

# KEMP: Keyframe-Based Hierarchical End-to-End Deep Model for Long-Term Trajectory Prediction

Qiuqing Lu<sup>\*,1,2,†</sup>, Weiqiao Han<sup>\*,1,3,†</sup>, Jeffrey Ling<sup>1</sup>, Minfa Wang<sup>1</sup>, Haoyu Chen<sup>1</sup>,  
Balakrishnan Varadarajan<sup>1</sup>, Paul Covington<sup>1</sup>  
<sup>1</sup>Waymo, <sup>2</sup>UCLA, <sup>3</sup>MIT

*Abstract*—Predicting future trajectories of road agents is a critical task for autonomous driving. Recent goal-based trajectory prediction methods, such as DenseTNT and PECNet [1, 2], have shown good performance on prediction tasks on public datasets. However, they usually require complicated goal-selection algorithms and optimization. In this work, we propose KEMP, a hierarchical end-to-end deep learning framework for trajectory prediction. At the core of our framework is *keyframe-based trajectory prediction*, where keyframes are representative states that trace out the general direction of the trajectory. KEMP first predicts keyframes conditioned on the road context, and then fills in intermediate states conditioned on the keyframes and the road context. Under our general framework, goal-conditioned methods are special cases in which the number of keyframes equal to one. Unlike goal-conditioned methods, our keyframe predictor is learned automatically and does not require hand-crafted goal-selection algorithms. We evaluate our model on public benchmarks and our model ranked 1st on Waymo Open Motion Dataset Leaderboard (as of September 1, 2021).

## I. INTRODUCTION

In order for robots to navigate safely in stochastic environments with multiple surrounding moving agents, predicting future trajectories of surrounding agents is a critical task. In the setting of autonomous driving, the road scene is highly complex, consisting of not only static objects, such as traffic lights and road fences, but also dynamic objects, such as vehicles, pedestrians and cyclists; any vehicle could choose to go straight and pass the intersection, or stop before the intersection and wait for the pedestrians to pass, or make turns. Predicting future trajectories of agents in the scene enables several downstream tasks, such as risk assessment of planned trajectories [3] and safe trajectory planning for autonomous vehicles with theoretical guarantees [4, 5].

Due to the dynamic, stochastic, and interactive nature of the environment, predicting future trajectories of agents based on past observations and the traffic scene is quite challenging. Traditional methods use hand-crafted features and manually-designed logic and models to predict trajectories [6, 7, 8], but they require a great deal of manual work and are brittle to edge cases. On the other hand, modern deep learning methods have successfully demonstrated the

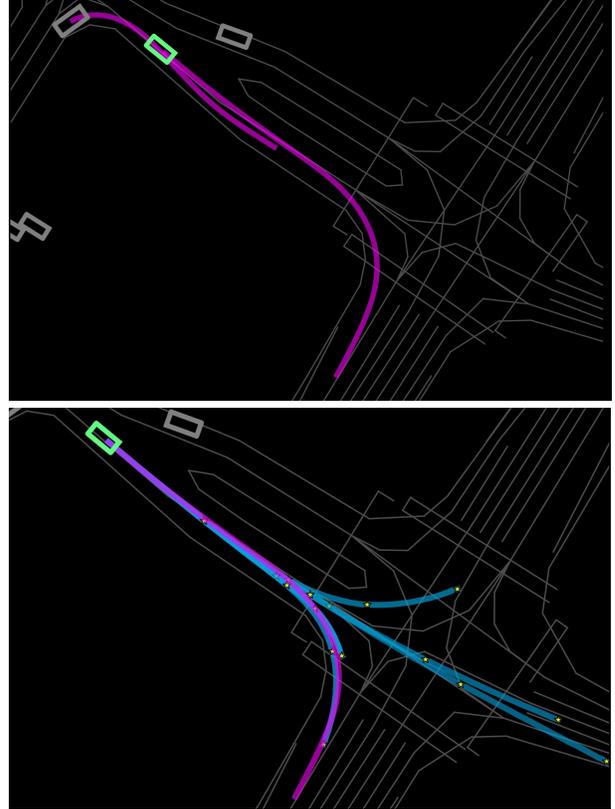


Fig. 1. Agents in an intersection scenario. Top: Agents with ground truth trajectories colored in magenta and the agent for the prediction task represented by a cyan box. Bottom: 6 predicted trajectories colored in blue and keyframes annotated with yellow stars.

ability to scale with larger datasets, including work on graph neural networks [9], long short-term memory (LSTM) [10], generative adversarial networks (GAN) [11], variational autoencoders (VAE) [12, 13], flows [14, 15], or transformers [16] to predict trajectories.

