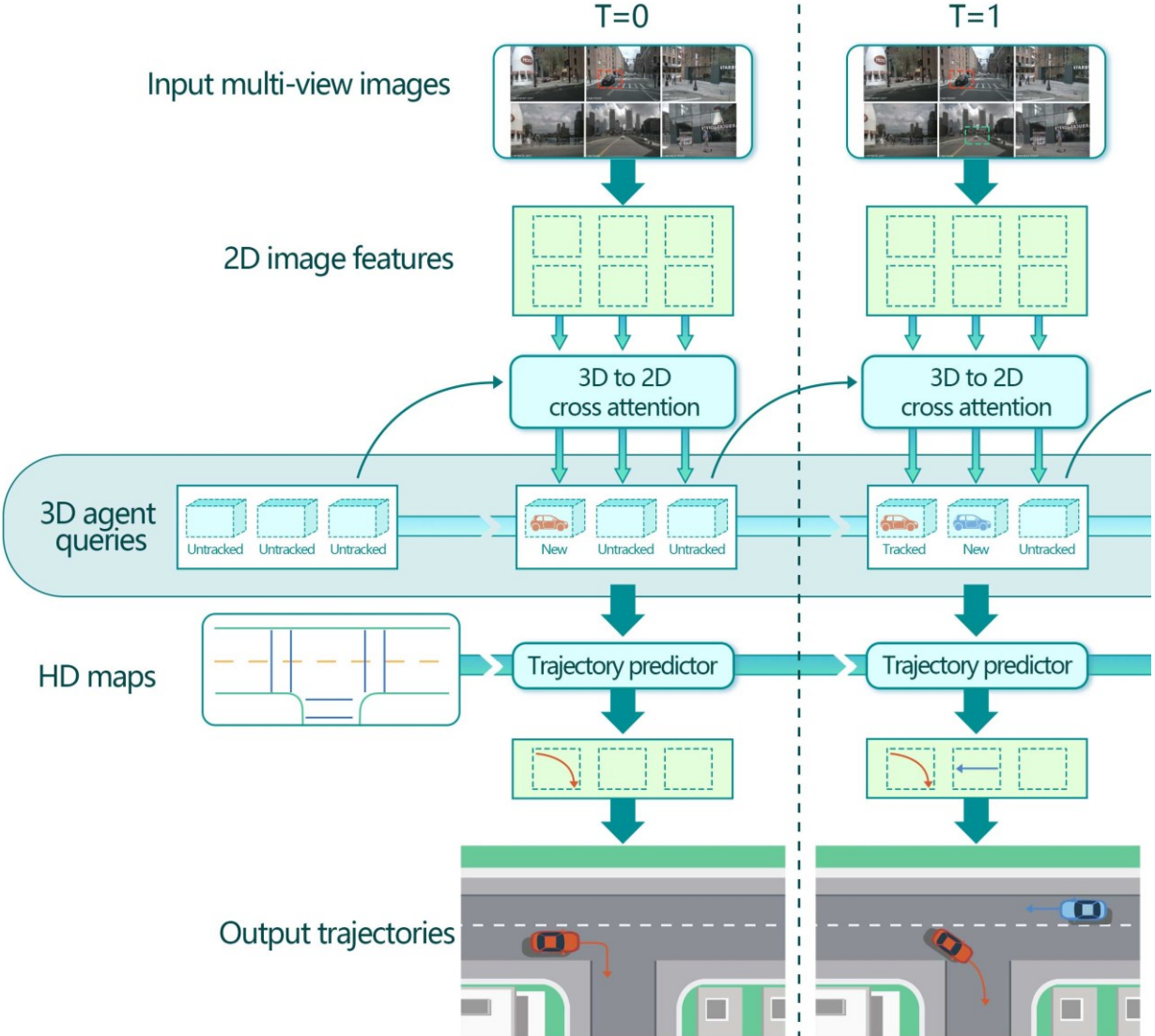


VectorNet

