# From Goals, Waypoints & Paths To Long Term Human Trajectory Forecasting

Karttikeya Mangalam[†*]    Yang An[§*]    Harshayu Girase[†]    Jitendra Malik[†]

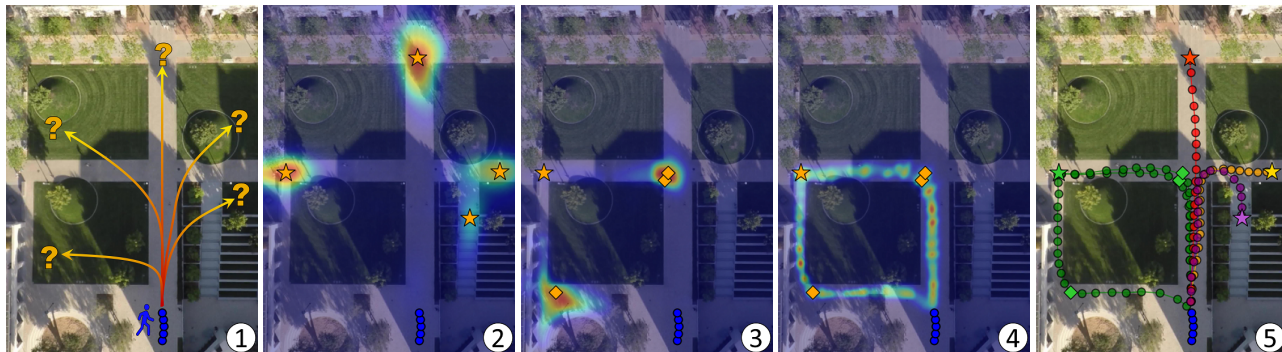[†] UC Berkeley    [§] Technical University of Munich

Figure 1: We tackle the problem of long term human trajectory forecasting. Given the past motion of an agent (blue) on a scene over the last five seconds, we aim to predict the multimodal future motion upto the next minute ①. To achieve this, we propose to factorizing overall multimodality into its *epistemic* and *aleatoric* factors. The *epistemic* factor is modeled with an estimated distribution over the long term goals ② while the *aleatoric* factor is modeled as a distribution over the intermediate waypoints ③ and trajectory ④. This is repeated for multiple goals and waypoints in parallel for scene-compliant multimodal human trajectory forecasting ⑤. Each color indicates predicted trajectories for different sampled goals.

## Abstract

*Human trajectory forecasting is an inherently multimodal problem. Uncertainty in future trajectories stems from two sources: (a) sources that are known to the agent but unknown to the model, such as long term goals and (b) sources that are unknown to both the agent & the model, such as intent of other agents & irreducible randomness in decisions. We propose to factorize this uncertainty into its epistemic & aleatoric sources. We model the epistemic uncertainty through multimodality in long term goals and the aleatoric uncertainty through multimodality in waypoints & paths. To exemplify this dichotomy, we also propose a novel long term trajectory forecasting setting, with prediction horizons upto a minute, an order of magnitude longer than prior works. Finally, we present Y-net, a scene compliant trajectory forecasting network that exploits the proposed epistemic & aleatoric structure for diverse trajectory predictions across long prediction horizons. Y-net significantly improves previous state-of-the-art performance on both (a) The well studied short prediction horizon settings on the Stanford Drone & ETH/UCY datasets and (b) The proposed long prediction horizon setting on the re-purposed Stanford Drone & Intersection Drone datasets.*

## 1. Introduction

Sequence prediction is a fundamental problem in several engineering disciplines such as signal processing, pattern recognition, control engineering and in virtually any domain concerned with temporal measurements. From the seminal work of A. A. Markov [29] on predicting the next syllable in the poem *Eugene Onegin* with Markov chains, to modern day autoregressive descendants like GPT-3 [6], next element prediction in a sequence has a long standing history. Time series forecasting is a key instantiation of the sequence prediction problem in the setting where the sequence is formed by elements sampled in time. Several classic techniques such as Autoregressive Moving Average Models (ARMA) [43] have been incorporated in deep learning architectures [41, 18] in modern day state-of-the-art time series forecasting methods [37].

However, humans are not inanimate Newtonian entities, slave to predetermined physical laws & forces. Predicting the future motion of a billiard ball smoothly rolling on a pool table under friction and physical constraints is a prob-

---

* indicates equal contribution.

lem of different nature from forecasting human motion and positions. Humans are goal conditioned agents that, unlike the ball, exert their will through actions to achieve a desired outcome [40]. Anticipating human motion is of fundamental importance to dynamic agents such as other humans, autonomous robots [3] and self-driving vehicles [39]. Human motion is inherently goal directed and is put in place by the agent to bring about a desired effect.

