

SSL-Lanes: Self-Supervised Learning for Motion Forecasting in Autonomous Driving

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Self-supervised learning (SSL) is an emerging technique that has been
2 successfully employed to train convolutional neural networks (CNNs) and graph
3 neural networks (GNNs) for more transferable, generalizable, and robust repre-
4 sentation learning. However its potential in motion forecasting for autonomous
5 driving has rarely been explored. In this study, we report the first systematic explo-
6 ration and assessment of incorporating self-supervision into motion forecasting.
7 We first propose to investigate four novel self-supervised learning tasks for motion
8 forecasting with theoretical rationale and quantitative and qualitative comparisons
9 on the challenging large-scale Argoverse dataset. Secondly, we point out that
10 our auxiliary SSL-based learning setup not only outperforms forecasting methods
11 which use transformers, complicated fusion mechanisms and sophisticated online
12 dense goal candidate optimization algorithms in terms of performance accuracy,
13 but also has low inference time and architectural complexity. Lastly, we conduct
14 several experiments to understand why SSL improves motion forecasting.

15 **Keywords:** Motion Forecasting, Autonomous Driving, Self-Supervised Learning

16 1 Introduction

17 Motion forecasting in a real-world urban environment is an important task for autonomous robots. It
18 involves predicting the future trajectories of traffic agents including vehicles and pedestrians. This is
19 absolutely crucial in the self-driving domain for safe, comfortable and efficient operation. However,
20 this is a very challenging problem. Difficulties include inherent stochasticity and multimodality
21 of driving behaviors, and that future motion can involve complicated maneuvers such as yielding,
22 nudging, lane-changing, turning and acceleration or deceleration.

23 The motion prediction task has traditionally been based on kinematic constraints and road map in-
24 formation with handcrafted rules. These approaches however fail to capture long-term behavior and
25 interactions with map structure and other traffic agents in complex scenarios. Tremendous progress
26 has been made with data-driven methods in motion forecasting [3, 4, 5, 6, 7, 8, 9, 10]. Recent
27 methods use a vector representation for HD maps and agent trajectories, including approaches like
28 Lane-GCN [2], Lane-RCNN [11], Vector-Net [12], TNT [5] and Dense-TNT [6]. More recently, the
29 enormous success of transformers [13] has been leveraged for forecasting in mm-Transformer [9],
30 Scene transformer [8], Multimodal transformer [14] and Latent Variable Sequential Transformers
31 [15]. Most of these methods however are extremely complex in terms of architecture and have low
32 inference speeds, which makes them unsuitable for real-world settings.

33 In this work, we extend ideas from self-supervised learning (SSL) to the motion forecasting task.
34 Self-supervision has seen huge interest in both natural language processing and computer vision
35 [16] to make use of freely available data without the need for annotations. It aims to assist the
36 model to learn more transferable and generalized representation from pseudo-labels via pretext tasks.
37 Given the recent success of self-supervision with CNNs, transformers, and GNNs, we are naturally
38 motivated to ask the question: *Can self-supervised learning improve accuracy and generalizability*
39 *of motion forecasting, without sacrificing inference speed or architectural simplicity?*

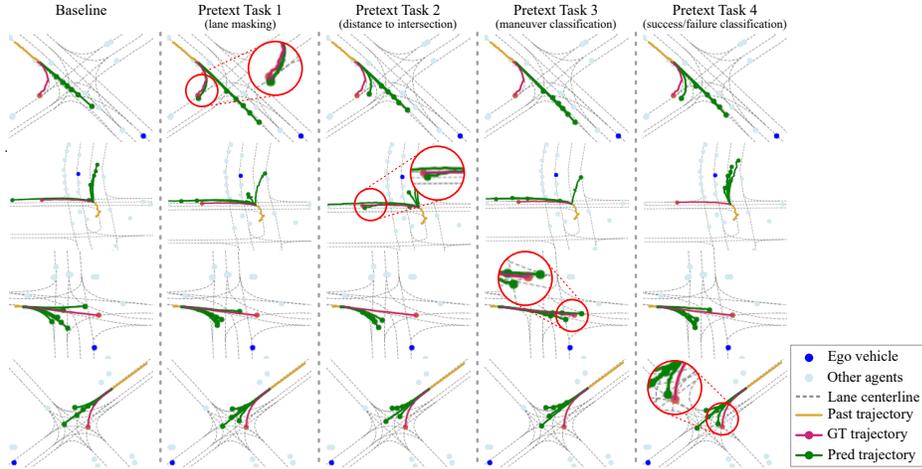


Figure 1: Motion forecasting on Argoverse [1] validation. We show four challenging scenarios at intersections. The baseline [2] misses all the predictions. In the first row, our proposed lane masking successfully captures the right-turn. For the second row, predicting distance to intersection helps the most in capturing the left turn. In the third row, acceleration at an intersection is best captured by the model that is made to classify maneuvers of traffic agents. Finally, in the fourth row, classifying successful final goal states is the most effective at capturing the left turn. These tasks are trained with pseudo-labels which are obtained for free from data. Please refer to Sec. 6.2 for details.

40 **Contributions:** Our work, SSL-Lanes, presents the first systematic study on how to incorporate
 41 self-supervision in a standard data-driven motion forecasting model. Our contributions are: (a)
 42 We demonstrate the effectiveness of incorporating self-supervised learning in motion forecasting.
 43 Since this does not add extra parameters or compute during inference, SSL-Lanes achieves the best
 44 accuracy-simplicity-efficiency trade-off on the challenging large-scale Argoverse [1] benchmark.
 45 (b) We propose four self-supervised tasks based on the nature of the motion forecasting problem.
 46 The key idea is to leverage easily accessible map/agent-level information to define domain-specific
 47 pretext tasks that encourage the standard model to capture more superior and generalizable represen-
 48 tations for forecasting in comparison to pure supervised learning. (c) We further design experiments
 49 to explore why forecasting benefits from SSL. We provide extensive results to hypothesize that
 50 SSL-Lanes learns richer features from the SSL training as compared to a model trained with vanilla
 51 supervised learning.

52 2 Related Work

53 **Motion Forecasting:** Traditional methods for motion forecasting primarily use Kalman filtering
 54 [17] with a prior from HD-maps to predict future motion states [18, 19]. With the huge success
 55 of deep learning, recent works use data-driven approaches for motion forecasting. These methods
 56 explore different architectures involving rasterized images and CNNs [3, 20, 21], vectorized repre-
 57 sentations and GNNs [12, 11, 22, 4, 7], point-cloud representations [23], transformers [8, 9, 15, 14]
 58 and sophisticated fusion mechanisms [2], to generate features that predict final output trajectories.
 59 While the focus of these works is to find more effective ways of feature extraction from HD-maps
 60 and interacting agents, they need huge model capacity, heavy parameterization, and extensive aug-
 61 mentations or large amounts of data to converge to a general solution. Other works [5, 10, 24, 25]
 62 build on them to incorporate prior knowledge in the form of predefined candidate trajectories ob-
 63 tained from sampling or clustering strategies from training data. However the disadvantage of these
 64 methods is that their performance is highly related to the quality of the trajectory proposals, which
 65 becomes an extra dependency. End-to-end solutions for optimizing end-points of these candidates
 66 trajectories are proposed by Dense-TNT [6] and HOME [26]. Dense-TNT has state-of-the-art accu-
 67 racy with a reasonable parameter budget, but its online dense goal candidate optimization strategy is
 68 computationally very expensive, which is unrealistic for real-time operations like autonomous driv-

