

THOMAS: TRAJECTORY HEATMAP OUTPUT WITH LEARNED MULTI-AGENT SAMPLING

Anonymous authors

Paper under double-blind review

ABSTRACT

In this paper, we propose THOMAS, a joint multi-agent trajectory prediction framework allowing for efficient and consistent prediction of multi-agent multi-modal trajectories. We present a unified model architecture for fast and simultaneous agent future heatmap estimation leveraging hierarchical and sparse image generation. We demonstrate that heatmap output enables a higher level of control on the predicted trajectories, allowing to incorporate additional constraints for tighter sampling or collision-free predictions in a deterministic way. However, we also highlight that generating scene-consistent predictions is a harder problem than just leveraging collisions, and design a scene-consistent recombination learned model that takes a set of predicted trajectories for each agent as input and outputs a consistent reordered recombination of the trajectories. We report our results on the Interaction multi-agent prediction challenge and rank 1st on the online test leaderboard.

1 INTRODUCTION

Motion forecasting is an essential step in the pipeline of an autonomous driving vehicle, transforming perception data into future prediction which are then leveraged to plan the future moves of the autonomous cars. The self-driving stacks needs to predict the future trajectories for all the neighbor agents, in a fast and coherent way.

The interactivity between agents plays an important role for accurate trajectory prediction. Agents need to be aware of their neighbors in order to adapt their speed, yield right of way and merge in neighbor lanes. To do so, different interaction mechanisms have been developed, such as social pooling Alahi et al. (2016); Lee et al. (2017); Deo & Trivedi (2018), graphs Salzmann et al. (2020); Zeng et al. (2021) or attention Mercat et al. (2020); Messaoud et al. (2020); Luo et al. (2020); Gao et al. (2020); Liang et al. (2020); Ngiam et al. (2021). These mechanisms allow agents to look at and share features with neighbors, to take them into account in their own predictions.

Multi-modality is another very important aspect of the possible future trajectories. A car can indeed chose to turn right or left, or decide to realise a certain maneuver in various ways. Uncertainty modeled as variance of Gaussians is insufficient to model these multiple cases, as it can only represent a continuous spread and can't show multiple discrete possibilities. Therefore, current state-of-the-art produces not one but K possible trajectories for each agent predicted, and most recent benchmarks Caesar et al. (2020); Chang et al. (2019); Zhan et al. (2019); Ettinger et al. (2021) include multi-modality in their metrics, taking only the minimum error over a predicted set of K trajectories.

However, up until very recently and the opening of multi-agent joint interaction challenges Ettinger et al. (2021); Zhan et al. (2019) 1.2, no motion forecasting prediction datasets were taking into account the coherence of modalities between different agents predicted at the same time. As a result, the first predicted modality of a given agent could crash with the first predicted modality of another agent without any check. Ettinger et al. (2021) introduces proposes a metric where models have to predict for two interacting agents jointly at each sample, and Zhan et al. (2019) 1.2 adds a general multi-agent track, where between 2 and 40 agents can be asked for joint prediction in a single sample.

Our THOMAS model encodes the past trajectories of all the agent presents in the scenes, as well as the HD-Map lanelet graph, and merges their information using self and cross attention. It then

predicts for each agent a sparse heatmap representing the future probability distribution at the last timestep. A deterministic sampling algorithm then iteratively selects the best K trajectory endpoints according to the heatmap for each agent. These endpoints are recombined to be scene-consistent by the COMBI model, and full trajectories are then generated for each endpoint.

Our contributions are summarized as follow:

- We present an efficient graph-based model enabling fast and efficient multi-agent future motion estimation
- We propose a variant heatmap sampling method that takes into account agent collisions
- We design a novel recombination model able to recombine the sampled endpoints in order to obtain scene-consistent samples across the agents