Nevertheless, even conditioned on the agent's past motion and overarching long term goals, is the future trajectory deterministic? Consider yourself standing at a crossing on a busy street, waiting for the pedestrian light to turn green. While you have every intention of crossing the street, the exact future trajectory remains stochastic as you might swerve to avoid other pedestrians, speed up your pace if the light is about to turn red, or pause abruptly if an unruly cyclist dashes by. Hence, even conditioned on the past observed motion and scene semantics, future human motion is inherently stochastic [16] owing to both *epistemic* uncertainty caused by latent decision variables like long term goals and *aleatoric* variability [11] stemming from random decision variables such as environmental factors. This dichotomy is even sharper in long term forecasting, since due to the increased uncertainty in future, the aleatoric randomness influences the trajectory much more strongly in long rather than short temporal horizons.

This motivates a factorized multimodal approach for human dynamics modeling where both factors of stochasticity are modeled hierarchically rather than lumped jointly. We hypothesize that the long term latent goals of the agent represent the *epistemic* uncertainty with motion prediction. This is motivated by the observation that while the agent has a goal in mind while planning and executing their trajectory, this is unknown to the prediction system. In physical terms, this is akin to the question of *where* the agent wants to go. Similarly, the *aleatoric* uncertainty is expressed in the stochasticity of the path leading to the goal, which encompasses factors like agent's handedness, environment variables such as other agents, partial scene information available to the agent and most importantly, the unconscious randomness in human decisions [19]. In physical terms, this is akin to the question of *how* the agent reaches the goal.

Hence, we propose to model the *epistemic* uncertainty first and then the *aleatoric* stochasticity conditioned on the obtained estimate. Concretely, with the RGB scene and the past motion history we first estimate an explicit probability distribution over the agent's final positions at the end of the trajectory, *i.e.* the agent's long term goals. This represents the *epistemic* uncertainty in the prediction system. We also estimate distributions over a few chosen future waypoint positions which along with the sampled goal points are used to obtain explicit probability maps over all the remaining trajectory positions. This represents the aleatoric uncertainty

in the prediction system. Together the samples from the *epistemic* goal and the *aleatoric* waypoint & trajectory distribution form the predicted future trajectory.

In summary, our contribution is threefold. **First**, we propose a novel long term prediction setting that extends upto a minute in the future which is about an order of magnitude longer than previous literature. We also benchmark performance of previous state-of-the-art short horizon prediction models on this setting along with simple baselines. **Second**, we propose Y-net, a scene-compliant long term trajectory prediction network that explicitly models both the *goal* and *path* multimodalities by making effective use of the scene semantics. **Third**, we show that the factorized multimodality modeling enables Y-net to improve the state-of-the-art both on the proposed long term settings and the well-studied short term prediction settings. We benchmark Y-net's performance on the Stanford Drone [31] and the ETH [30]/UCY[23] benchmark in the short term setting, where it outperforms previous approaches by significant margins of 26.9% & 5.6% respectively, on ADE and by 34.0% and 51.9% respectively, on FDE metric. Further, we also study Y-net's performance in the proposed long term prediction setting on the Stanford Drone & the Intersection Drone Dataset [5] where it substantially improves the performance of state-of-the-art short term methods by over 50.6% and 35.0% respectively, on ADE and 77.1% and 55.9% respectively, on FDE metric.

## 2. Related Works

Several recent studies have investigated human trajectory prediction in different settings. Broadly, these approaches can be grouped on the basis of the proposed formulation for multimodality in forecasting, inputs signals available to the prediction model and the nature and form of prediction results furnished by the model. Several diverse input signals such as agent's past motion history [17], human pose [27], RGB scene image [14, 35, 8, 22, 26], scene semantic cues [8], location [36, 24, 4] & gaze of other pedestrian [27, 46] in the scene, moving vehicles such as cars [36] and also latent inferred signals such as agent's goals [28] have been used. The form of prediction results produced are also diverse with multimodality [26] and scene-compliant forecasting being central to the prior works.