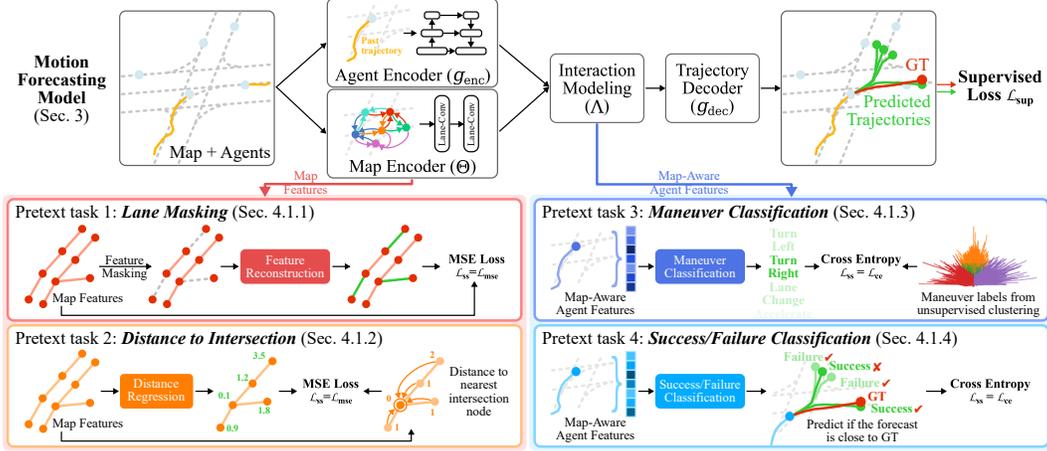


Figure 2: Illustration of the overall SSL-Lanes framework for self-supervision on motion forecasting through joint training. SSL-Lanes improves upon a standard-motion forecasting baseline, that consists of an agent encoder, map encoder, interaction model and a trajectory decoder, trained using a supervised loss \mathcal{L}_{sup} . SSL-Lanes proposes four pretext tasks: (1) Lane Masking: which recovers feature information from the perturbed lane graphs. (2) Distance to Intersection: which predicts the distance (in terms of shortest path length) from all lane nodes to intersection nodes. (3) Maneuver Classification: predicts the form of a ‘maneuver’ the agent-of-interest intends to execute (4) Success/Failure Classification: which trains an agent specialized at achieving end-point goals.

69 ing. Lately, ensembling techniques like MultiPath++ [27] and DCMS [28] have been proposed and
70 while they have high forecasting performance, a major disadvantage is their high memory cost for
71 training and heavy computational cost at inference. We also refer the reader to the supplementary
72 for a detailed discussion of how SSL-Lanes differs from methods like Vector-Net [12], CS-LSTM
73 [29] and MultiPath[3].

74 **Self-supervised Learning:** SSL is a rapidly emerging learning framework that generates additional
75 supervised signals to train deep learning models through carefully designed pretext tasks. In the
76 image domain, various self-supervised learning techniques have been developed for learning high-
77 level image representations, including predicting the relative locations of image patches [30], jigsaw
78 puzzle [31], image rotation [32], image clustering [33], image inpainting [34], image colorization
79 [35] and segmentation prediction [36]. In the domain of graphs and graph neural networks, pretext
80 tasks include graph partitioning, node clustering, context prediction and graph completion [37, 38,
81 39, 40]. To the best of our knowledge, this is the first principled approach that explores motion
82 forecasting for autonomous driving with self-supervision.

83 3 Background

84 **Problem Formulation:** We are given the past motion of N actors. The i -th actor is denoted as a
85 set of center locations over the past L time-steps. We pre-process it to represent each trajectory as
86 a sequence of displacements $\mathcal{P}_i = \{\Delta p_i^{-L+1}, \dots, \Delta p_i^{-1}, \Delta p_i^0\}$, where p_i^l is the 2D displacement
87 from time step $l - 1$ to l . We are also given a high-definition (HD) map, which contains lanes and
88 semantic attributes. Each lane is composed of many consecutive lane nodes, with a total of M nodes.
89 $\mathbf{X} \in \mathbb{R}^{M \times F}$ denotes the lane node feature matrix, where $x_j = \mathbf{X}[j, :]^T$ is the F -dimensional lane
90 node vector. Following the connections between lane centerlines (i.e., predecessor, successor, left
91 neighbour and right neighbour), we represent the connectivity within the lane nodes with 4 adjacency
92 matrices $\{\mathbf{A}_f\}_{f \in \{\text{pre}, \text{suc}, \text{left}, \text{right}\}}$, with $\mathbf{A}_f \in \mathbb{R}^{M \times M}$. This implies that if $\mathbf{A}_{f,gh} = 1$, then node h is
93 an f -type neighbor of node g . Our goal is to forecast the future motions of all actors in the scene
94 $\mathcal{O}_{GT}^{1:T} = \{(x_i^1, y_i^1), \dots, (x_i^T, y_i^T) | i = 1, \dots, N\}$, where T is our prediction horizon.

SSL Task	Property Level	Primary Assumption	Type
Lane-Masking	Map features	Local map structure	Aux. auto-encoder
Distance to Intersection		Global map structure	Aux. regression
Maneuver Classification	Map-aware agent features	Agent feature similarity	Aux. classification
Success/Failure Classification		Distance to success state	

Table 1: Overview of our proposed self-supervised (SSL) tasks

95 **Standard Motion Forecasting Model:** We briefly introduce a standard data-driven motion fore-
96 casting framework, consisting of a feature encoder, interaction-modeler and prediction header.

97 *Feature Encoding:* We first encode the agent and map inputs similar to Lane-GCN [2]. The agent
98 encoder includes a 1D convolution with a feature pyramid network, parameterized by g_{enc} , as given
99 by Eq. (1). For map-encoding, we adopt two Lane-Conv residual blocks, parameterized by $\Theta =$
100 $\{\mathbf{W}_0, \mathbf{W}_{\text{left}}, \mathbf{W}_{\text{right}}, \mathbf{W}_{\text{pre},k}, \mathbf{W}_{\text{suc},k}\}$, where $k \in \{1, 2, 4, 8, 16, 32\}$, as given by Eq. (2).

$$101 \quad \hat{\mathbf{p}}_i = g_{\text{enc}}(\mathcal{P}_i) \quad (1)$$

$$102 \quad \mathbf{Y} = \mathbf{X}\mathbf{W}_0 + \sum_{j \in \{\text{left}, \text{right}\}} \mathbf{A}_j \mathbf{X}\mathbf{W}_j + \sum_k \mathbf{A}_{\text{pre}}^k \mathbf{X}\mathbf{W}_{\text{pre},k} + \mathbf{A}_{\text{suc}}^k \mathbf{X}\mathbf{W}_{\text{suc},k} \quad (2)$$

102 *Modeling Interactions:* Since the behavior of agents depends on map topology and social consis-
103 tency, each encoded agent i subsequently aggregates context from the surrounding map features and
104 its neighboring agent features, via spatial attention [41] as given by Eq. (3):

$$105 \quad \tilde{\mathbf{p}}_i = \hat{\mathbf{p}}_i \mathbf{W}_{\text{M2A}} + \sum_j \phi(\text{concat}(\hat{\mathbf{p}}_i, \Delta_{i,j}, \mathbf{y}_j) \mathbf{W}_1) \mathbf{W}_2 \quad (3)$$

$$106 \quad \hat{\mathbf{p}}_i = \tilde{\mathbf{p}}_i \mathbf{W}_{\text{A2A}} + \sum_j \phi(\text{concat}(\tilde{\mathbf{p}}_i, \Delta_{i,j}, \hat{\mathbf{p}}_j) \mathbf{W}_3) \mathbf{W}_4$$

105 Here, \mathbf{y}_j is the feature of the j -th node, $\hat{\mathbf{p}}_i$ is the feature of the i -th agent, ϕ the composition of layer
106 normalization and ReLU, and $\Delta_{i,j} = \text{MLP}(\mathbf{v}_j - \mathbf{v}_i)$, where \mathbf{v} denotes the (x, y) 2-D **bird's-eye-view**
107 (BEV) location of the agent or the lane node. The parameters for map and agent feature aggregation
108 is represented by $\Lambda = \{\mathbf{W}_{\text{M2A}}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_{\text{A2A}}, \mathbf{W}_3, \mathbf{W}_4\}$.

109 *Trajectory Prediction:* Finally, we decode the future trajectories from the features $\hat{\mathbf{p}}_i$ corresponding
110 to the agents of interest as given by: $\mathcal{O}_{\text{pred}}^{1:T} = \{g_{\text{dec}}(\hat{\mathbf{p}}_i) | i = 1, \dots, N\}$, where g_{dec} is the parameter-
111 ized trajectory decoder. The parameters for the motion forecasting model are learned by minimizing
112 the supervised loss (\mathcal{L}_{sup}) calculated between the predicted output and the ground-truth future tra-
113 jectories ($\mathcal{O}_{\text{GT}}^{1:T}$), as given by Eq. (4):

$$114 \quad g_{\text{enc}}^*, \Theta^*, \Lambda^*, g_{\text{dec}}^* = \arg \min_{g_{\text{enc}}, \Theta, \Lambda, g_{\text{dec}}} \mathcal{L}_{\text{sup}}(\mathcal{O}_{\text{pred}}^{1:T}, \mathcal{O}_{\text{GT}}^{1:T}) \quad (4)$$

114 4 SSL-Lanes

115 The goal of our proposed SSL-Lanes framework is to improve the performance of the primary
116 motion forecasting baseline by learning simultaneously with various self-supervised tasks. Fig. 2
117 shows the pipeline of our proposed approach, and Tab. 1 summarizes the self-supervised tasks.