Recently, anchor-based and goal-conditioned methods [17, 18, 1, 14] have received much attention as they directly consider the intention of agents and are more interpretable. However, when making long-term predictions (for example, in Waymo Open Motion Dataset [19], where one needs to predict 8 seconds into the future based on 1 second of past trajectories), only modeling a single high-level goal or intent may not be enough. For one thing, the goal prediction

\* Equal contribution

† Work done during internship at Waymo. Corresponding to weiqiaoh@mit.edu, qiuqing@ucla.edu

This article solely reflects the opinions and conclusions of its authors and not Waymo or any other Waymo entity.

for long trajectories may not be accurate, and for another, trajectories can vary significantly between the fixed starting and goal points. To address this problem, we draw ideas from long-term motion planning and hierarchical reinforcement learning literature [20, 21, 22, 23], where in order for the robot to reach a goal far away or accomplish a complex task, a high level model generates subgoals that are easier for the robot to reach, and a low level model generates control inputs that enable the robot to navigate between two subgoals.

In this paper, we propose a hierarchical end-to-end deep learning framework for autonomous driving trajectory prediction: Keyframe MultiPath (KEMP). At the core of our framework is the keyframe-based trajectory prediction. In this framework, the model first predicts several *keyframes*, which are representative states in the trajectory that trace out the general direction of the trajectory, conditioned on the road context. The model then fills in the gaps between keyframes by predicting intermediate states conditioned on the keyframes and the road context. To our best knowledge, it is the first time that keyframe-based hierarchical prediction is applied to trajectory prediction for autonomous vehicles. Our framework is in some sense a generalization of goal-conditioned trajectory prediction models. In particular, goal-conditioned trajectory prediction models, such as TNT [18], DenseTNT [1], and PECNet [2], can be viewed as special cases of keyframe-based trajectory prediction models where the number of keyframes equals to 1, but unlike these models, we allow the model to learn to predict keyframes instead of manually selecting goals. Other trajectory prediction models that predict trajectories in one shot without conditioning on the final goal can be viewed as special cases of keyframe-based trajectory prediction models where the number of keyframes equals to 0. Our model is not only *more general* than previous methods but also *simpler* as keyframe prediction is learned automatically. Finally, our model achieves *state-of-the-art performance* in autonomous driving trajectory prediction tasks, ranking 1st on Waymo Open Dataset Motion Prediction Leaderboard (as of September 1, 2021).

## II. RELATED WORK

### Latent-variable-sampling-based trajectory prediction.

A popular approach for trajectory prediction is sampling from latent variables. DESIRE [12] generates trajectory samples via a conditional VAE-based RNN encoder-decoder. R2P2 [15] and PRECOG [14] use flows to predict agent futures. SocialGAN [11] uses recurrent generative adversarial networks to predict future trajectories. These methods require stochastic sampling from latent distributions to produce implicit trajectories. The latent variables are not fully interpretable and hence do not work in combination with external prior knowledge.

**Intention-based trajectory prediction.** IntentNet [24] predicts intentions of drivers to guide trajectory prediction. They classify intentions into 8 classes, including keep lane, turn left, turn right, and so on. The method requires a great deal of manual engineering and might miss special cases on

large datasets. Multipath [17] first predicts intents as a set of anchors and then fix the anchors and learn to predict the residual with respect to the anchors.

**Goal-conditioned trajectory prediction.** Goal-conditioned trajectory prediction models are a promising way to develop interpretable autonomous vehicle systems. PECNet [2] predicts the goal as a latent variable and predicts the trajectory conditioned on this latent variable. TNT [18] and DenseTNT [1] predict a set of targets directly and then predict trajectories conditioned on the targets. Compared to latent-variable-sampling-based methods and intention-based methods, goal-conditioned methods such as are more interpretable, because the predicted goal is part of the trajectory instead of a latent variable. Our method can be viewed as a generalization of this line of work, where we predict not only the goal but also other keyframes in the trajectory. Unlike DenseTNT, in which there is a complicated goal-selection algorithm, our method automatically learns to predict keyframes without any hand-crafted engineering.

## III. METHOD

Our method consists of three steps. First, we extract the features from the scene and encode them as context using multiple encoders. Second, we predict keyframes of the output trajectory using a keyframe predictor conditioned on the context. Third, we predict intermediate states conditioned on keyframes and the context using the whole trajectory predictor. The whole model is trained end-to-end.

### A. Context Encoding

To encode the road context, previous work uses rasterized encoding methods [25, 17, 26, 27, 28]. This method renders trajectories of moving agents and road context information as birds-eye view images and encodes them with CNNs. Recently, vector-based representations, which represent the road and agents as polylines, has been more effective in capturing the structural features of high-definition maps [29, 18, 1]. We adopt this vector-based sparse encoding representation. Our context encoding mostly follows the methods in Multipath++ [30].

We use a deep neural network consisting of multiple copies of multi-layer perceptrons (MLPs) and max pooling layers to extract geometric features from road polylines and their connections with each agent. We use PointNet [31, 32] to encode features from 2D points around each agent. Each agent’s raw state, including past positions, velocity, and heading, is encoded using an MLP. Interactions between agents are captured by encoding relative positions and speeds between pairs of agents using MLP and max pooling. All these features are mixed by going through several Multi-Context Gating (MCG) encoders, an efficient mechanism for fusing information [30]. In the end we concatenate all outputs from MCG encoders and get the context embedding  $\mathbf{c}$ .