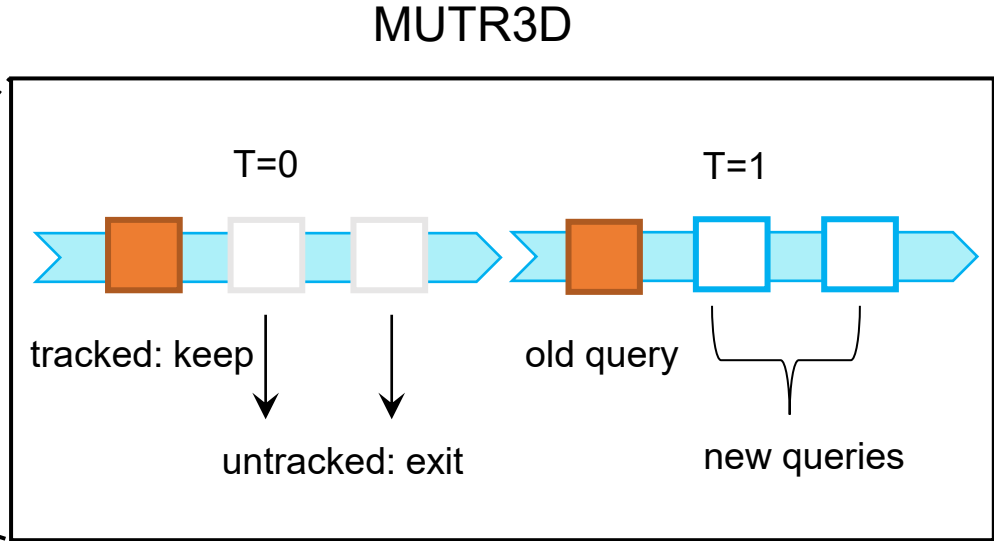
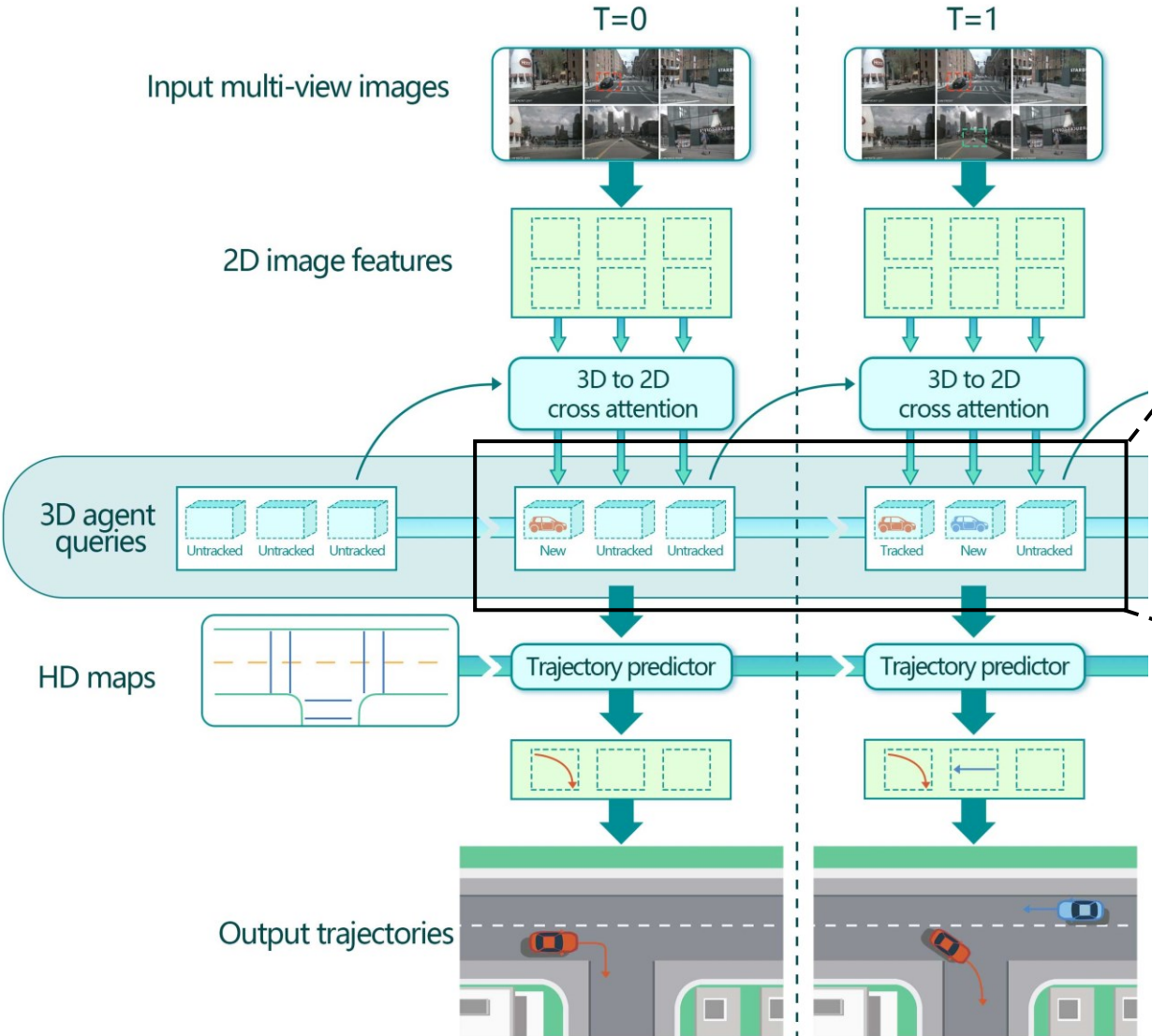