118 **Self-Supervision meets Motion Forecasting:** Considering our motion forecasting task and a self-
119 supervised task, the output and the training process can be formulated as:

$$120 \quad \Psi^*, \Omega^*, \Theta_{\text{ss}}^* = \arg \min_{\Psi, \Omega, \Theta_{\text{ss}}} \alpha_1 \mathcal{L}_{\text{sup}}(\Psi, \Omega) + \alpha_2 \mathcal{L}_{\text{ss}}(\Psi, \Theta_{\text{ss}}) \quad (5)$$

120 where, $\mathcal{L}_{\text{ss}}(\cdot, \cdot)$ is the loss function of the self-supervised task, Θ_{ss} parameterizes the corresponding
121 task-specific layers, and $\alpha_1, \alpha_2 \in \mathbb{R}_{>0}$ are the weights for the supervised and self-supervised losses.
122 If the pretext task only focuses on the map encoder, then $\Psi = \{\Theta\}$ and $\Omega = \{g_{\text{enc}}, \Lambda, g_{\text{dec}}\}$. Other-
123 wise, $\Psi = \{g_{\text{enc}}, \Theta, \Lambda\}$ and $\Omega = \{g_{\text{dec}}\}$. Henceforth, we also define the following representations.
124 We will represent the primary task encoder as function f_{Ψ} , parameterized by Ψ . Furthermore, given
125 a pretext task, which we will design in the next section, the pretext decoder $p_{\Theta_{\text{ss}}}$ is a function that
126 predicts pseudo-labels and is parameterized by Θ_{ss} .

Method	minADE ₁	minFDE ₁	MR ₁	minADE ₆	minFDE ₆	MR ₆
Baseline	1.42	3.18	51.35	0.73	1.12	11.07
Lane-Masking	1.36	2.96	49.45	0.70	1.02	8.82
Distance to Intersection	1.38	3.02	49.53	0.71	1.04	8.93
Maneuver Classification	1.33	2.90	49.26	0.72	1.05	9.36
Success/Failure Classification	1.35	2.93	48.54	0.70	1.01	8.59

Table 2: Motion forecasting performance on Argoverse validation with our proposed pretext tasks

127 4.1 Pretext tasks for Motion Forecasting

128 At the core of our SSL-Lanes approach is defining pretext tasks based upon self-supervised informa-
129 tion from the underlying map structure *and* the overall temporal prediction problem itself (Tab. 1).

130 4.1.1 Lane-Masking

131 The goal of the *Lane-Masking* pretext task is to encourage the map encoder $\Psi = \{\Theta\}$ to learn local
132 structure information in addition to the forecasting task that is being optimized. Specifically, we
133 randomly mask (i.e., set equal to zero) the features of m_a percent of nodes per lane and then ask the
134 self-supervised decoder to reconstruct these features.

$$\Psi^*, \Theta_{ss}^* = \arg \min_{\Psi, \Theta_{ss}} \frac{1}{m_a} \sum_{i=1}^{m_a} \mathcal{L}_{\text{mse}} \left(p_{\Theta_{ss}}([f_{\Psi}(\tilde{\mathbf{X}}, \mathbf{A}_f)]_{v_i}), \mathbf{X}_i \right) \quad (6)$$

135 Here, $\tilde{\mathbf{X}}$ is the node feature matrix corrupted with random masking, i.e., some rows of \mathbf{X} corre-
136 sponding to nodes v_i are set to zero. $p_{\Theta_{ss}}$ is a fully connected network that maps the representations
137 to the reconstructed features. \mathcal{L}_{mse} is the mean squared error (MSE) loss function penalizing the
138 distance between the reconstructed map features $p_{\Theta_{ss}}([f_{\Psi}(\tilde{\mathbf{X}}, \mathbf{A}_f)]_{v_i})$ for node v_i and its actual
139 features \mathbf{X}_i .

140 4.1.2 Distance to Intersection

141 *Distance-to-Intersection* pretext task is proposed to guide the map-encoder, $\Psi = \{\Theta\}$, to maintain
142 global topology information by predicting the distance (in terms of shortest path length) from all
143 lane nodes to intersection nodes. We aim to regress the distances from each lane node to pre-
144 labeled intersection nodes annotated as part of the dataset. Given K labeled intersection nodes
145 $\mathcal{V}_{\text{intersection}} = \{v_{\text{intersection},k} | k = 1, \dots, K\}$, we first generate reliable pseudo labels using breadth-first
146 search (BFS). Specifically, BFS calculates the shortest distance $d_i \in \mathbb{R}$ for every lane node v_i from
147 the given set $\mathcal{V}_{\text{intersection}}$. The target of this task is to predict the pseudo-labeled distances using a
148 pretext decoder. If $p_{\Theta_{ss}}([f_{\Psi}(\mathbf{X}, \mathbf{A}_f)]_{v_i})$ is the prediction of node v_i , and \mathcal{L}_{mse} is the mean-squared
149 error loss function for regression, then the loss formulation for this SSL pretext task is as follows:

$$\Psi^*, \Theta_{ss}^* = \arg \min_{\Psi, \Theta_{ss}} \frac{1}{M} \sum_{i=1}^M \mathcal{L}_{\text{mse}} \left(p_{\Theta_{ss}}([f_{\Psi}(\mathbf{X}, \mathbf{A}_f)]_{v_i}), d_i \right) \quad (7)$$

150 4.1.3 Maneuver Classification

151 We propose *Maneuver Classification*, and we expect it to provide prior regularization to $\Psi =$
152 $\{g_{\text{enc}}, \Theta, \Lambda\}$, based on driving modes of agents. We aim to construct pseudo label to di-
153 vide agents into different clusters according to their driving behavior and thus explore un-
154 supervised clustering algorithms to acquire the maneuver for each agent. We find that us-
155 ing naive k -Means (on agent end-points) or DBSCAN (on Hausdorff distance between entire
156 trajectories [42]) leads to noisy clustering. We find that constrained k -means [43] on agent
157 end-points works best to divide trajectory samples into C clusters equally. We define $C =$
158 $\{\text{maintain-speed, accelerate, decelerate, turn-left, turn-right, lane-change}\}$ and the clustering func-
159 tion as ρ . If $p_{\Theta_{ss}}(f_{\Psi}(\mathcal{P}_i, \mathbf{X}, \mathbf{A}_f))$ is the prediction of agent i 's intention and $E_i = (x_{i,\text{GT}}^T, y_{i,\text{GT}}^T)$

Method	minADE ₁	minFDE ₁	MR ₁	minADE ₆	minFDE ₆	MR ₆	b-FDE ₆
NN + Map [1]	3.65	8.12	94.0	2.08	4.02	58.0	-
Jean [4]	1.74	4.24	68.56	0.98	1.42	13.08	2.12
Lane-GCN [2]	1.71	3.78	58.77	0.87	1.36	16.20	2.05
LaneRCNN [11]	<u>1.68</u>	3.69	56.85	0.90	1.45	<u>12.32</u>	2.15
TNT [5]	<u>1.77</u>	3.91	59.70	0.94	1.54	13.30	2.14
DenseTNT [6]	1.68	3.63	58.43	0.88	1.28	12.58	1.97
PRIME [24]	1.91	<u>3.82</u>	58.67	1.22	1.55	11.50	2.09
WIMP [7]	1.82	4.03	62.88	0.90	1.42	16.69	2.11
TPCN [23]	1.66	3.69	58.80	0.87	1.38	15.80	1.92
HOME [26]	1.70	3.68	57.23	0.89	1.29	8.46	1.86
mmTransformer [9]	1.77	4.00	61.78	0.87	1.34	<u>15.40</u>	<u>2.03</u>
MultiModalTransformer [14]	1.74	3.90	60.23	<u>0.84</u>	1.29	14.29	1.94
LatentVariableTransformer [15]	-	-	-	0.89	1.41	16.00	-
SceneTransformer [8]	1.81	4.06	59.21	0.80	1.23	12.55	<u>1.88</u>
Success/Failure Classification (Ours)	1.63	3.56	56.71	<u>0.84</u>	<u>1.25</u>	13.26	1.94

Table 3: Comparison of our (best) proposed model and top approaches on the Argoverse Test. The best results are in bold and underlined, and the second best is also underlined.