### 2.1. Unimodal Forecasting

Early trajectory forecasting work focused on unimodal predictions of the future. Social Forces [17] proposes modeling interactions as attractive and repulsive forces and future trajectory as a deterministic path evolving under these forces. Social LSTM [1] focuses on other agents in the scene and models their effects through a novel pooling module. [46] tackles motion forecasting in ego-centric views and develops a system that exploits subtle cues like body

pose and gaze along with camera wearer's ego-motion for other agent's future location prediction. [42] proposes to use attention to model current agent's interaction with other agent's. [27] predicts trajectory as the 'global' branch for pose prediction and proposes to condition downstream tasks such as pose prediction on predicted unimodal trajectories.

## 2.2. Multimodality through Generative Modeling

A line of work aims to model the stochasticity inherent in future prediction through a latent variable with a defined prior distribution through approaches such as conditional variational auto-encoders [20]. Lee *et al*. [22] propose DE-SIRE, an inverse reinforcement learning based approach that uses multimodality in sampling of a latent variable that is ranked and optimized with a refinement module. CF-VAE [4] uses normalizing flows with VAEs for modeling structure in sequences such as trajectories. [27] introduces the use of a CVAE for capturing multimodality in the final position of the pedestrians conditioned on the past motion history. Trajectron++ [36] represents agent's trajectories in a graph structured recurrent network for scene complaint trajectroy forecasting, taking into account the interaction with a diverse set of agents. CGNS [24] uses variational divergence minimization procedure in multimodal latent space to learn feasible regions for future trajectories.

A different line of work includes Social GAN [14] which uses adversarial losses [13] for incorporating multimodality in predictions. SoPhie [35] further incorporates attention modules to model agent's interactions with the environment and other agents.

While such generative approaches do produce diverse trajectories, overall coverage of critical modes cannot be guaranteed and little control is afforded over the properties of predicted trajectories such as direction, number of samples *etc*. In contrast, our method, Y-net, estimates explicit probability maps which allow easily incorporating spatial constraints for a downstream task.

## 2.3. Multimodality through spatial probability estimates

Another line of work obtains multimodality via estimated probability maps. Activity Forecasting from Kitani *et al*. [21] proposes to use a hidden Markov Decision process for modeling the future paths. However, in contrast to our work the future predictions in [21] are conditioned on the activity label such as 'approach car', 'depart car' *etc*. More recently, some works have used a grid based scene representation for estimating probabilities for future time steps [25, 26, 10]. Relatedly, some prior works such as [27, 47, 8] propose a goal-conditioned trajectory forecasting method. However, no prior works have proposed factorized modeling of *epistemic* uncertainty or goals & *aleatoric* uncertainty or paths as Y-net uses.

## 3. Proposed Method

The problem of multimodal trajectory prediction can be formulated formally as follows. Given a RGB scene image $\mathcal{I}$ and past positions of a pedestrian in the scene $\mathcal{I}$ denoted by $\{\mathbf{u}_n\}_{n=1}^{n_p}$ for the past $t_p = n_p/\texttt{FPS}$ seconds sampled at the frame rate $\texttt{FPS}$, the model aims to predict the position of the pedestrian for the next $t_f$ seconds in the future, denoted by $\{\mathbf{u}_n^i\}_{n=n_p}^{n_p+n_f}$ where $t_f = n_f/\texttt{FPS}$.

Since the future is stochastic, multiple predictions for the future trajectories are produced. In this work, we factorize the overall stochasticity into two modes. First are the modes relating to *epistemic* uncertainty *i.e.* multimodality in the final destination of the agent for which the module produces $K_e$ predictions. Second are the modes relating to the *aleatoric* uncertainty *i.e.* multimodality in the path taken to the destination stemming from uncontrolled randomness given the goal, for which the module produces $K_a$ separate predictions for each given destination. In the short temporal horizon limit, since the overall path length is small, the options for paths to a given goal are limited and similar to each other. Hence, this results in setting $K_a = 1$ and so the total paths predicted ($K$ in prior works) is the same as $K_e$. However, for longer temporal horizons, there are several paths to the same goal and hence $K_a > 1$. Next, we describe in detail the working of our model, Y-net and its three sub-networks $\mathbf{U}_e$, $\mathbf{U}_g$ and $\mathbf{U}_t$ (Section 3.1) followed by details of the non-parametric sampling process (Section 3.2) and loss functions (Section 3.3) used.