Vectornet: Encoding HD maps and agent dynamics from vectorized representation, Jiyang Gao, Chen Sun, Hang Zhao, Yi Shen, Dragomir Anguelov, Congcong Li, Cordelia Schmid, CVPR 2020

# Model Architecture

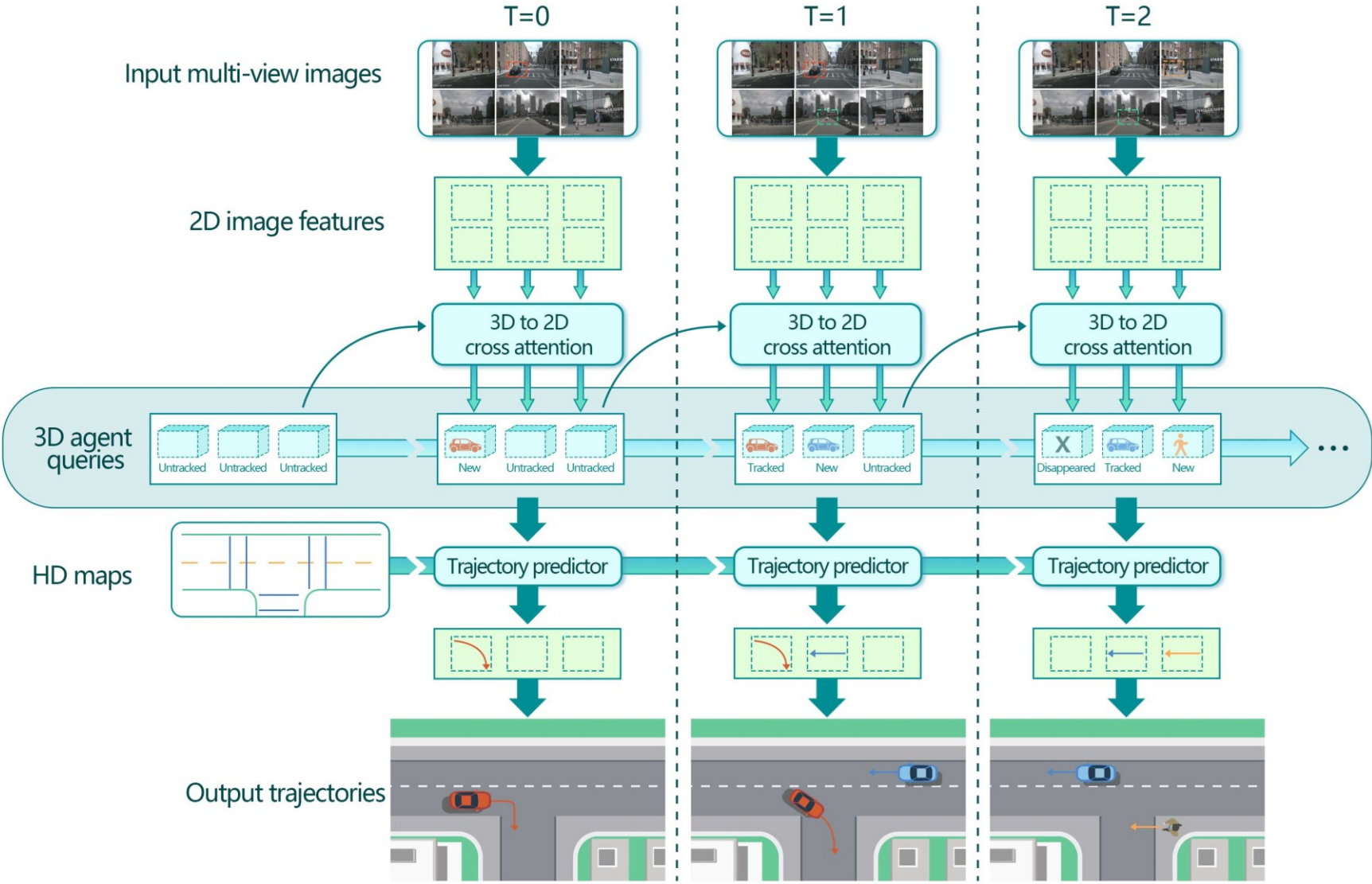


# Model Architecture



MUTR3D: A Multi-camera Tracking Framework via 3D-to-2D Queries, Tianyuan Zhang, Xuanyao Chen, Yue Wang, Yilun Wang, Hang Zhao, CVPRW 2022

# Model Architecture



**Streaming**  
 Queries initialized, updated and discarded over time...

# Experiments

- Comparison with traditional pipeline and the SOTA model PnPNet.

		Traditional		PnPNet-vision [30]		ViP3D (Ours)
Architecture	detector	DETR3D		DETR3D		DETR3D
	detector-tracker interface	boxes		boxes		queries
	tracker	Kalman Filter	CenterPoint	Kalman Filter	CenterPoint	query-based
	tracker-predictor interface	trajectories		cropped features		queries
	predictor	regression-based		regression-based		regression-based
Metrics	minADE↓	2.07	2.06	2.04	2.04	<b>2.03</b>
	minFDE↓	3.10	3.02	3.08	3.03	<b>2.90</b>
	MR↓	0.289	0.277	0.277	0.271	<b>0.239</b>
	EPA↑	0.191	0.209	0.198	0.213	<b>0.236</b>

# Ablation Study

- Differentiability is important
- Historical trajectories are no longer necessary, since agent queries captured these information implicitly

	Prediction inputs	Differentiable	minADE ↓	minFDE ↓	MR ↓	EPA ↑
	Agent trajectories	✗	2.30	3.33	0.282	0.186
	Agent trajectories + Agent queries	✗	2.20	3.19	0.274	0.211
ViP3D	Agent queries	✓	<b>2.03</b>	<b>2.90</b>	<b>0.239</b>	<b>0.236</b>