160 is its ground-truth end-point, then the learning objective is to classify each agent maneuver into its
161 corresponding cluster using cross-entropy loss \mathcal{L}_{ce} as:

$$\Psi^*, \Theta_{ss}^* = \arg \min_{\Psi, \Theta_{ss}} \mathcal{L}_{ce} \left(p_{\Theta_{ss}}(f_{\Psi}(\mathcal{P}_i, \mathbf{X}, \mathbf{A}_f)), \rho(E_i) \right) \quad (8)$$

162 4.1.4 Forecasting Success/Failure Classification

163 We propose a pretext task called *Success/Failure Classification*, which trains an agent specialized at
164 achieving end-point goals and thus links directly to the forecasting task. We expect this to constrain
165 $\Psi = \{g_{enc}, \Theta, \Lambda\}$ to predict trajectories ϵ distance away from the correct final end-point. Similar to
166 maneuver classification, we wish to create pseudo-labels for our data samples. We label trajectory
167 predictions as successful ($c = 1$) if the final prediction $(x_{i,pred}^T, y_{i,pred}^T)$ is within $\epsilon < 2m$ of the
168 final end-point E_i , and as failure ($c = 0$) otherwise. We choose $2m$ as our ϵ threshold because it is
169 also used for miss-rate calculation (Sec. 5). If the pretext decoder predicts agent i 's final-endpoint
170 as $p_{\Theta_{ss}}(f_{\Psi}(\mathcal{P}_i, \mathbf{X}, \mathbf{A}_f))$ and, given the ground-truth end-point E_i , its success or failure label is c_i ,
171 then the pretext loss can be formulated as:

$$\Psi^*, \Theta_{ss}^* = \arg \min_{\Psi, \Theta_{ss}} \mathcal{L}_{ce} \left(p_{\Theta_{ss}}(f_{\Psi}(\mathcal{P}_i, \mathbf{X}, \mathbf{A}_f)), c_i \right) \quad (9)$$

172 4.2 Learning

173 As all the modules are differentiable, we can train the model in an end-to-end way. We use the sum
174 of classification, regression and self-supervised losses to train the model. Specifically, we use:

$$\mathcal{L} = \mathcal{L}_{cls} + \mathcal{L}_{reg} + \mathcal{L}_{terminal} + \mathcal{L}_{ss} \quad (10)$$

175 For classification and regression loss design, we adopt the formulation proposed in [2]. $\mathcal{L}_{terminal} =$
176 $\frac{1}{N} \sum_{i=1}^N L2 \left((x_{i,pred}^T, y_{i,pred}^T), (x_{i,GT}^T, y_{i,GT}^T) \right)$ is a simple L2 loss that minimizes the distance between
177 predicted final-endpoints and the ground-truth. This is because \mathcal{L}_{reg} is averaged across all time-
178 points $1 : T$, and from a practical end user perspective, minimizing the endpoint loss is much more
179 important than weighting loss from all time-steps equally. Our proposed pretext tasks contribute to
180 \mathcal{L}_{ss} . During evaluation, we study each pretext task separately, and their corresponding loss formula-
181 tions defined in Eq. (6), Eq. (7), Eq. (8), Eq. (9) are used for joint training.

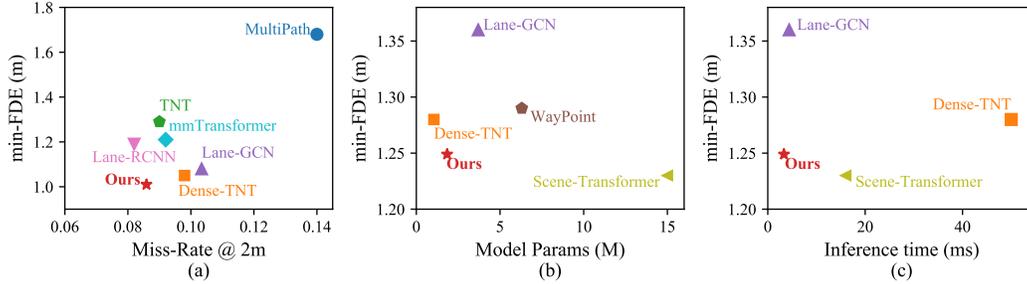


Figure 3: (a) min-FDE₆ - Miss-Rate₆ trade-off on Argoverse Validation. Lower-left is better. We optimize both successfully in comparison to other popular approaches. (b) and (c) We plot min-FDE on Argoverse Test against number of model parameters (in millions) and inference time (in milliseconds). We find that there is a trade-off between min-FDE performance, architectural complexity (as measured by number of parameters) and computational efficiency (as measured by inference time). Our work achieves the best trade-off (lower-left).

182 5 Experiments

183 **Dataset:** Argoverse provides a large-scale dataset, where the task is to forecast 3 seconds of future
 184 motions, given 2 seconds of past observations. It has more than 300K real-world driving sequences
 185 collected in Miami (MIA) and Pittsburgh (PIT). Those sequences are further split into train, val-
 186 idation, and test sets, without any geographical overlap. Each of them has 205,942, 39,472, and
 187 78,143 sequences respectively. In particular, each sequence contains the positions of all actors in
 188 a scene within the past 2 seconds history, annotated at 10Hz. It also specifies one actor of interest
 189 in the scene, with type ‘agent’, whose future 3 seconds of motion are used for the evaluation. The
 190 train and validation splits additionally provide future locations of all actors within 3 second hori-
 191 zon labeled at 10Hz, while annotations for test sequences are withheld from the public and used
 192 for the leaderboard evaluation. HD map information is available for all sequences. We have two
 193 main requirements for the dataset: (a) **Scale of Data:** Modern motion forecasting methods and self-
 194 supervised learning systems require a large amount of training data to imitate human maneuvers in
 195 complex real-world scenarios. Thus, the dataset should be large-scale and diverse, such that it has
 196 a wide range of behaviors and trajectory shapes across different geometries represented in the data.
 197 (b) **Interesting Scenarios for Forecasting Evaluation:** The dataset should be collected for inter-
 198 esting behaviours by biasing sampling towards complex observed behaviours (e.g., lane changes,
 199 turns) and road features (e.g., intersections), since we wish to focus on these cases. We find that
 200 on the basis of these requirements, as well as its popularity in the the motion forecasting commu-
 201 nity, Argoverse [1] is the best candidate to showcase our method. Please refer to the supplementary
 202 for more details regarding why we choose to focus on it in comparison to other motion forecasting
 203 benchmarks.

204 **Metrics:** ADE is defined as the average displacement error between ground-truth trajectories and
 205 predicted trajectories over all time steps. FDE is defined as displacement error between ground-truth
 206 trajectories and predicted trajectories at the final time step. We compute K likely trajectories for
 207 each scenario with the ground truth label, where $K = 1$ and $K = 6$ are used. Therefore, minADE
 208 and minFDE are minimum ADE and FDE over the top K predictions, respectively. Miss rate (MR)
 209 is defined as the percentage of the best-predicted trajectories whose FDE is within a threshold (2 m).
 210 Brier-minFDE is the minFDE plus $(1 - p)^2$, where p is the corresponding trajectory probability.

211 **Experimental Details:** To normalize the data, we translate and rotate the coordinate system of each
 212 sequence so that the origin is at current position $t = 0$ of ‘agent’ actor and x-axis is aligned with its
 213 current direction, i.e., orientation from the agent location at $t = -1$ to the agent location at $t = 0$
 214 is the positive x axis. We use all actors and lanes whose distance from the agent is smaller than
 215 100 meters as the input. We train the model on 4 TITAN-X GPUs using a batch size of 128 with
 216 the Adam [44] optimizer with an initial learning rate of 1×10^{-3} , which is decayed to 1×10^{-4} at

217 100,000 steps. The training process finishes at 128,000 steps and takes about 10 hours to complete.
218 We provide more implementation details in the supplementary.