2 RELATED WORK

Learning-based models have quickly overtaken physics-based methods for trajectory prediction, as the sequential nature of trajectories is a logical application for recurrent architectures (Alahi et al., 2016; Altché & de La Fortelle, 2017; Lee et al., 2017; Mercat et al., 2020; Khandelwal et al., 2020), and convolutional layers can easily be applied to bird-view rasters of the map context (Lee et al., 2017; Tang & Salakhutdinov, 2019; Cui et al., 2019; Hong et al., 2019; Salzmann et al., 2020; Chai et al., 2020; Gilles et al., 2021b), benefiting from the latest progresses in computer vision. Surrounding HD-Maps, usually formalized as connected lanelets, can also be encoded using Graph Neural Networks (Gao et al., 2020; Liang et al., 2020; Zeng et al., 2021; Gilles et al., 2021a), in order to get a more compact representation closer to the trajectory space. Finally, some point-based approaches (Ye et al., 2021) can be applied in a broader way to trajectory prediction, as both lanes and trajectories can be considered as ordered set of points.

Multi-modality in prediction can be obtained simply through a multiple prediction head in the model (Cui et al., 2019; Liang et al., 2020; Ngiam et al., 2021; Deo et al., 2021). However some methods rather adopt a candidate-based approach, where potential endpoints are obtained either from anchor trajectories obtained through clustering (Chai et al., 2020; Phan-Minh et al., 2020) or a model-based generator (Song et al., 2021). Other approaches use a broader set of candidates from the context graph (Zhang et al., 2020; Zhao et al., 2020; Zeng et al., 2021; Kim et al., 2021) or a dense grid around the target agent (Deo & Trivedi, 2020; Gu et al., 2021; Gilles et al., 2021b;a). Another family of approaches use variational inference to generate diverse predictions, through latent variables (Lee et al., 2017; Rhinehart et al., 2018; Tang & Salakhutdinov, 2019; Casas et al., 2020) or GAN (Gupta et al., 2018; Rhinehart et al., 2018; Sadeghian et al., 2019) but the sampling of these trajectories is random and doesn't provide any probability value for each sample.

While very little work has directly tackled multi-agent prediction and evaluation so far, multiple methods hint at the ability to predict multiple agents at the same time (Liang et al., 2020; Zeng et al., 2021) even if they then focus on a more single-agent oriented framework. SceneTransformer (Ngiam et al. (2021)) repeats each agent features across possible modalities, and performs self-attention operations inside each modality before using a joint loss to train a model and evaluate on the WOMD (Ettinger et al., 2021) interaction track, but this evaluation focuses on dual agent prediction and doesn't broaden to a large number of agents at the same time. ILVM (Casas et al., 2020) uses scene latent representations conditioned on all agents to generate scene-consistent samples, but its variational inference doesn't provide confidence score for each modality. AIR² (Wu & Wu, 2021) extends Multipath (Chai et al., 2020) and produces a cross-distribution for two agents along all possible trajectory anchors, but it scales exponentially with the number of agents, making impractical for a real-time implementation that could encounter more than 10 agents at the same time.

3 METHOD

3.1 MODEL BACKBONE

Our goal is to predict the future F timesteps of A agents using their past history made of T timesteps and the HD-Map context. Similar to recent work (Zhao et al., 2020; Zeng et al., 2021; Gu et al.,

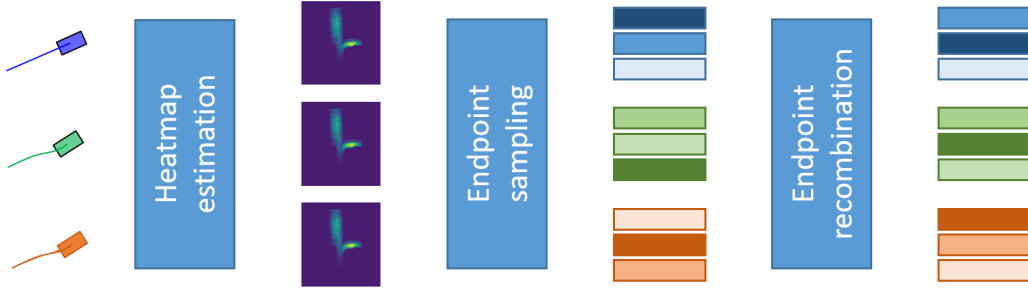


Figure 1: Illustration of the THOMAS multi-agent prediction pipeline

2021), we will divide the problem into goal-based prediction followed by full trajectory reconstruction. We first encode each agent trajectory and the HD-Map context graph into a common representation. We then decode a future probability heatmap for each agent in the scene, which we sample heuristically to maximize coverage. Finally, we re-combine the sampled endpoints into scene-consistent modalities across agents, and build the full trajectories for each agent.