### B. Keyframe-Based Hierarchical Trajectory Prediction

In the prediction part, given the context  $\mathbf{c}$ , the goal is to predict  $N$  trajectories  $\ell_1, \dots, \ell_N$ . We follow the formulation

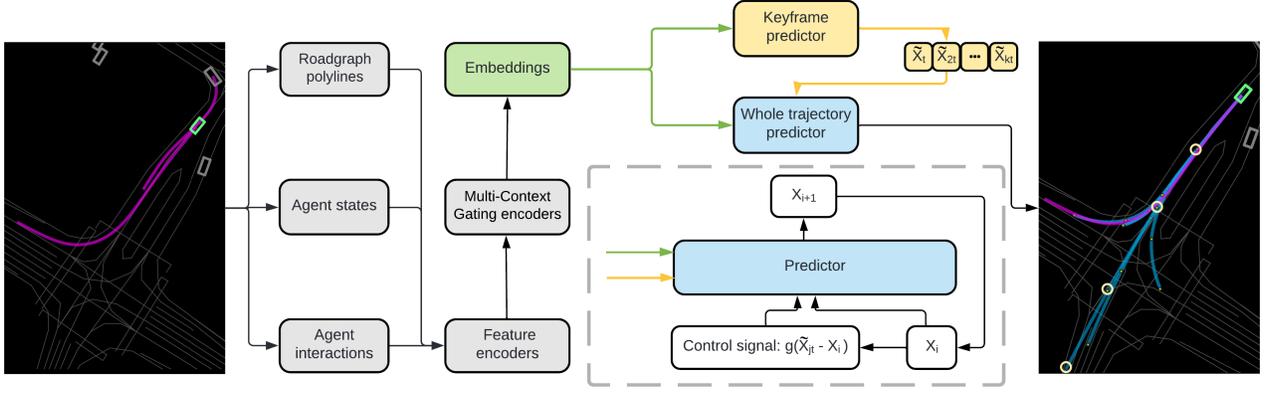


Fig. 2. KEMP architecture. Context features are extracted from scenario inputs with multiple agent historical tracks by multiple encoders. They are then sent to our hierarchical decoders for generating predicted trajectories. The decoder consists of two parts: the keyframe decoder for the generation of keyframe locations and whole trajectory decoder for producing the final whole trajectory based on the previously decoded keyframe locations and context embeddings. In the predictor, we can feed in a control signal  $g(\tilde{X}_{jt} - X_i)$  as a function of the distance to the subgoal.

in MultiPath [17]. Each trajectory is the union of  $T$  states  $\ell_i = \{X_1, \dots, X_T\}$ , and each state  $X_i$  is the tuple  $(\mu_i, \Sigma_i)$ , where  $\mu_i$  is the expectation of the  $(x, y)$  position of the agent at time  $i$ , and  $\Sigma_i$  is the covariance matrix of the position prediction at time  $i$ .

In our proposed method, the keyframe predictor predicts several *keyframes*, which are defined as representative states in the trajectory that trace out the general direction of the trajectory, conditioned on the context  $\mathbf{c}$ . In this paper we focus on *evenly spaced keyframes*. More precisely, suppose  $T = kt$ , where  $T$  is the total number of time steps for the prediction task and  $k, t$  are two positive integers. Then the keyframe predictor predicts  $k$  keyframes  $\tilde{X}_t, \tilde{X}_{2t}, \dots, \tilde{X}_{kt}$  conditioned on the context  $\mathbf{c}$ .

We model the keyframes using a joint distribution

$$\tilde{X}_t, \tilde{X}_{2t}, \dots, \tilde{X}_{kt} \sim p(x_t, x_{2t}, \dots, x_{kt} | \mathbf{c}).$$

We can either use an autoregressive formulation

$$\tilde{X}_{(i+1)t} \sim p(x_{(i+1)t} | \mathbf{c}, x_t, x_{2t}, \dots, x_{it}), i = 0, \dots, k-1.$$

or assume conditional independence between the keyframes

$$\tilde{X}_{(i+1)t} \sim p(x_{(i+1)t} | \mathbf{c}), i = 0, \dots, k-1.$$

In the former, an autoregressive predictor can be implemented with an LSTM, and in the latter, a non-autoregressive predictor can be implemented with a single MLP over all time steps.

Given the  $k$  keyframes  $\tilde{X}_t, \tilde{X}_{2t}, \dots, \tilde{X}_{kt}$ , we consider two ways to generate final trajectories.