### 3.1. Y-net Sub-Networks

To effectively use the scene information in the semantic space with the trajectory information expressed in coordinates, alignment needs to be created between the different signals. Some prior works [35] achieve this by encoding the two-dimensional RGB image $\mathcal{I}$ as a one-dimensional hidden state vector extracted from some pretrained network. While this provides the network with scene information, any meaningful spatial signal gets highly conflated when flattened into a vector and pixel alignment is destroyed. This is highlighted in [28] which establish previous state-of-the-art without any RGB information underlining the misuse of image information in prior works. In this work, we adopt a trajectory on scene heatmap representation that solved the alignment issue by representing the trajectory spatially in the same two-dimensional space as image $\mathcal{I}$.

#### 3.1.1 Trajectory on Scene Heatmap Representation

The RGB image $\mathcal{I}$ is first processed with a semantic segmentation network such as U-net [33] that produces segmentation map $\mathbf{S}$ of $\mathcal{I}$ comprising of $N_c$ classes determined according to the affordance provided by the surface to an agent for actions such as walking, standing, running etc. In
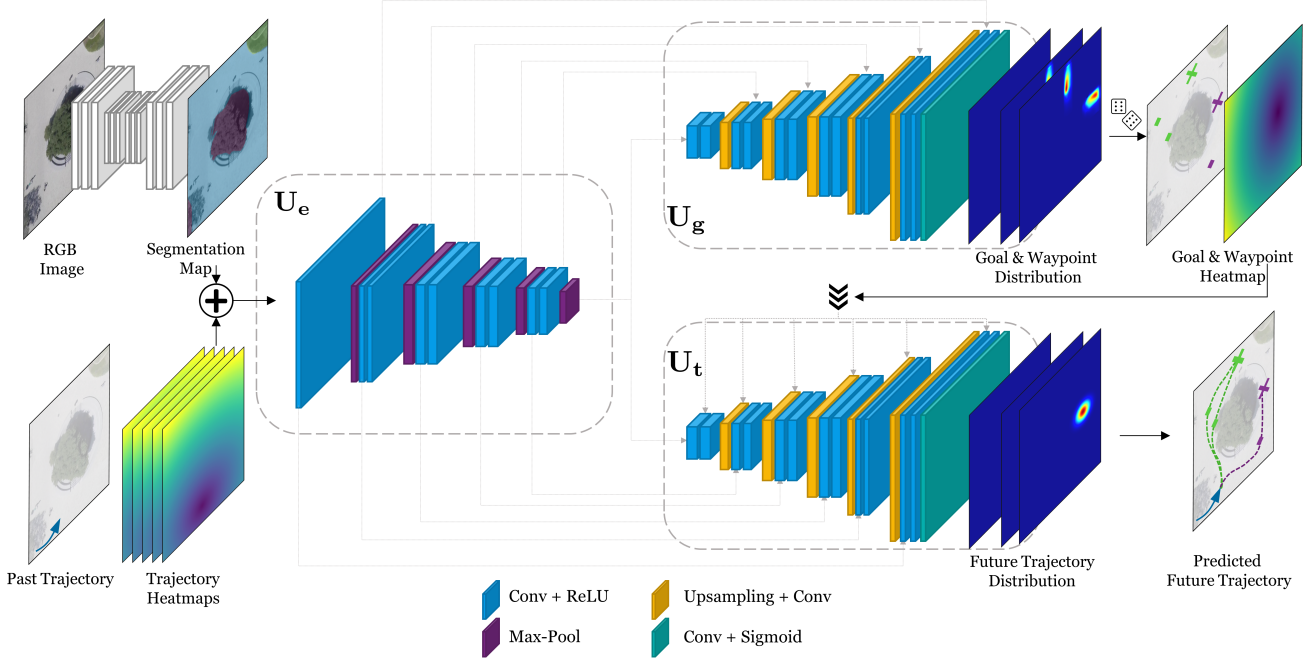
Figure 2: **Model Architecture**: Y-net comprises of three sub-networks $\mathbf{U}_e$, $\mathbf{U}_g$ & $\mathbf{U}_t$ modeled after the Unet architecture [33] (Section 3.1). Y-net adopts a factorized approach to multimodality, expressing the stochasticity in goals & waypoints through estimated distributions furnished by $\mathbf{U}_g$. And multimodality in paths is achieved through estimated probability distributions obtained by $\mathbf{U}_t$ conditioned on samples from $\mathbf{U}_g$ for predicting diverse multimodal scene-compliant futures.

a parallel branch, the past motion history $\{\mathbf{u}_n\}_{n=1}^{n_p}$ of agent $p$ is converted to a trajectory heatmap $\mathbf{H}$ of spatial sizes of $\mathcal{I}$ and $n_p$ channels corresponding to the past $t_p$ seconds sampled at the frame rate. Mathematically,

$$\mathbf{H}(n, i, j) = 2 \frac{\|(i,j) - \mathbf{u}_n\|}{\max\limits_{(x,y)\in\mathcal{I}} \|(x,y) - \mathbf{u}_n\|}$$

The heatmap trajectory representation is then concatenated with the semantic map $\mathcal{S}$ along the channel dimension producing the trajectory on scene heatmap tensor $\mathbf{H}_{\mathcal{S}}$ a $H \times W \times (N_c + n_p)$ dimensional input tensor which is passed to the encoder network $\mathbf{U}_e$.