3.1.1 GRAPH ENCODER

We use the same encoder as the GOHOME model Gilles et al. (2021a). The agent trajectories are encoded through *TrajEncoder* using a 1D CNN followed by a UGRU recurrent layer, and the HD-Map is encoded as a lanelet graph using a GNN *GraphEncoder* made of graph convolutions. The optional future trajectories are also encoded with a separate 1D CNN + UGRU encoder *FutureEncoder*. The future encoding is then concatenated to the past encoding in order to obtain a full encoding for each agent. We then run cross-attention *Lanes2Agents* to add context information to the agent features, followed by self-attention *Agents2Agents* to observe interaction between agents. The final result is an encoding F_a for each agent, where history, future, context and interactions have been summarized. This encoding F_a is used in the next decoder operations, but is also stored to be potentially used in modality re-ranking described in Sec. 3.2.2

3.1.2 HIERARCHICAL GRID DECODER

Our aim here is to decode each agent encoding into a heatmap representing its future probability distribution at prediction horizon T . Since we create this heatmap for each agent in the scene, the decoding process has to be fast so that it can be applied to a great number of agents in parallel. HOME (Gilles et al., 2021b) generates a heatmap through CNN operations, but these are costly and don't scale them with prediction range. DenseTNT (Gu et al., 2021) uses attention on dense gridpoints sampled around the lanes only, while GOHOME (Gilles et al., 2021a) create curvilinear rasters for a subselection of lanes before merging them together, but both these approaches neglects possible endpoints outside the drivable area so they can achieve reasonable inference times. We

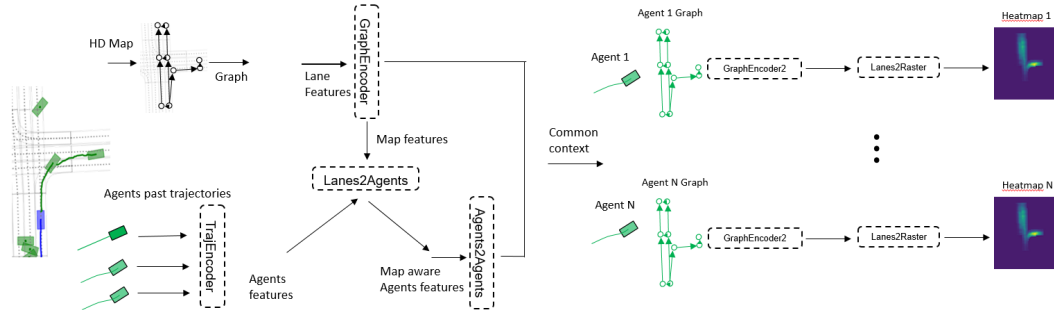


Figure 2: Model architecture for multi-agent prediction with shared backbone

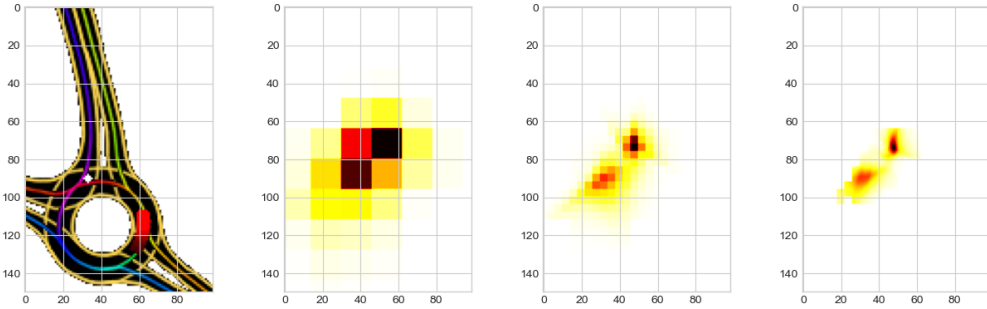


Figure 3: Hierarchical iterative refinement of the grid probabilities. First, the full grid is evaluated at a very low resolution, then the highest cells are up-sampled and evaluated at a higher resolution, until final resolution is reached.

inspire ourselves from these, with a few modifications so that we can predict endpoints anywhere on the map while smartly restricting computing to only part of the full image.