1) *Interpolation Model*: For any interval  $[\tilde{X}_{it}, \tilde{X}_{(i+1)t}]$ , the whole trajectory predictor predicts the states inside the interval  $X_{it+1}, \dots, X_{(i+1)t-1}$  conditioned on  $\tilde{X}_{it}$  and  $\tilde{X}_{(i+1)t}$ , as well as the context  $\mathbf{c}$ . This gives us a complete trajectory

$$X_1, \dots, X_{t-1}, \tilde{X}_t, X_{t+1}, \dots, X_{2t-1}, \tilde{X}_{2t}, X_{2t+1}, \dots, \tilde{X}_{kt}.$$

The predictors predict  $N$  trajectories  $\ell_1, \dots, \ell_N$ . We assign a probability to each trajectory  $p_i = p(\ell_i | \mathbf{c}) = \frac{\exp f(\ell_i | \mathbf{c})}{\sum_j \exp f(\ell_j | \mathbf{c})}$ , where  $f(\ell | \mathbf{c})$  is implemented by a deep neural network. Therefore, our prediction is a mixture of Gaussian distribution. We impose the negative log-likelihood loss on the predicted trajectory

$$L_{traj}(\theta) = - \sum_{j=1}^N I(j=r) [\log p(\ell_j | \mathbf{c}; \theta) + \sum_{i=1}^T \log \mathcal{N}(\bar{\mu}_i | \mu_i, \Sigma_i; \theta)],$$

where  $\{\bar{\mu}_1, \dots, \bar{\mu}_T\}$  represents the ground truth trajectory,  $\theta$  represents the parameter to be learned, which is all the weights inside predictor models implemented by deep neural networks, including the whole trajectory predictor and the probability predictor.  $r$  denotes the index of the trajectory that is closest to the ground truth measured by the  $\ell_2$  distance.

2) *Separable Model*: In the interpolation model, the keyframes in the final trajectory are predicted by the keyframe predictor. The whole trajectory predictor does not predict keyframes. In the separable model, the whole trajectory predictor predicts the intermediate states  $X_{it+1}, \dots, X_{(i+1)t}$ , including the keyframes, for any interval  $[\tilde{X}_{it}, \tilde{X}_{(i+1)t}]$ . This gives us a complete trajectory  $X_1, \dots, X_T$  predicted by the whole trajectory predictor. As an aside, when generating intermediate states, the whole trajectory predictor could condition on, in addition to the keyframes, some other manually defined control signals, such as a function of the distance to the subgoal  $g(\tilde{X}_{jt} - X_i)$ ; in practice we use  $g$  as the identity function.

As in the interpolation model, we impose the negative log-likelihood loss  $L_{traj}(\theta)$  on the trajectories predicted by the whole trajectory predictor. Different from the interpolation model, we also impose the consistency loss on the keyframes predicted by the keyframe predictor and the keyframes

predicted by the whole trajectory predictor

$$L_{cons} = \sum_{i=1}^k \|X_{it} - \tilde{X}_{it}\|_2^2.$$

In addition, we impose the negative log-likelihood loss on the keyframes

$$L_{key}(\theta) = - \sum_{j=1}^N I(j=r) \sum_{i=1}^k \log \mathcal{N}(\bar{\mu}_{it} | \mu_{it}, \Sigma_{it}; \theta).$$

The total loss function is a weighted sum of the losses above

$$L = L_{traj} + \alpha L_{cons} + \beta L_{key},$$

where  $\alpha$  and  $\beta$  are weights.

## IV. EXPERIMENTS

### A. Datasets

We evaluate our method on two large-scale real world datasets, the Argoverse Forecasting Dataset and the Waymo Open Motion Dataset.

**Argoverse Forecasting Dataset:** The Argoverse Forecasting Dataset [33] includes 324,557 five seconds tracked scenarios (2s for the past and 3s for the prediction) collected from 1006 driving hours across both Miami and Pittsburgh. Each motion sequence contains the 2D bird’s eye view centroid of each tracked object sampled at 10 Hz. It covers diverse scenarios such as vehicles at intersection, taking turns, changing lanes, and dense traffic. Only one challenging vehicle trajectory is selected as the focus of the forecasting task in each scenario.

**Waymo Open Motion Dataset:** The Waymo Open Motion Dataset (WOMD) [34] is by far the largest interactive motion dataset with multiple types of agents: vehicles, pedestrians and cyclists. It consists of 104,000 run segments with over 1,750 km of roadways and 7.64 million unique agent tracks with 20 seconds duration and sampled at 10 Hz. Each segment is further broken into 9 seconds windows (1s for the past and 8 seconds of future data) with 5 seconds overlap.

### B. Metrics

Given one historical motion,  $K$  predictions are output from a model to compare with the ground truth motion. We used both standard metrics and dataset-specific ranking metrics to evaluate our model’s performance.  $L_2$  distance between a predicted trajectory and the corresponding ground truth is widely used to quantify the displacement error. For the multiple predictions setting, minimum average displacement error (minADE) among all predictions is computed for performance comparison among models. Similarly, minFDE is computed as the minimal  $L_2$  distance among the predicted trajectories and ground truth at the last time step (endpoint). Besides these two standard metrics, Miss Rate (MR) is additionally evaluated, which is the number of scenarios where none of the predicted trajectories are within a certain distance of the ground truth according to the endpoint error divided by the total number of predictions. For WOMD, mean Average Precision (mAP) is designed to measure

precision-recall performance of the future predictions with a normalized averaged over different types of behavior; we use mAP as the primary model metric.