219 6 Results

220 6.1 Ablation Studies

221 We first examine the effect of incorporating our proposed pretext tasks (Sec. 4) with the standard
222 data-driven motion forecasting baseline (Sec. 3). While evaluating the importance of our proposed
223 pretext tasks, we wish to underline that motion prediction for autonomous driving is a safety-critical
224 task, especially at intersections where most of our data is collected, and most accidents also happen.
225 We thus posit that in this situation, even a small error in predicting final locations (FDE) for a given
226 agent can lead to dangerous potential collision scenarios. Results in Tab. 2 show that all proposed
227 pretext tasks improve motion forecasting performance for Argoverse. Specifically, the Lane Mask
228 pretext task improves min-FDE by 8.9% and MR@2m by 20.3%. Distance to Intersection improves
229 min-FDE by 7.1% and 19.3%. Maneuver classification improves min-FDE by 6.3% and MR@2m
230 by 15.4%. We expect that improving the quality of clustering for maneuvers and thus creating better
231 pseudo-labels will improve this further. Finally, Success/Failure classification improves min-FDE
232 by 9.8% and, perhaps expectedly, MR@2m by 22.4%. Moreover, since pretext tasks are not used for
233 inference and only for training, they also do not add any extra parameters or FLOPs to the baseline,
234 thereby increasing accuracy but at no cost to computational efficiency or architectural complexity.

235 6.2 Comparison with State-of-the-Art

236 **Performance:** We compare our approach with top entries on Argoverse [1] in Tab. 3. SSL-Lanes
237 improves the metrics for $K = 1$ convincingly and outperforms existing approaches w.r.t. $min-$
238 ADE_1 , $min-FDE_1$ and MR_1 . We are strongly competitive w.r.t. $min-ADE_6$, $min-FDE_6$ and MR_6 .
239 with a relatively simple architecture.

240 **Trade-off between min-FDE and Miss-Rate:** $min-FDE_6$ and MR_6 are both important for au-
241 tonomous robots to optimize. Ideally we wish for both of these metrics to be low. However, there
242 exists a frequent trade-off between them. We compare this trade-off in Fig. 3(a) w.r.t 6 other pop-
243 ular motion forecasting models (in terms of citations and GitHub stars), namely: Lane-GCN [2],
244 Lane-RCNN [2], MultiPath [3], mm-Transformer [9], TNT [5] and Dense-TNT [6] on the Argov-
245 erse Validation Set. We are on the lowest-left of Fig. 3(a), meaning we optimize both $min-FDE_6$ and
246 MR_6 successfully in comparison to other top models.

247 **Trade-off between accuracy, efficiency and complexity:** We are the first to point out a trade-off
248 that exists for current state-of-the-art motion forecasting models between forecasting performance,
249 architectural complexity and inference speed, in this work. This is illustrated in Fig. 3(b, c). In
250 contrast to the popular models, our approach has high accuracy ($min-FDE_6$: 1.25m, MR_6 : 13.3%),
251 while also having low architectural complexity (1.84M parameters) and high inference speed (3.3
252 ms). Thus it provides a great balance for application to real-time safety-critical autonomous robots.

253 **Qualitative Results:** We present some multi-modal prediction trajectories on several hard cases
254 shown in Fig. 1. The yellow trajectory represents the observed 2s. Red represents ground truth for
255 the next 3s and green represents the multiple forecasted trajectories for those 3s. In Row 1, the agent
256 turns right at the intersection. The baseline misses this mode completely, despite having access to
257 the map. The model trained with lane-masking successfully predicts this right turn within 2m of the
258 ground-truth end-point. In Row 2, the agent has a noisy past history and accelerates while turning
259 left at the intersection. The pretext task distance-to-intersection can correctly capture this, while the
260 baseline has only one trajectory covering this mode but vastly overshoots the ground-truth. Inter-
261 estingly, we note that the success/failure pretext task is unable to capture this mode. We believe
262 this is due to a stronger prior imposed by the model during learning. In Row 3, we have an agent
263 accelerating while going straight at an intersection. We find that the maneuver classification pretext
264 task is the only model that correctly predicts trajectories aligned with the ground-truth. In Row 4,

Description	Experimental Setup		Method	minADE ₆	minFDE ₆	MR ₆
	Training	Validation				
Effects of limited training data	25% of train	All	Baseline	0.82	1.33	14.66
			Ours	0.78	1.22	12.63
Effects of new domain	100% PIT + 20% MIA	MIA val	Baseline	0.88	1.46	17.21
			Ours	0.85	1.34	14.96
Performance on difficult maneuvers	All	Turning & lane changing	Baseline	0.90	1.53	19.90
			Ours	0.84	1.34	14.93
Effects of imbalanced data	2x straight 1x other maneuvers	Turning & lane changing	Baseline	0.94	1.65	21.53
			Ours	0.90	1.49	17.97
Effects of noisy data	All	Gaussian noise ($\sigma = 0.2$) with $p = 0.25$	Baseline	1.01	1.37	15.59
			Ours	0.96	1.24	11.98
Effects of noisy data	All	Gaussian noise ($\sigma = 0.2$) with $p = 0.5$	Baseline	1.19	1.56	20.64
			Ours	1.13	1.40	15.65

Table 4: Different experimental settings for SSL-based training

265 we have an agent turning left at an intersection. Most of the predictions of other models predicts
266 that the agent will go straight. The success/failure pretext task however picks up on the left-turn,
267 possibly due to the priors imposed upon it by end-point conditioning.
268 Overall, SSL-Lanes can capture left and right turns better, while also being able to discern accelera-
269 tion at intersections. Our pretext tasks provide priors for the model and provides data-regularization
270 for free. We believe this can improve forecasting through better understanding of map topology,
271 agent context with respect to the map, and generalization with respect to imbalance implicitly present
272 in data.

273 6.3 When does SSL help Motion Forecasting?

274 We design 6 different training and testing setups as shown in Tab. 4. We use Success/Failure classifi-
275 cation as the pretext task, and all models are trained for 50,000 steps. We initialize the map-encoder
276 with the parameters from a model trained with the lane-masking pretext task.

277 We hypothesize that training with SSL pretext tasks helps motion forecasting in the following ways:
278 (a) Topology-based context prediction assumes feature similarity or smoothness in small neighbor-
279 hoods of maps, and the resulting feature representation may improve prediction performance. This
280 is mainly expected to help in the first and second settings, which requires generalizing to new topolo-
281 gies. (b) Clustering and classification assumes that feature similarity implies target-label similarity
282 and can group distant nodes with similar features together, leading to better generalization. This is
283 mainly expected to help with dataset imbalance and performance on difficult maneuvers, which re-
284 quires generalizing to hard cases. (c) Supervised learning with imbalanced datasets sees significant
285 degradation in performance. Although most of the data samples in Argoverse are at an intersection,
286 a significantly large number involve driving straight while maintaining speed. Recent studies [45]
287 have shown that SSL tends to learn richer features from more frequent classes, which also allows it
288 to generalize to other classes better. We expect this to help with imbalanced data, limited training
289 data and noisy data.

290 **SSL leads to better generalization compared to pure supervised learning:** To provide evidence
291 for our hypotheses, we design 6 different training and testing setups as shown in Tab. 4. We use
292 Success/Failure classification as the pretext task, and all models are trained for 500,000 steps. We
293 initialize the map-encoder with the parameters from a model trained with the lane mask pretext task.
294 Our *first* setting is to train with 25% of the total data available for training and testing on the full
295 validation set. Our *second* setting assumes that SSL also generalizes to topology from different
296 cities and trains on 100% of data from Pittsburgh (PIT) but only 20% of data from Miami (MIA).
297 For evaluation, we only test on data examples taken from the city of MIA. For our *third* setting,
298 we assume that SSL learns superior features and can thus perform better in difficult cases like lane-
299 changes and turning cases. For evaluation, we only test on data examples which involves these

300 difficult cases. In our *fourth* setting, we choose to explicitly train with data that contains $2\times$ ‘straight-
301 with-same-speed’ maneuver and $1\times$ all other maneuvers. We test only on lane-changes and turning
302 cases from validation. Finally in order to test the effect of noise on motion forecasting performance,
303 we take two models already trained on full data. We now take the full validation set, randomly select
304 agent trajectories or map nodes with probability $p = 0.25$ and $p = 0.5$, and then add Gaussian noise
305 with zero mean and 0.2 variance to their features. There is strong evidence from our experiments that
306 SSL-based tasks provide better generalization and can thus be more effective than pure supervised
307 training.

308 7 Discussion: Potential of this Work

309 We expect this work to influence real world deployment of SSL forecasting methods for autonomous
310 driving. Another use case for this work is realistic behavior generation in traffic simulation. The
311 general construction of the prediction problem, inspired by [2], enables a generic understanding
312 of how an object moves in a given environment without memorizing the training data. A neural
313 network may learn to associate particular areas of a scene with certain motion patterns. To prevent
314 this, we centre around the agent of interest and normalize all other trajectory and map coordinates
315 with respect to it. We predict relative motion as opposed to absolute motion for the future trajectory.
316 This helps to learn general motion patterns. Reconstructing the map or predicting distances from
317 map elements are conducted in a frame-of-reference relative to the agent of interest. This helps
318 in learning general map connectivity. Following work in pedestrian trajectory prediction, we also
319 additionally add random rotations to the training trajectories to reduce directional bias. Furthermore,
320 we provide strong evidence that SSL-based tasks provide better generalization compared to pure
321 supervised training, thereby having the ability to effectively reuse the same prediction model across
322 different scenarios.