We leverage drivable area sparsity so that our model chooses the area to define with more precision while still having the opportunity to predict over the whole dense grid. We use hierarchical predictions at lower resolution levels first in order to process only a relatively low number of grid points. This hierarchical process is illustrated in Fig. 3.

We first predict a full dense grid probability at resolution 8mx8m by pixel. We then select the 16 highest ranking grid points, and upsample only these points to a 2mx2m intermediary grid. We repeat the process to select the top 64 points of this grid and upsample them to a final 0.5mx0.5m grid for the definitive heatmap prediction. The total prediction range is set to 192m. At each step, the grid points features are computed by a 2-layer MLP applied on the point coordinates, they are then concatenated to the agent encoding followed by a linear layer, and finally refined by a 2-layer cross-attention on the graph lane features. This hierarchical process allows the model to only operate on $24 \times 24 + 16 \times 4 \times 4 + 64 \times 4 \times 4 = 1856$ grid points instead of the 147 456 available.

3.1.3 FULL TRAJECTORY GENERATION

From each heatmap, we decode K end points using the same MR optimization algorithm as Gilles et al. (2021b). We then generate the full trajectories for each end point the same way, using a fully-connected MLP.

3.2 MULTI-AGENT CONSISTENT PREDICTION

The difficulty in multimodal outputs for multi-agents prediction comes from having coherent modalities between each agent. Since the modalities are considered scene-wise, the first predicted modality has to match with the first prediction of the other agents, and so on. Moreover, these modalities must not collide with each other, as they should represent realistic scenarios.

3.2.1 INITIAL COLLISION-FREE SAMPLING

We use the same sampling algorithm as Gilles et al. (2021a) based on MR optimization, but add a sequential iteration over the agents for each modalities. For a single modality k , we predict the possible endpoint of a first agent a by taking the maximum accumulated predicted probability under an area of radius r . We then not only set to zero the heatmap values of this agent heatmap \mathcal{I}_k^a around the sampled location so not to sample it in the next modalities k' , but we also set to zero the same area on the heatmaps $\mathcal{I}_k^{a'}$ of the other agents a' on the same modality k , so that these other agents cannot be sampled at the same position for this modality. This way, we try to enforce collision-free endpoints, and expect that considering collisions brings logic to improve the overall consistency of the predictions. However, as will be highlighted in Sec. 4.3, this is not sufficient to generate joint predictions with this sampling that is otherwise mostly independent across agents even if it significantly improves the collision rate.