### C. Baseline Algorithms

As our keyframe-based model can be viewed as a generalization of goal-conditioned models, we contrast its performance with the currently top ranked goal-conditioned model, DenseTNT, as well as other strong baseline models.

### D. Implementation Details

**Multiple future predictions:** To obtain more diverse candidate trajectories, our model predicts  $m$  trajectories, more than the required number of predictions. To make the predictions more robust, we trained  $n$  models independently and ensemble their predictions. At inference time, among all  $nm$  candidate trajectories, the top  $K$  trajectories are selected using non-maximum suppression algorithm (NMS), where top ranked trajectories (with highest estimated likelihood) are selected greedily while their nearby trajectories are rejected.

**Model variants:** The keyframe predictor and the whole trajectory predictor both output sequences of states. They can either predict the sequence in one shot, or predict the sequence iteratively in an autoregressive fashion. For both cases, we can use an MLP as the predictor (either one-shot or autoregressive), or an LSTM for the autoregressive case.

**Training details:** Our model is trained with a batch size of 256 on WOMD training dataset and a batch size of 128 on Argoverse. We set loss weights  $\alpha = 10, \beta = 1$  for all models. Network is trained by ADAM optimizer with learning rate  $3 \times 10^{-4}$ , with an exponential decay of 0.5 every 200k steps for WOMD and 100k for Argoverse. Both models are trained on TPU custom hardware accelerator [35] and converged in 3 days on WOMD and 2 days on Argoverse Dataset.

## V. RESULTS

### A. Results on benchmarks

Benchmarks on both datasets are listed in Tables I and V. **Waymo Open Motion Dataset:** This is a more challenging dataset compared to Argoverse Dataset due to the longer prediction duration and more complex scenarios. The models are ranked by mAP. Our best models are:

- 1) KEMP-I-LSTM in Table I: An interpolation model, where the keyframe predictor is implemented by LSTM, and the whole trajectory predictor is implemented by MLP. The number of keyframes is 4.
- 2) KEMP-I-MLP in Table I: An interpolation model, where the keyframe predictor and the whole trajectory predictor are implemented by MLP. The number of keyframes is 4.
- 3) KEMP-S in Table I: A separable model, where the keyframe predictor and the whole trajectory predictor are implemented by LSTM. The number of keyframes is 4.

The first five rows in Table I show the top 5 methods on the Waymo Open Motion Dataset Leaderboard as of September 1st, 2021. KEMP-I-LSTM outperforms baseline models in all



Fig. 3. Samples from WOMB dataset. The agent to be predicted is shown in cyan with its ground truth trajectory shown in magenta. 6 predicted trajectories are shown in blue with yellow stars annotating the keyframes. We compare KEMP-I-LSTM with Multipath. 1st and 3rd rows: Multipath; 2nd and 4th rows: KEMP-I-LSTM.

TABLE I  
MODEL PERFORMANCE ON WAYMO OPEN DATASET (LEADERBOARD)

Model	minADE ↓	minFDE ↓	MR ↓	mAP ↑	mAP(3s) ↑	mAP(5s) ↑	mAP(8s) ↑
DenseTNT <sup>5th</sup> [1]	1.0387	1.5514	0.1779	0.3281	0.4059	0.3195	0.2589
TVN <sup>4th</sup>	0.7498	1.5840	0.1833	0.3341	0.3888	0.3284	0.2852
Scene-Transformer(M+NMS) <sup>3rd</sup>	0.6784	1.3762	0.1977	0.3370	0.3984	0.3317	0.2809
Kraken-NMS <sup>2nd</sup>	0.7407	1.5786	0.2074	0.3561	0.4339	0.3591	0.2754
Multipath++ <sup>1st</sup>	0.5749	1.2117	0.1475	0.3952	0.4710	0.4024	0.3123
KEMP-I-LSTM (ours)	<b>0.5733</b>	<b>1.2088</b>	<b>0.1453</b>	<b>0.3977</b>	<b>0.4729</b>	<b>0.4042</b>	<b>0.3160</b>
KEMP-I-MLP (ours)	<b>0.5723</b>	<b>1.2048</b>	<b>0.1450</b>	0.3968	0.4683	0.4080	0.3141
KEMP-S (ours)	<b>0.5714</b>	<b>1.1986</b>	<b>0.1453</b>	0.3942	0.4729	0.4018	0.3080

TABLE II  
ABLATION STUDY ON WAYMO OPEN DATASET (VALIDATION SET)