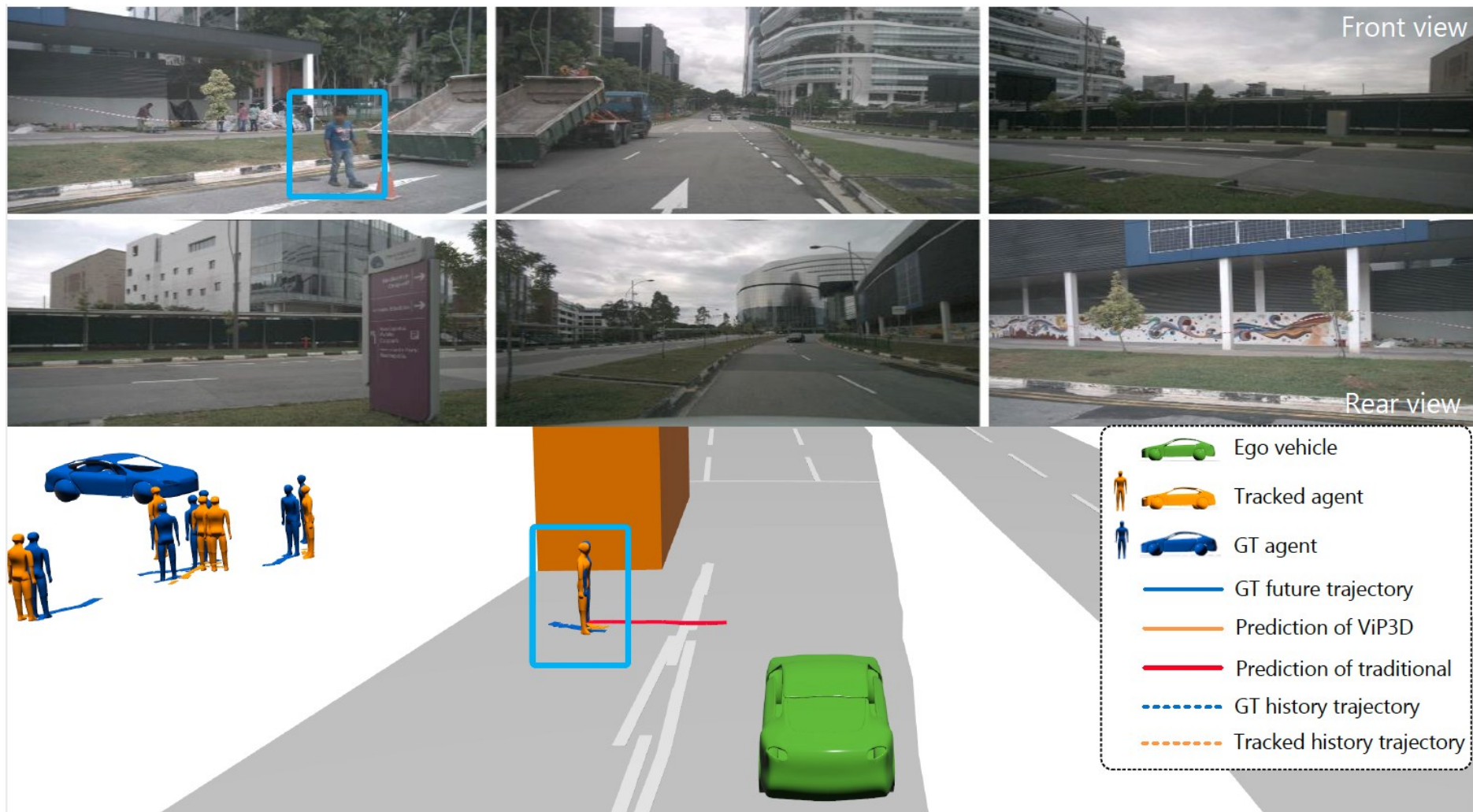


# Ablation Study

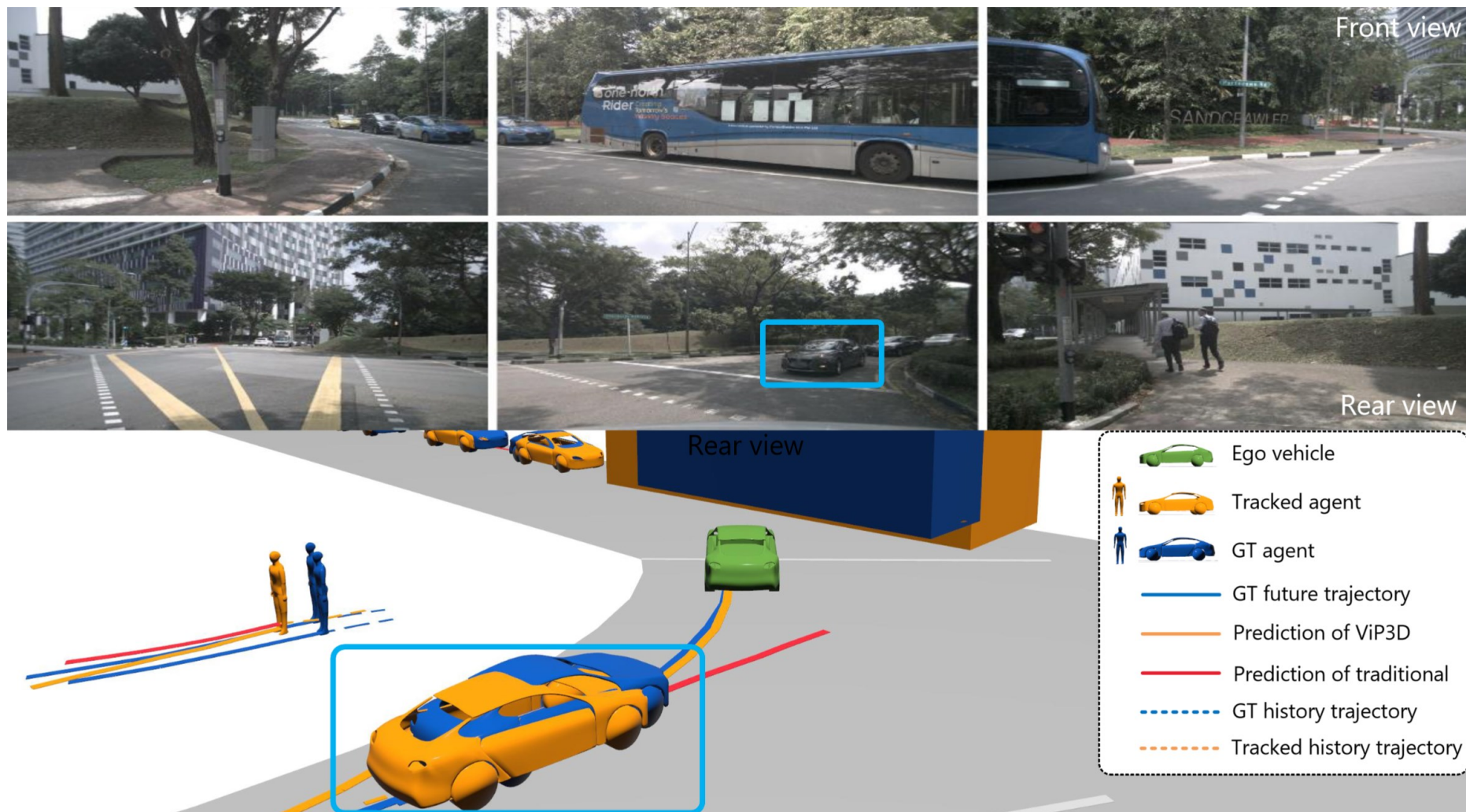
- ViP3D is compatible with multiple trajectory decoders.
- Such as TNT (goal-based) and HOME (heatmap-based).