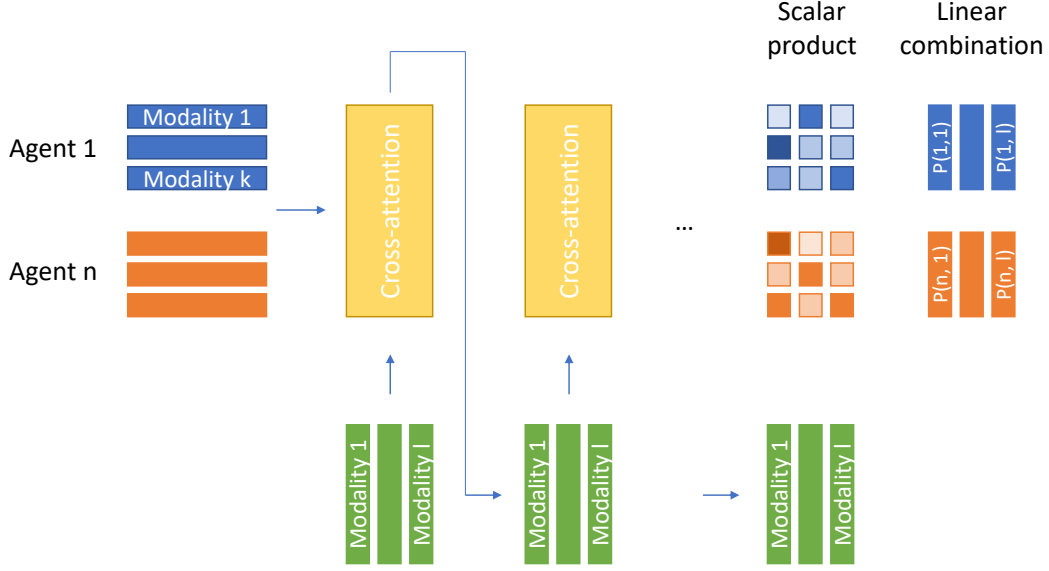


Figure 4: Illustration of THOMAS methods for generation of scene-consistent agent modalities

3.2.2 MODALITY COMBINATION RANKING

We address the scene consistency after the initial agent sampling. As single agents metrics are usually noticeably better than joint ones, this hint us that a potential solution already exists in the individually sampled modalities, and only a re-ordering of these is needed.

Our goal is to output a joint prediction $\mathcal{J} = (K, A)$ from the marginal prediction $\mathcal{M} = (A, K)$, where each scene modality k_s selects the combination of agent modalities k_a that are the most coherent together. Effectively, \mathcal{J} could just be a reordering of \mathcal{M} , as the matching modalities of each agent just need to be aligned together.

We illustrate our scene modality generation process in Fig. 4. We initialize L scene modality vectors S_l of D features each. $K \times A$ agent modality vectors A_k^a are also derived from each agent modality position. These vectors are obtained through a 2-layer MLP applied on the agent modality coordinates p_k^a , to which the stored agent encoding F_a (previously described in Sec. 3.1.1) is concatenated in order to help the model recognising modalities from the same agent. The scene modality vectors S_l are enriched through cross-attention layers on the agent modality encodings A_k^a . Then, for each scene modality l , for each agent a , a matching score is computed for each agent modality k as a scalar dot product between the scene modality vector S_l and the agent modality vector A_k^a :

$$s_l^{k_a} = S_l \cdot A_k^a$$

Since we can't use argmax as it is non-differentiable, we then use a soft argmax as a weighted linear combination of the agent modalities p_k^a using a softmax on the $s_l^{k_a}$ scores:

$$p_l^a = \text{softmax}(s_l^{k_a}) p_k^a$$

4 EXPERIMENTS

4.1 EXPERIMENTAL SETTINGS

4.1.1 DATASET

We use the Interaction v1.2 dataset that has recently opened a new multi-agent track in the context of its Interpret challenge. It contains 40 000 training cases, 11 794 validation cases and 20 000 testing cases, with each case containing between 1 and 40 agents to predict simultaneously.

4.1.2 METRICS

We report the marginal usually defined minFDE and MissRate that would be averaged over the agents after the minimum:

$$\min FDE_k = \frac{1}{A} \sum_a \min_k \|p_k^a - \hat{p}^a\|_2$$

$$MR_k = \frac{1}{A} \sum_a \min_k \mathbb{1}_{miss k}^a$$

As in Ettinger et al. (2021) and Zhan et al. (2019), a miss is counted when the prediction is closer than a lateral (1m) and an longitudinal threshold with regard to speed:

$$\text{Threshold}_{lon} = \begin{cases} 1 & v < 1.4m/s \\ 1 + \frac{v-1.4}{v-11} & 1.4m/s \leq v \leq 11m/s \\ 2 & v \geq 11m/s \end{cases}$$

For consistent scene multi-agent prediction, we report the joint metrics, where the average operation over the agents is done before the minimum operator:

$$\min JointFDE_k = \min_k \frac{1}{A} \sum_a \|p_k^a - \hat{p}^a\|_2$$

$$JointMR_k = \min_k \frac{1}{A} \sum_a \mathbb{1}_{miss k}^a$$

We also report *CrossCollisionrate* which is the percentage of modalities where two or more agents collide together, and *ConsistentMinJointMR*, which is *JointMR* where colliding modalities are also counted as misses even if they are closer than the defined threshold.