323 8 Conclusion

324 We propose SSL-Lanes to leverage supervisory signals generated from data for free in the form of
325 pseudo-labels and integrate it with a standard motion forecasting model. We design four pretext tasks
326 that can take advantage of map-structure and similarities between agent dynamics to generate these
327 pseudo-labels, namely: lane masking, distance to intersection prediction, maneuver classification
328 and success/failure classification. We validate our proposed approach by achieving competitive
329 results on the challenging large-scale Argoverse benchmark. The main advantage of SSL-Lanes is
330 that it has high accuracy combined with low architectural complexity and high inference speed. We
331 further demonstrate that each proposed SSL pretext task improves upon the baseline, especially in
332 difficult cases like left/right turns and acceleration/deceleration. We also provide hypotheses and
333 experiments on why SSL-Lanes can improve motion forecasting.

334 **Limitations:** A limitation of our framework is that it uses the different losses for our formulation
335 only in a 1:1 ratio without tuning them. We also use only one pretext task at a time and do not
336 explore the combination of these different tasks. For our future work, we plan to incorporate meta-
337 learning [46] to identify an effective combination of pretext tasks and automatically balance them—
338 we expect that this will lead to more gains in terms of forecasting performance. Another limitation
339 is that we report improvements with SSL-pretext tasks in scenarios without specifically considering
340 multiple heavily interacting agents. In the future we would like to explore how the interactions
341 between road agents can influence our SSL losses on the interaction split of the Waymo Open Motion
342 dataset (WOMD) [47]. Finally, we explore generalization in terms of implicit data imbalance only in
343 comparison to pure supervised training on the same dataset from which training samples are derived.
344 We would like to study the generalization of our work to other datasets without re-training.

References

- [1] M. Chang, J. Lambert, P. Sangkloy, J. Singh, S. Bak, A. Hartnett, D. Wang, P. Carr, S. Lucey, D. Ramanan, and J. Hays. Argoverse: 3d tracking and forecasting with rich maps. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 8748–8757. Computer Vision Foundation / IEEE, 2019. doi:10.1109/CVPR.2019.00895. URL http://openaccess.thecvf.com/content_CVPR_2019/html/Chang_Argoverse_3D_Tracking_and_Forecasting_With_Rich_Maps_CVPR_2019_paper.html. 2, 6, 7, 8
- [2] M. Liang, B. Yang, R. Hu, Y. Chen, R. Liao, S. Feng, and R. Urtasun. Learning lane graph representations for motion forecasting. In *ECCV*, 2020. 1, 2, 4, 6, 8, 10
- [3] Y. Chai, B. Sapp, M. Bansal, and D. Anguelov. Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction. In L. P. Kaelbling, D. Kragic, and K. Sugiura, editors, *3rd Annual Conference on Robot Learning, CoRL 2019, Osaka, Japan, October 30 - November 1, 2019, Proceedings*, volume 100 of *Proceedings of Machine Learning Research*, pages 86–99. PMLR, 2019. URL <http://proceedings.mlr.press/v100/chai20a.html>. 1, 2, 3, 8
- [4] J. Mercat, T. Gilles, N. E. Zoghby, G. Sandou, D. Beauvois, and G. P. Gil. Multi-head attention for multi-modal joint vehicle motion forecasting. In *2020 IEEE International Conference on Robotics and Automation, ICRA 2020, Paris, France, May 31 - August 31, 2020*, pages 9638–9644. IEEE, 2020. doi:10.1109/ICRA40945.2020.9197340. URL <https://doi.org/10.1109/ICRA40945.2020.9197340>. 1, 2, 6
- [5] H. Zhao, J. Gao, T. Lan, C. Sun, B. Sapp, B. Varadarajan, Y. Shen, Y. Shen, Y. Chai, C. Schmid, C. Li, and D. Anguelov. TNT: target-driven trajectory prediction. In J. Kober, F. Ramos, and C. J. Tomlin, editors, *4th Conference on Robot Learning, CoRL 2020, 16-18 November 2020, Virtual Event / Cambridge, MA, USA*, volume 155 of *Proceedings of Machine Learning Research*, pages 895–904. PMLR, 2020. URL <https://proceedings.mlr.press/v155/zhao21b.html>. 1, 2, 6, 8
- [6] J. Gu, C. Sun, and H. Zhao. Densetnt: End-to-end trajectory prediction from dense goal sets. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 15283–15292. IEEE, 2021. doi:10.1109/ICCV48922.2021.01502. URL <https://doi.org/10.1109/ICCV48922.2021.01502>. 1, 2, 6, 8
- [7] S. Khandelwal, W. Qi, J. Singh, A. Hartnett, and D. Ramanan. What-if motion prediction for autonomous driving. *CoRR*, abs/2008.10587, 2020. URL <https://arxiv.org/abs/2008.10587>. 1, 2, 6
- [8] J. Ngiam, B. Caine, V. Vasudevan, Z. Zhang, H. L. Chiang, J. Ling, R. Roelofs, A. Bewley, C. Liu, A. Venugopal, D. Weiss, B. Sapp, Z. Chen, and J. Shlens. Scene transformer: A unified multi-task model for behavior prediction and planning. *CoRR*, abs/2106.08417, 2021. URL <https://arxiv.org/abs/2106.08417>. 1, 2, 6
- [9] Y. Liu, J. Zhang, L. Fang, Q. Jiang, and B. Zhou. Multimodal motion prediction with stacked transformers. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 7577–7586. Computer Vision Foundation / IEEE, 2021. URL https://openaccess.thecvf.com/content/CVPR2021/html/Liu_Multimodal_Motion_Prediction_With_Stacked_Transformers_CVPR_2021_paper.html. 1, 2, 6, 8
- [10] S. Casas, W. Luo, and R. Urtasun. Intentnet: Learning to predict intention from raw sensor data. In *2nd Annual Conference on Robot Learning, CoRL 2018, Zürich, Switzerland, 29-31 October 2018, Proceedings*, volume 87 of *Proceedings of Machine Learning Research*, pages 947–956. PMLR, 2018. URL <http://proceedings.mlr.press/v87/casas18a.html>. 1, 2