Model	minADE ↓	minFDE ↓	MR ↓	mAP ↑	mAP(3s) ↑	mAP(5s) ↑	mAP(8s) ↑
KEMP-I-LSTM	0.5718	1.2061	0.1470	0.3881	0.4735	0.3904	0.3004
KEMP-I-MLP	0.5758	1.2164	0.1487	0.3922	0.4780	0.3995	0.2991
LSTM (No keyframes)	0.5724	1.2099	0.1482	0.3837	0.4676	0.3879	0.2955
MLP (No keyframes)	0.5736	1.2157	0.1493	0.3828	0.4656	0.3892	0.2935

TABLE III  
ABLATION STUDY ON WAYMO OPEN DATASET (VALIDATION SET)

Model	minADE ↓	minFDE ↓	MR ↓	mAP ↑	mAP(3s) ↑	mAP(5s) ↑	mAP(8s) ↑
KEMP-S	0.5691	1.1993	0.1458	0.3940	0.4791	0.3959	0.3071
KEMP-S without $L_{cons}$ loss	0.5698	1.2021	0.1476	0.3949	0.4785	0.4019	0.3043
KEMP-S without $L_{key}$ loss	0.5710	1.2074	0.1467	0.3955	0.4783	0.4009	0.3074
KEMP-S without $L_{cons}$ and $L_{key}$ losses	0.5723	1.2103	0.1484	0.3942	0.4801	0.4018	0.3008

metrics. Additionally, the higher values we achieved in the breakdown of mAP from 3 seconds, 5 seconds and 8 seconds indicate the effectiveness of keyframes as a guidance to the

whole trajectory and demonstrate that our model is able to predict trajectories more accurately in the long-term task. KEMP-S is better than KEMP-I-LSTM in terms of minADE

TABLE IV  
EFFECT OF NUMBER OF KEYFRAMES ON WOMD (VALIDATION SET)

Number of keyframes	minADE↓	minFDE↓	MR ↓	mAP ↑	mAP(3s) ↑	mAP(5s)↑	mAP(8s)↑
0	0.5724	1.2099	0.1482	0.3837	0.4676	0.3879	0.2955
1	0.5720	1.2117	0.1466	0.3945	0.4781	<b>0.4019</b>	0.3034
2	<b>0.5678</b>	<b>1.1993</b>	<b>0.1454</b>	0.3881	0.4726	0.3921	0.2995
4	0.5691	1.1993	0.1458	0.3940	0.4791	0.3959	0.3071
8	0.5715	1.2082	0.1490	<b>0.3963</b>	<b>0.4800</b>	0.3977	<b>0.3110</b>
16	0.5735	1.2136	0.1501	0.3894	0.4742	0.3927	0.3012
40	0.5737	1.2164	0.1487	0.3915	0.4780	0.3959	0.3006

TABLE V  
MODEL PERFORMANCE ON ARGOVERSE DATASET (LEADERBOARD)

Model	minADE ↓	minFDE ↓	MR ↓
TNT [18]	0.94	1.54	13.3%
LaneRCNN [9]	0.90	1.45	12.3%
SenseTime_AP	0.87	<b>1.36</b>	12.0%
Poly	0.87	1.47	12.0%
PRIME [36]	1.22	1.56	11.5%
DenseTNT [1]	0.94	1.49	<b>10.5%</b>
KEMP-I-LSTM	<b>0.85</b>	1.38	12.9%

and minFDE, though it has lower mAP.

Figure 3 shows qualitative results in different scenarios from WOMD validation set. We look at the four examples in the top two rows in detail. In the first case, both models are able to predict diverse modes (turning left, going straight), while the baseline model fails to predict turning right, which is the agent’s actual behavior in the next 8 seconds. In the second case, both models have predicted the correct intent of turning left, but our model has a more natural prediction with keyframes closely aligned to the ground truth. In the third case, our model is able to propose more diverse and reasonable possibilities in the future without missing the mode that agent actually follows. In the fourth case, although both models have diverse predictions spanning the roadgraph space, our model has more reasonable predictions. Compared to the baseline model, KEMP is able to produce more accurate predictions and recall more diverse modes. We believe these good properties are brought by the design of the keyframe architecture. By focusing on the keypoints first to ease the burden of predicting intermediate points, patterns between trajectory and the environment are more easily learned.

**Argoverse Forecasting Dataset:** Table V shows several popular methods on Argoverse Dataset Leaderboard, including TNT, DenseTNT, LaneRCNN, and PRIME. Our model achieves lower minADE and minFDE compared to DenseTNT. In general, we achieve the lowest minADE and second-lowest minFDE among all baseline models, which indicates that KEMP is able to produce realistic trajectories that are very close to the ground truth. However, as the trajectories in Argoverse are fairly short (3 seconds future), our keyframe model does not have a significant advantage over other models.

### B. Ablation studies

First, we compare KEMP against non-keyframe models on the validation set of the Waymo Open Dataset. As shown in

Table II, KEMP-I-MLP has the highest mAP on validation set, though it has also the highest minADE and minFDE. We also run LSTM and MLP models without keyframe prediction as baselines, and observe that their mAP is more than 1% lower than those of KEMP models. This suggests that keyframes have a positive effect on model quality.