4.2 COMPARISON WITH STATE-OF-THE-ART

We compare our THOMAS model performance with other joint predictions methods ILVM (Casas et al., 2020) and SceneTransformer (Ngiam et al., 2021). For fair comparison, we use a GOHOME encoder for each of the method, and adapt them according so that they predict only endpoints similar to our method:

- ILVM (Casas et al., 2020) uses variational inference to learn a latent representation of the scene conditioned of each agent with a Scene Interaction Module, and decodes it with a similar Scene Interaction Module. We use a GOHOME encoder for the prior, posterior and decoder Scene Interaction Modules. We weight the KL term with $\beta = 1$ which worked best according to our experiments.
- SceneTransformer (Ngiam et al., 2021) applies a transformer architecture on a $[F, A, T, D]$ tensor where F is the potential modality dimension, A the agent dimension and T the time dimension, with D the feature embedding. They apply factorized self-attention to the agent and time dimensions separately, so that agents can look at each-other inside a specific scene modality. The resulting output is optimized using a jointly formalized loss. For our implementation, we get rid of the T dimension as we focus on endpoint prediction and coherence between the A agents. The initial encoded $[A, D]$ tensor is obtained with a GOHOME encoder, multiplied across the F futures and concatenated with a modality-specific one-hot encoding as in Ngiam et al. (2021) to obtain the $[F, A, D]$ tensor. We then apply two layers of agent self-attention similar to the original paper, before decoding the endpoints through a MLP.

The results are reported in Tab. 1. While having comparable marginal distance performance, THOMAS significantly outperforms other methods on every joint metric.

We also report our numbers from the Interpret multi-agent track challenge online leaderboard in Tab. 2. We are currently first of the challenge, with no competition in sight.

Table 1: Comparison of consistent solutions on Interpret multi-agent validation track

	Marginal metrics			Joint metrics			
	mADE	mFDE	MR	mFDE	MR	Col	cMR
ILVM (Casas et al., 2020)	0.30	0.62	10.8	0.84	19.8	5.7	21.3
SceneTranformer (Ngiam et al., 2021)	0.29	0.59	10.5	0.84	15.7	3.4	17.3
THOMAS	0.31	0.60	8.2	0.76	11.8	2.4	12.7

Table 2: Results on Interpret multi-agent leaderboard (test set)

	jointMinADE	JointMinFDE	JointMR	CrossCollision	ConsistentJointMR
THOMAS	0.36	0.85	17.2	25.7	18.7

4.3 ABLATION STUDIES

4.3.1 BASELINES

We establish the following baselines to assess the effects our THOMAS recombination

- Scalar output: we train a model with the GOHOME graph encoder and a multimodal scalar regression head similar to Liang et al. (2020); Ngiam et al. (2021). We optimize it with both marginal and joint formulation.
- Heatmap output with deterministic sampling: we try various sampling methods applied on the heatmap:
 - Deterministic sampling: we apply the sampling algorithm as described in Gilles et al. (2021a). We also evaluate a joint variant as described in Sec .
 - Learned sampling: we train a model to regress itself the sampled modalities from the input heatmap. This model is also optimized using either of the marginal or joint loss formulation

Compared to these baselines, THOMAS can be seen as an hybrid sampling method that takes the result of deterministic sampling as input and recombine it into a more coherent solution.

We report the comparison between the aforementioned baselines and THOMAS in Tab. 3.

Table 3: Comparison of consistent solutions on Interpret multi-agent validation track

Output	Sampling	Objective	Marginal metrics			Joint metrics			
			mADE	mFDE	MR	mFDE	MR	Col	cMR
Scalar	-	Marg	0.28	0.59	10.4	1.04	23.7	6.4	24.9
Scalar	-	Joint	0.30	0.65	12.8	0.88	17.4	4.5	18.8
Heat	Learned	Marg	0.26	0.46	4.9	0.98	20.9	4.1	21.9
Heat	Learned	Joint	0.29	0.58	9.8	0.88	15.2	3.0	16.4
Heat	Algo	Marg	0.29	0.54	3.8	0.83	14.8	7.2	15.9
Heat	Algo	Joint	0.29	0.54	3.8	0.83	14.8	2.6	15.6
Heat	Combi	Joint	0.31	0.60	8.2	0.76	11.8	2.4	12.7