- 394 [11] W. Zeng, M. Liang, R. Liao, and R. Urtasun. Lanercnn: Distributed representations for graph-
395 centric motion forecasting. In *IEEE/RSJ International Conference on Intelligent Robots and*
396 *Systems, IROS 2021, Prague, Czech Republic, September 27 - Oct. 1, 2021*, pages 532–539.
397 IEEE, 2021. doi:10.1109/IROS51168.2021.9636035. URL <https://doi.org/10.1109/IROS51168.2021.9636035>. 1, 2, 6
- 399 [12] J. Gao, C. Sun, H. Zhao, Y. Shen, D. Anguelov, C. Li, and C. Schmid. Vectornet: En-
400 coding HD maps and agent dynamics from vectorized representation. In *2020 IEEE/CVF*
401 *Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA,*
402 *June 13-19, 2020*, pages 11522–11530. Computer Vision Foundation / IEEE, 2020. doi:
403 10.1109/CVPR42600.2020.01154. URL https://openaccess.thecvf.com/content_CVPR_2020/html/Gao_VectorNet_Encoding_HD_Maps_and_Agent_Dynamics_From_Vectorized_Representation_CVPR_2020_paper.html. 1, 2, 3
- 406 [13] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser,
407 and I. Polosukhin. Attention is all you need. In I. Guyon, U. von Luxburg, S. Ben-
408 gio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, editors, *Ad-*
409 *vances in Neural Information Processing Systems 30: Annual Conference on Neural*
410 *Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA,*
411 pages 5998–6008, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>. 1
- 413 [14] Z. Huang, X. Mo, and C. Lv. Multi-modal motion prediction with transformer-based neural
414 network for autonomous driving. *CoRR*, abs/2109.06446, 2021. URL <https://arxiv.org/abs/2109.06446>. 1, 2, 6
- 416 [15] R. Girgis, F. Golemo, F. Codevilla, M. Weiss, J. A. D’Souza, S. E. Kahou, F. Heide, and
417 C. Pal. Latent variable sequential set transformers for joint multi-agent motion prediction.
418 In *10th International Conference on Learning Representations, ICLR 2022, Conference Track*
419 *Proceedings, 2022*. 1, 2, 6
- 420 [16] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging
421 properties in self-supervised vision transformers. In *2021 IEEE/CVF International Confer-*
422 *ence on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages
423 9630–9640. IEEE, 2021. doi:10.1109/ICCV48922.2021.00951. URL <https://doi.org/10.1109/ICCV48922.2021.00951>. 1
- 425 [17] R. E. Kalman and Others. A new approach to linear filtering and prediction problems. *Journal*
426 *of basic Engineering*, 82(1):35–45, 1960. 2
- 427 [18] G. Xie, H. Gao, L. Qian, B. Huang, K. Li, and J. Wang. Vehicle trajectory prediction by
428 integrating physics- and maneuver-based approaches using interactive multiple models. *IEEE*
429 *Trans. Ind. Electron.*, 65(7):5999–6008, 2018. doi:10.1109/TIE.2017.2782236. URL <https://doi.org/10.1109/TIE.2017.2782236>. 2
- 431 [19] A. Houenou, P. Bonnifait, V. Cherfaoui, and W. Yao. Vehicle trajectory prediction based on
432 motion model and maneuver recognition. In *2013 IEEE/RSJ International Conference on*
433 *Intelligent Robots and Systems, Tokyo, Japan, November 3-7, 2013*, pages 4363–4369. IEEE,
434 2013. doi:10.1109/IROS.2013.6696982. URL <https://doi.org/10.1109/IROS.2013.6696982>. 2
- 436 [20] M. Bansal, A. Krizhevsky, and A. S. Ogale. Chauffeurnet: Learning to drive by imitat-
437 ing the best and synthesizing the worst. In A. Bicchi, H. Kress-Gazit, and S. Hutchin-
438 son, editors, *Robotics: Science and Systems XV, University of Freiburg, Freiburg im Breis-*
439 *gau, Germany, June 22-26, 2019*, 2019. doi:10.15607/RSS.2019.XV.031. URL <https://doi.org/10.15607/RSS.2019.XV.031>. 2

- 441 [21] T. Phan-Minh, E. C. Grigore, F. A. Boulton, O. Beijbom, and E. M. Wolff. Covernet:
442 Multimodal behavior prediction using trajectory sets. In *2020 IEEE/CVF Conference*
443 *on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June*
444 *13-19, 2020*, pages 14062–14071. Computer Vision Foundation / IEEE, 2020. doi:
445 10.1109/CVPR42600.2020.01408. URL [https://openaccess.thecvf.com/content_](https://openaccess.thecvf.com/content_CVPR_2020/html/Phan-Minh_CoverNet_Multimodal_Behavior_Prediction_Using_Trajectory_Sets_CVPR_2020_paper.html)
446 [CVPR_2020/html/Phan-Minh_CoverNet_Multimodal_Behavior_Prediction_Using_](https://openaccess.thecvf.com/content_CVPR_2020/html/Phan-Minh_CoverNet_Multimodal_Behavior_Prediction_Using_Trajectory_Sets_CVPR_2020_paper.html)
447 [Trajectory_Sets_CVPR_2020_paper.html](https://openaccess.thecvf.com/content_CVPR_2020/html/Phan-Minh_CoverNet_Multimodal_Behavior_Prediction_Using_Trajectory_Sets_CVPR_2020_paper.html). 2
- 448 [22] A. A. Mohamed, K. Qian, M. Elhoseiny, and C. G. Claudel. Social-stgcn: A social
449 spatio-temporal graph convolutional neural network for human trajectory prediction. In *2020*
450 *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle,*
451 *WA, USA, June 13-19, 2020*, pages 14412–14420. Computer Vision Foundation / IEEE,
452 2020. doi:10.1109/CVPR42600.2020.01443. URL [https://openaccess.thecvf.com/](https://openaccess.thecvf.com/content_CVPR_2020/html/Mohamed_Social-STGCNN_A_Social_Spatio-Temporal_Graph_Convolutional_Neural_Network_for_Human_CVPR_2020_paper.html)
453 [content_CVPR_2020/html/Mohamed_Social-STGCNN_A_Social_Spatio-Temporal_](https://openaccess.thecvf.com/content_CVPR_2020/html/Mohamed_Social-STGCNN_A_Social_Spatio-Temporal_Graph_Convolutional_Neural_Network_for_Human_CVPR_2020_paper.html)
454 [Graph_Convolutional_Neural_Network_for_Human_CVPR_2020_paper.html](https://openaccess.thecvf.com/content_CVPR_2020/html/Mohamed_Social-STGCNN_A_Social_Spatio-Temporal_Graph_Convolutional_Neural_Network_for_Human_CVPR_2020_paper.html). 2
- 455 [23] M. Ye, T. Cao, and Q. Chen. TPCN: temporal point cloud networks for motion forecasting.
456 In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual,*
457 *June 19-25, 2021*, pages 11318–11327. Computer Vision Foundation / IEEE, 2021. URL
458 [https://openaccess.thecvf.com/content/CVPR2021/html/Ye_TPCN_Temporal_](https://openaccess.thecvf.com/content/CVPR2021/html/Ye_TPCN_Temporal_Point_Cloud_Networks_for_Motion_Forecasting_CVPR_2021_paper.html)
459 [Point_Cloud_Networks_for_Motion_Forecasting_CVPR_2021_paper.html](https://openaccess.thecvf.com/content/CVPR2021/html/Ye_TPCN_Temporal_Point_Cloud_Networks_for_Motion_Forecasting_CVPR_2021_paper.html). 2, 6
- 460 [24] H. Song, D. Luan, W. Ding, M. Y. Wang, and Q. Chen. Learning to predict vehicle trajectories
461 with model-based planning. In A. Faust, D. Hsu, and G. Neumann, editors, *Conference on*
462 *Robot Learning, 8-11 November 2021, London, UK*, volume 164 of *Proceedings of Machine*
463 *Learning Research*, pages 1035–1045. PMLR, 2021. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v164/song22a.html)
464 [press/v164/song22a.html](https://proceedings.mlr.press/v164/song22a.html). 2, 6
- 465 [25] W. Zeng, W. Luo, S. Suo, A. Sadat, B. Yang, S. Casas, and R. Urtasun. End-
466 to-end interpretable neural motion planner. In *IEEE Conference on Computer Vi-*
467 *sion and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*,
468 pages 8660–8669. Computer Vision Foundation / IEEE, 2019. doi:10.1109/CVPR.
469 2019.00886. URL [http://openaccess.thecvf.com/content_CVPR_2019/html/Zeng_](http://openaccess.thecvf.com/content_CVPR_2019/html/Zeng_End-To-End_Interpretable_Neural_Motion_Planner_CVPR_2019_paper.html)
470 [End-To-End_Interpretable_Neural_Motion_Planner_CVPR_2019_paper.html](http://openaccess.thecvf.com/content_CVPR_2019/html/Zeng_End-To-End_Interpretable_Neural_Motion_Planner_CVPR_2019_paper.html). 2
- 471 [26] T. Gilles, S. Sabatini, D. Tsishkou, B. Stanciulescu, and F. Moutarde. HOME: heatmap out-
472 put for future motion estimation. In *24th IEEE International Intelligent Transportation Sys-*
473 *tems Conference, ITSC 2021, Indianapolis, IN, USA, September 19-22, 2021*, pages 500–507.
474 IEEE, 2021. doi:10.1109/ITSC48978.2021.9564944. URL [https://doi.org/10.1109/](https://doi.org/10.1109/ITSC48978.2021.9564944)
475 [ITSC48978.2021.9564944](https://doi.org/10.1109/ITSC48978.2021.9564944). 2, 6
- 476 [27] B. Varadarajan, A. Hefny, A. Srivastava, K. S. Refaat, N. Nayakanti, A. Cornman, K. Chen,
477 B. Douillard, C. Lam, D. Anguelov, and B. Sapp. Multipath++: Efficient information fusion
478 and trajectory aggregation for behavior prediction. *CoRR*, abs/2111.14973, 2021. URL [https:](https://arxiv.org/abs/2111.14973)
479 [//arxiv.org/abs/2111.14973](https://arxiv.org/abs/2111.14973). 3
- 480 [28] M. Ye, J. Xu, X. Xu, T. Cao, and Q. Chen. DCMS: motion forecasting with dual consistency
481 and multi-pseudo-target supervision. *CoRR*, abs/2204.05859, 2022. doi:10.48550/arXiv.2204.
482 05859. URL <https://doi.org/10.48550/arXiv.2204.05859>. 3
- 483 [29] N. Deo and M. M. Trivedi. Multi-modal trajectory prediction of surrounding vehicles with
484 maneuver based lstms. In *2018 IEEE Intelligent Vehicles Symposium, IV 2018, Changshu,*
485 *Suzhou, China, June 26-30, 2018*, pages 1179–1184. IEEE, 2018. doi:10.1109/IVS.2018.
486 8500493. URL <https://doi.org/10.1109/IVS.2018.8500493>. 3
- 487 [30] C. Doersch, A. Gupta, and A. A. Efros. Unsupervised visual representation learning by context
488 prediction. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago,*