Second, we ablate losses from the KEMP-S model in Table III. The fluctuation of mAP among different models is minor: less than 0.3%. The reason might be that the consistency loss and the keyframe loss are complementary given the whole trajectory loss – getting rid of either does not affect the model much. Even after removing both losses, the keyframe predictor can still learn certain features or latent keyframes, because the whole trajectory predictor predicts segments conditioned on the output of the keyframe predictor.

Finally, we vary the number of keyframes in the KEMP-S model as shown in Table IV. Note that our keyframes are equally spaced. So when the number of keyframes is 1, the model becomes goal-conditioned and hence resembles TNT. We observe that 2 keyframes attains the best recall, as the minADE, minFDE, and MR metrics are best. However, the mAP metrics are generally better with 8 keyframes. This indicates that there is a tradeoff that can be made between fewer keyframes, which may increase diversity at the cost of precision, and more keyframes, which provide finer granularity and hence better precision. With too many keyframes, the model may not be taking advantage of the hierarchical structure of the trajectory prediction problem – when we go up to 40 keyframes, for example, metrics become worse. Therefore, depending on whether we care more about precision or recall, the keyframe number can be tuned accordingly.

## VI. CONCLUSION AND FUTURE WORK

In this paper, we proposed a keyframe-based hierarchical end-to-end deep model for long-term trajectory prediction. Our framework generalizes goal-based trajectory prediction methods. Our predictors are automatically learned and does not require hand-crafted algorithms. Our model achieved state-of-the-art performance on the Waymo Open Motion Dataset. Future work could try more complicated structure for the keyframe predictor and the whole trajectory predictor for better performance. Another important direction could be a different definition of keyframes. Currently in our model the keyframes are evenly-spaced states. One could try unevenly-spaced states as keyframes.