REFERENCES

- Alexandre Alahi, Kratarth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *CVPR*, 2016.
- Florent Althé and Arnaud de La Fortelle. An lstm network for highway trajectory prediction. In *ITSC*, 2017.
- Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *CVPR*, 2020.
- Sergio Casas, Cole Gulino, Simon Suo, Katie Luo, Renjie Liao, and Raquel Urtasun. Implicit latent variable model for scene-consistent motion forecasting. In *ECCV*, 2020.
- Yuning Chai, Benjamin Sapp, Mayank Bansal, and Dragomir Anguelov. Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction. In *CoRL*, 2020.
- Ming-Fang Chang, John Lambert, Patsorn Sangkloy, Jagjeet Singh, Slawomir Bak, Andrew Hartnett, De Wang, Peter Carr, Simon Lucey, Deva Ramanan, et al. Argoverse: 3d tracking and forecasting with rich maps. In *CVPR*, 2019.
- Henggang Cui, Vladan Radosavljevic, Fang-Chieh Chou, Tsung-Han Lin, Thi Nguyen, Tzu-Kuo Huang, Jeff Schneider, and Nemanja Djuric. Multimodal trajectory predictions for autonomous driving using deep convolutional networks. In *ICRA*, 2019.
- Nachiket Deo and Mohan M Trivedi. Convolutional social pooling for vehicle trajectory prediction. In *CVPR*, 2018.
- Nachiket Deo and Mohan M Trivedi. Trajectory forecasts in unknown environments conditioned on grid-based plans. *arXiv:2001.00735*, 2020.
- Nachiket Deo, Eric M Wolff, and Oscar Beijbom. Multimodal trajectory prediction conditioned on lane-graph traversals. *arXiv:2106.15004*, 2021.
- Scott Ettinger, Shuyang Cheng, Benjamin Caine, Chenxi Liu, Hang Zhao, Sabeek Pradhan, Yuning Chai, Ben Sapp, Charles Qi, Yin Zhou, et al. Large scale interactive motion forecasting for autonomous driving: The waymo open motion dataset. *arXiv:2104.10133*, 2021.
- Jiyang Gao, Chen Sun, Hang Zhao, Yi Shen, Dragomir Anguelov, Congcong Li, and Cordelia Schmid. Vectornet: Encoding hd maps and agent dynamics from vectorized representation. In *CVPR*, 2020.
- Thomas Gilles, Stefano Sabatini, Dzmitry Tsishkou, Bogdan Stanciulescu, and Fabien Moutarde. Gohome: Graph-oriented heatmap output for future motion estimation. *arXiv preprint arXiv:2108.09640*, 2021a.
- Thomas Gilles, Stefano Sabatini, Dzmitry Tsishkou, Bogdan Stanciulescu, and Fabien Moutarde. Home: Heatmap output for future motion estimation. In *ITSC*, 2021b.
- Junru Gu, Chen Sun, and Hang Zhao. Densetnt: End-to-end trajectory prediction from dense goal sets. In *ICCV*, 2021.
- Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social gan: Socially acceptable trajectories with generative adversarial networks. In *CVPR*, 2018.
- Joey Hong, Benjamin Sapp, and James Philbin. Rules of the road: Predicting driving behavior with a convolutional model of semantic interactions. In *CVPR*, 2019.
- Siddhesh Khandelwal, William Qi, Jagjeet Singh, Andrew Hartnett, and Deva Ramanan. What-if motion prediction for autonomous driving. *arXiv:2008.10587*, 2020.
- ByeoungDo Kim, Seong Hyeon Park, Seokhwan Lee, Elbek Khoshimjonov, Dongsuk Kum, Junsoo Kim, Jeong Soo Kim, and Jun Won Choi. Lapred: Lane-aware prediction of multi-modal future trajectories of dynamic agents. In *CVPR*, 2021.