- 489 Chile, December 7-13, 2015, pages 1422–1430. IEEE Computer Society, 2015. doi:10.1109/
490 ICCV.2015.167. URL <https://doi.org/10.1109/ICCV.2015.167>. 3
- 491 [31] M. Noroozi and P. Favaro. Unsupervised learning of visual representations by solving jigsaw
492 puzzles. In B. Leibe, J. Matas, N. Sebe, and M. Welling, editors, *Computer Vision - ECCV 2016*
493 *- 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings,*
494 *Part VI*, volume 9910 of *Lecture Notes in Computer Science*, pages 69–84. Springer, 2016. doi:
495 10.1007/978-3-319-46466-4_5. URL [https://doi.org/10.1007/978-3-319-46466-4_](https://doi.org/10.1007/978-3-319-46466-4_5)
496 5. 3
- 497 [32] S. Gidaris, P. Singh, and N. Komodakis. Unsupervised representation learning by predicting
498 image rotations. In *6th International Conference on Learning Representations, ICLR 2018,*
499 *Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenRe-
500 view.net, 2018. URL <https://openreview.net/forum?id=S1v4N210->. 3
- 501 [33] M. Caron, P. Bojanowski, A. Joulin, and M. Douze. Deep clustering for unsupervised learn-
502 ing of visual features. In V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, editors,
503 *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, Septem-*
504 *ber 8-14, 2018, Proceedings, Part XIV*, volume 11218 of *Lecture Notes in Computer Sci-*
505 *ence*, pages 139–156. Springer, 2018. doi:10.1007/978-3-030-01264-9_9. URL https://doi.org/10.1007/978-3-030-01264-9_9. 3
- 507 [34] D. Pathak, P. Krähenbühl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature
508 learning by inpainting. In *2016 IEEE Conference on Computer Vision and Pattern Recognition,*
509 *CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2536–2544. IEEE Computer So-
510 ciety, 2016. doi:10.1109/CVPR.2016.278. URL [https://doi.org/10.1109/CVPR.2016.](https://doi.org/10.1109/CVPR.2016.278)
511 278. 3
- 512 [35] R. Zhang, P. Isola, and A. A. Efros. Colorful image colorization. In B. Leibe, J. Matas,
513 N. Sebe, and M. Welling, editors, *Computer Vision - ECCV 2016 - 14th European Con-*
514 *ference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part III*, volume
515 9907 of *Lecture Notes in Computer Science*, pages 649–666. Springer, 2016. doi:10.1007/
516 978-3-319-46487-9_40. URL https://doi.org/10.1007/978-3-319-46487-9_40. 3
- 517 [36] D. Pathak, R. B. Girshick, P. Dollár, T. Darrell, and B. Hariharan. Learning features by watch-
518 ing objects move. In *2017 IEEE Conference on Computer Vision and Pattern Recognition,*
519 *CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 6024–6033. IEEE Computer Society,
520 2017. doi:10.1109/CVPR.2017.638. URL <https://doi.org/10.1109/CVPR.2017.638>. 3
- 521 [37] Y. You, T. Chen, Z. Wang, and Y. Shen. When does self-supervision help graph convolutional
522 networks? In *Proceedings of the 37th International Conference on Machine Learning, ICML*
523 *2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Re-*
524 *search*, pages 10871–10880. PMLR, 2020. URL [http://proceedings.mlr.press/v119/](http://proceedings.mlr.press/v119/you20a.html)
525 you20a.html. 3
- 526 [38] W. Jin, T. Derr, H. Liu, Y. Wang, S. Wang, Z. Liu, and J. Tang. Self-supervised learning
527 on graphs: Deep insights and new direction. *CoRR*, abs/2006.10141, 2020. URL <https://arxiv.org/abs/2006.10141>. 3
- 529 [39] Y. Liu, S. Pan, M. Jin, C. Zhou, F. Xia, and P. S. Yu. Graph self-supervised learning: A survey.
530 *CoRR*, abs/2103.00111, 2021. URL <https://arxiv.org/abs/2103.00111>. 3
- 531 [40] W. Hu, B. Liu, J. Gomes, M. Zitnik, P. Liang, V. S. Pande, and J. Leskovec. Strategies for pre-
532 training graph neural networks. In *8th International Conference on Learning Representations,*
533 *ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL <https://openreview.net/forum?id=HJlWWJSFDH>. 3

- 535 [41] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser,
536 and I. Polosukhin. Attention is all you need. In I. Guyon, U. von Luxburg, S. Ben-
537 gio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, editors, *Ad-
538 vances in Neural Information Processing Systems 30: Annual Conference on Neural
539 Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*,
540 pages 5998–6008, 2017. URL [https://proceedings.neurips.cc/paper/2017/hash/
541 3f5ee243547dee91fbd053c1c4a845aa-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html). 4
- 542 [42] E. A. Abolfathi, M. Rohani, E. Banijamali, J. Luo, and P. Poupart. Self-supervised simultane-
543 ous multi-step prediction of road dynamics and cost map. In *IEEE Conference on Computer
544 Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 8494–8503.
545 Computer Vision Foundation / IEEE, 2021. URL [https://openaccess.thecvf.com/
546 content/CVPR2021/html/Amirloo_Self-Supervised_Simultaneous_Multi-Step_
547 Prediction_of_Road_Dynamics_and_Cost_Map_CVPR_2021_paper.html](https://openaccess.thecvf.com/content/CVPR2021/html/Amirloo_Self-Supervised_Simultaneous_Multi-Step_Prediction_of_Road_Dynamics_and_Cost_Map_CVPR_2021_paper.html). 5
- 548 [43] K. Wagstaff, C. Cardie, S. Rogers, and S. Schrödl. Constrained k-means clustering with back-
549 ground knowledge. In C. E. Brodley and A. P. Danyluk, editors, *Proceedings of the Eighteenth
550 International Conference on Machine Learning (ICML 2001), Williams College, Williamstown,
551 MA, USA, June 28 - July 1, 2001*, pages 577–584. Morgan Kaufmann, 2001. 5
- 552 [44] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In Y. Bengio and
553 Y. LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015,
554 San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL [http:
555 //arxiv.org/abs/1412.6980](http://arxiv.org/abs/1412.6980). 7
- 556 [45] H. Liu, J. Z. HaoChen, A. Gaidon, and T. Ma. Self-supervised learning is more robust to dataset
557 imbalance. *CoRR*, abs/2110.05025, 2021. URL <https://arxiv.org/abs/2110.05025>. 9
- 558 [46] D. Hwang, J. Park, S. Kwon, K. Kim, J. Ha, and H. J. Kim. Self-supervised auxiliary learn-
559 ing with meta-paths for heterogeneous graphs. In H. Larochelle, M. Ranzato, R. Hadsell,
560 M. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems 33: An-
561 nual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December
562 6-12, 2020, virtual*, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/
563 74de5f915765ea59816e770a8e686f38-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/74de5f915765ea59816e770a8e686f38-Abstract.html). 10
- 564 [47] S. Ettinger, S. Cheng, B. Caine, C. Liu, H. Zhao, S. Pradhan, Y. Chai, B. Sapp, C. R. Qi,
565 Y. Zhou, Z. Yang, A. Chouard, P. Sun, J. Ngiam, V. Vasudevan, A. McCauley, J. Shlens, and
566 D. Anguelov. Large scale interactive motion forecasting for autonomous driving : The waymo
567 open motion dataset. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV
568 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 9690–9699. IEEE, 2021. doi:
569 [10.1109/ICCV48922.2021.00957](https://doi.org/10.1109/ICCV48922.2021.00957). URL [https://doi.org/10.1109/ICCV48922.2021.
570 00957](https://doi.org/10.1109/ICCV48922.2021.00957). 10