## REFERENCES

- [1] J. Gu, C. Sun, and H. Zhao, "DenseTNT: End-to-end Trajectory Prediction from Dense Goal Sets," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 15303–15312.
- [2] K. Mangalam, H. Girase, S. Agarwal, K.-H. Lee, E. Adeli, J. Malik, and A. Gaidon, "It is not the Journey but the Destination: Endpoint Conditioned Trajectory Prediction," in *European Conference on Computer Vision*. Springer, 2020, pp. 759–776.
- [3] A. Wang, X. Huang, A. Jasour, and B. Williams, "Fast Risk Assessment for Autonomous Vehicles Using Learned Models of Agent Futures," *Robotics: Science and Systems (RSS)*, 2020.
- [4] T. Lew, R. Bonalli, and M. Pavone, "Chance-constrained Sequential Convex Programming for Robust Trajectory Optimization," in *2020 European Control Conference (ECC)*. IEEE, 2020, pp. 1871–1878.
- [5] A. Jasour, W. Han, and B. Williams, "Convex Risk Bounded Continuous-Time Trajectory Planning in Uncertain Nonconvex Environments," *Robotics: Science and Systems (RSS)*, 2021.
- [6] N. Deo, A. Rangesh, and M. M. Trivedi, "How would surround vehicles move? A Unified Framework for Maneuver Classification and Motion Prediction," *IEEE Transactions on Intelligent Vehicles*, vol. 3, no. 2, pp. 129–140, 2018.
- [7] K. Yamaguchi, A. C. Berg, L. E. Ortiz, and T. L. Berg, "Who are you with and where are you going?" in *2011 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2011*. IEEE Computer Society, 2011, pp. 1345–1352.
- [8] W.-C. Ma, D.-A. Huang, N. Lee, and K. M. Kitani, "Forecasting Interactive Dynamics of Pedestrians with Fictitious Play," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 774–782.
- [9] W. Zeng, M. Liang, R. Liao, and R. Urtasun, "LaneRCNN: Distributed Representations for Graph-Centric Motion Forecasting," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 532–539.
- [10] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, "Social LSTM: Human Trajectory Prediction in Crowded Spaces," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 961–971.
- [11] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi, "Social GAN: Socially Acceptable Trajectories with Generative Adversarial Networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2255–2264.
- [12] N. Lee, W. Choi, P. Vernaza, C. B. Choy, P. H. Torr, and M. Chandraker, "DESIRE: Distant Future Prediction in Dynamic Scenes with Interacting Agents," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 336–345.
- [13] Y. Yuan and K. M. Kitani, "Diverse Trajectory Forecasting with Determinantal Point Processes," in *International Conference on Learning Representations*, 2019.
- [14] N. Rhinehart, R. McAllister, K. Kitani, and S. Levine, "PRECOG: Prediction Conditioned on Goals in Visual Multi-Agent Settings," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2821–2830.
- [15] N. Rhinehart, K. M. Kitani, and P. Vernaza, "R2P2: A Reparameterized Pushforward Policy for Diverse, Precise Generative Path Forecasting," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 772–788.
- [16] J. Ngiam, B. Caine, V. Vasudevan, Z. Zhang, H.-T. L. Chiang, J. Ling, R. Roelofs, A. Bewley, C. Liu, A. Venugopal *et al.*, "Scene Transformer: A unified multi-task model for behavior prediction and planning," *arXiv preprint arXiv:2106.08417*, 2021.
- [17] Y. Chai, B. Sapp, M. Bansal, and D. Anguelov, "Multipath: Multiple Probabilistic Anchor Trajectory Hypotheses for Behavior Prediction," in *Conference on Robot Learning*. PMLR, 2020, pp. 86–99.
- [18] H. Zhao, J. Gao, T. Lan, C. Sun, B. Sapp, B. Varadarajan, Y. Shen, Y. Shen, Y. Chai, C. Schmid *et al.*, "TNT: Target-driveN Trajectory Prediction," in *Conference on Robot Learning*. PMLR, 2021, pp. 895–904.
- [19] S. Ettinger, S. Cheng, B. Caine, C. Liu, H. Zhao, S. Pradhan, Y. Chai, B. Sapp, C. R. Qi, Y. Zhou *et al.*, "Large Scale Interactive Motion Forecasting for Autonomous Driving: The WAYMO OPEN MOTION DATASET," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 9710–9719.
- [20] A. Wang, T. Kurutach, K. Liu, P. Abbeel, and A. Tamar, "Learning Robotic Manipulation through Visual Planning and Acting," in *Robotics: science and systems*, 2019.
- [21] S. Nair and C. Finn, "Hierarchical Foresight: Self-Supervised Learning of Long-Horizon Tasks via Visual Subgoal Generation," in *International Conference on Learning Representations*, 2019.
- [22] T. Kurutach, A. Tamar, G. Yang, S. J. Russell, and P. Abbeel, "Learning Plannable Representations with Causal InfoGAN," *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [23] O. Nachum, S. S. Gu, H. Lee, and S. Levine, "Data-efficient Hierarchical Reinforcement Learning," *Advances in neural information processing systems*, vol. 31, 2018.
- [24] S. Casas, W. Luo, and R. Urtasun, "InterNet: Learning to Predict Intention from Raw Sensor Data," in *Conference on Robot Learning*. PMLR, 2018, pp. 947–956.
- [25] M. Bansal, A. Krizhevsky, and A. Ogale, "ChauffeurNet: Learning to Drive by Imitating the Best and Synthesizing the Worst," *arXiv preprint arXiv:1812.03079*, 2018.
- [26] F.-C. Chou, T.-H. Lin, H. Cui, V. Radosavljevic, T. Nguyen, T.-K. Huang, M. Niedoba, J. Schneider, and N. Djuric, "Predicting Motion of Vulnerable Road Users using High-Definition Maps and Efficient ConvNets," in *2020 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2020, pp. 1655–1662.
- [27] H. Cui, V. Radosavljevic, F.-C. Chou, T.-H. Lin, T. Nguyen, T.-K. Huang, J. Schneider, and N. Djuric, "Multimodal Trajectory Predictions for Autonomous Driving using Deep Convolutional Networks," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 2090–2096.
- [28] J. Hong, B. Sapp, and J. Philbin, "Rules of the Road: Predicting Driving Behavior with a Convolutional model of Semantic Interactions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8454–8462.
- [29] J. Gao, C. Sun, H. Zhao, Y. Shen, D. Anguelov, C. Li, and C. Schmid, "VectorNet: Encoding HD Maps and Agent Dynamics from Vectorized Representation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11525–11533.
- [30] B. Varadarajan, A. Hefny, A. Srivastava, K. S. Refaat, N. Nayakanti, A. Cornman, K. Chen, B. Douillard, C. P. Lam, D. Anguelov *et al.*, "Multipath++: Efficient Information Fusion and Trajectory Aggregation for Behavior Prediction," *arXiv preprint arXiv:2111.14973*, 2021.
- [31] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
- [32] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space," *Advances in neural information processing systems*, vol. 30, 2017.
- [33] M.-F. Chang, J. Lambert, P. Sangkloy, J. Singh, S. Bak, A. Hartnett, D. Wang, P. Carr, S. Lucey, D. Ramanan *et al.*, "Argoverse: 3D Tracking and Forecasting with Rich Maps," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8748–8757.
- [34] P. Sun, H. Kretschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine *et al.*, "Scalability in Perception for Autonomous Driving: Waymo Open Dataset," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2446–2454.
- [35] N. P. Jouppi, C. Young, N. Patil, D. Patterson, G. Agrawal, R. Bajwa, S. Bates, S. Bhatia, N. Boden, A. Borchers *et al.*, "In-datacenter Performance Analysis of a Tensor Processing Unit," in *Proceedings of the 44th annual international symposium on computer architecture*, 2017, pp. 1–12.
- [36] H. Song, D. Luan, W. Ding, M. Y. Wang, and Q. Chen, "Learning to Predict Vehicle Trajectories with Model-based Planning," in *Conference on Robot Learning*. PMLR, 2022, pp. 1035–1045.