- Namhoon Lee, Wongun Choi, Paul Vernaza, Christopher B Choy, Philip HS Torr, and Manmohan Chandraker. Desire: Distant future prediction in dynamic scenes with interacting agents. In *CVPR*, 2017.
- Ming Liang, Bin Yang, Rui Hu, Yun Chen, Renjie Liao, Song Feng, and Raquel Urtasun. Learning lane graph representations for motion forecasting. In *ECCV*, 2020.
- Chenxu Luo, Lin Sun, Dariush Dabiri, and Alan Yuille. Probabilistic multi-modal trajectory prediction with lane attention for autonomous vehicles. *arXiv:2007.02574*, 2020.
- Jean Mercat, Thomas Gilles, Nicole El Zoghby, Guillaume Sandou, Dominique Beauvois, and Guillermo Pita Gil. Multi-head attention for multi-modal joint vehicle motion forecasting. In *ICRA*, 2020.
- Kaouther Messaoud, Nachiket Deo, Mohan M Trivedi, and Fawzi Nashashibi. Multi-head attention with joint agent-map representation for trajectory prediction in autonomous driving. *arXiv:2005.02545*, 2020.
- Jiquan Ngiam, Benjamin Caine, Vijay Vasudevan, Zhengdong Zhang, Hao-Tien Lewis Chiang, Jeffrey Ling, Rebecca Roelofs, Alex Bewley, Chenxi Liu, Ashish Venugopal, et al. Scene transformer: A unified multi-task model for behavior prediction and planning. *arXiv preprint arXiv:2106.08417*, 2021.
- Tung Phan-Minh, Elena Corina Grigore, Freddy A Boulton, Oscar Beijbom, and Eric M Wolff. Covernet: Multimodal behavior prediction using trajectory sets. In *CVPR*, 2020.
- Nicholas Rhinehart, Kris M Kitani, and Paul Vernaza. R2p2: A reparameterized pushforward policy for diverse, precise generative path forecasting. In *ECCV*, 2018.
- Amir Sadeghian, Vineet Kosaraju, Ali Sadeghian, Noriaki Hirose, Hamid Rezatofighi, and Silvio Savarese. Sophie: An attentive gan for predicting paths compliant to social and physical constraints. In *CVPR*, 2019.
- Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and Marco Pavone. Trajectron++: Dynamically-feasible trajectory forecasting with heterogeneous data. In *ECCV*, 2020.
- Haoran Song, Di Luan, Wenchao Ding, Michael Yu Wang, and Qifeng Chen. Learning to predict vehicle trajectories with model-based planning. *arXiv:2103.04027*, 2021.
- Yichuan Charlie Tang and Ruslan Salakhutdinov. Multiple futures prediction. In *NeurIPS*, 2019.
- David Wu and Yunnan Wu. Air² for interaction prediction. Waymo Open Dataset Challenges Reports, CVPR Workshop on Autonomous Driving, <http://cvpr2021.wad.vision/>, 2021.
- Maosheng Ye, Tongyi Cao, and Qifeng Chen. Tpcn: Temporal point cloud networks for motion forecasting. *arXiv:2103.03067*, 2021.
- Wenyuan Zeng, Ming Liang, Renjie Liao, and Raquel Urtasun. Lanercnn: Distributed representations for graph-centric motion forecasting. *arXiv:2101.06653*, 2021.
- Wei Zhan, Liting Sun, Di Wang, Haojie Shi, Aubrey Clausse, Maximilian Naumann, Julius Kummerle, Hendrik Konigshof, Christoph Stiller, Arnaud de La Fortelle, et al. Interaction dataset: An international, adversarial and cooperative motion dataset in interactive driving scenarios with semantic maps. *arXiv:1910.03088*, 2019.
- Lingyao Zhang, Po-Hsun Su, Jerrick Hoang, Galen Clark Haynes, and Micol Marchetti-Bowick. Map-adaptive goal-based trajectory prediction. In *CoRL*, 2020.
- Hang Zhao, Jiyang Gao, Tian Lan, Chen Sun, Benjamin Sapp, Balakrishnan Varadarajan, Yue Shen, Yi Shen, Yuning Chai, Cordelia Schmid, et al. Tnt: Target-driven trajectory prediction. *CoRL*, 2020.

A APPENDIX

You may include other additional sections here.