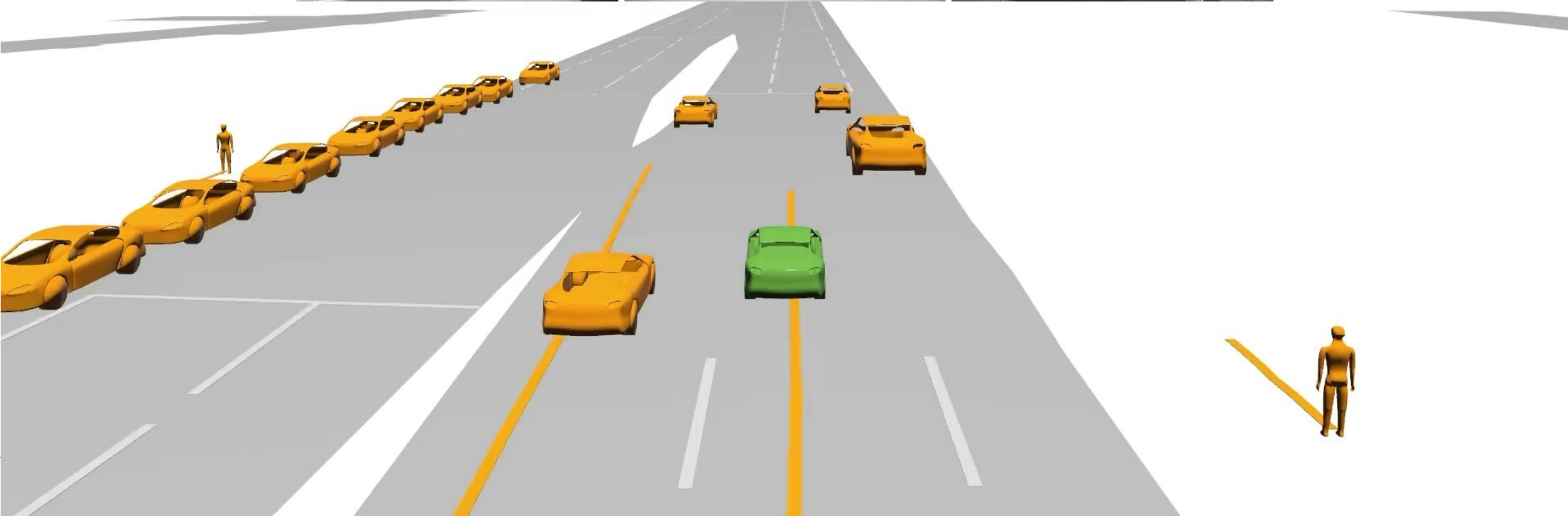
Decoder	Pipeline	mADE	mFDE	MR	EPA
Goal [60]	Traditional	2.50	3.93	0.266	0.195
	ViP3D	<b>2.24</b>	<b>3.33</b>	<b>0.238</b>	<b>0.219</b>
Heatmap [14]	Traditional	2.53	3.81	0.264	0.197
	ViP3D	<b>2.33</b>	<b>3.42</b>	<b>0.218</b>	<b>0.214</b>

# Qualitative Results



# Qualitative Results







MARS Lab  
THE END THANKS



Project Page