### 3.1.2 Trajectory on Scene Heatmap Encoder

The tensor $\mathbf{H}_{\mathcal{S}}$ is processed with the encoder $\mathbf{U}_e$ designed as a U-net encoder [33] as shown in Figure 2. The encoder $\mathbf{U}_e$ consists of a total of $N_{\mathbf{U}_e}$ blocks, where the spatial dimensions are reduced from $H \times W$ to $H_{\mathbf{U}} \times W_{\mathbf{U}}$ halving after every block using max pooling and channel depth is increased sequentially from $N_c + n_p$ to $C_{\mathbf{U}_e}$ doubling after a certain number of blocks. The final spatially compact and deep representation $\mathbf{H}_{\mathbf{U}_e}$ along with the $N_{\mathbf{U}_e} - 1$ intermediate feature tensors of varying spatial resolution are passed

onto the goal decoder $\mathbf{U}_g$ and the trajectory decoder $\mathbf{U}_t$ as discussed below.

### 3.1.3 Goal & Waypoint Heatmap Decoder

The feature maps of $U_e$ at various spatial resolutions are passed onto the goal & waypoint heatmap decoder $\mathbf{U}_g$ which is modeled after the expansion arm in the U-net architecture [33]. After a center block consisting of two convolutional layers which takes the most compact feature map $\mathbf{H}_{\mathbf{U}_e}$ from $\mathbf{U}_e$, each block of the expansion arm starts by expanding the previous feature map, spatially doubling the resolution in every block using bilinear up-sampling and convolution (together forming Deconvolution [33]).

Further, after every Deconvolution, the corresponding intermediate representations from $\mathbf{U}_e$ are merged and the features are passed through two convolution layers. Merging intermediate high resolution feature maps from $\mathbf{U}_e$ is necessary since just using $\mathbf{H}_{\mathbf{U}_e}$ would severely limit the final resolution of the goal heatmap, thus missing fine spatial details that are preserved in the intermediate feature maps.

Deconvolution, feature merging and convolution form a U-net block. $\mathbf{U}_g$ consists of $N_{\mathbf{U}_g}$ blocks, followed by a per-pixel sigmoid that produces $N^w + 1$ estimated spatial heatmaps. This includes the probability distribution of

the final goal of the agent *i.e.* the non-parametric distribution $\hat{\mathbb{P}}(\mathbf{u}_{n_p+n_f})$ and $N^w$ intermediate waypoint probability heatmaps at frame steps $w_i \in \{n_p, n_p + 1, \ldots, n_p + n_f\} \; \forall \; i = 1, \ldots, N^w$ represented as non-paramateric distributions $\hat{\mathbb{P}}(\mathbf{u}_{w_i})$.

### 3.1.4 Trajectory Heatmap Decoder

The estimated goal and waypoint distributions are sampled as described in Section 3.2 and the obtained sample $\hat{\mathbf{u}}_{n_p+n_f}$ for the goal and $\{\hat{\mathbf{u}}_{w_i}\}_{i=1}^{N^w}$ for the intermediate waypoints are converted to a heatmap representation as described in Section 3.1. The obtained conditioning tensor $\mathbf{H}_{\mathbf{U}_g}$ is spatially downsampled to match the corresponding block's size. Those are passed along with the past motion and scene representation $\mathbf{H}_{\mathbf{U}_e}$ and the $N_{\mathbf{U}_e} - 1$ intermediate high resolution tensors to the trajectory decoder network $\mathbf{U}_t$. $\mathbf{U}_t$ is modeled after the expansion arm of a U-net as well and proceeds in a similar fashion as $\mathbf{U}_g$, as described in Section 3.1.3 with a total of $N_{\mathbf{U}_t}$ expansion blocks. The final trajectory distribution is obtained with a channel independent per-pixel sigmoid that produces for each future timestep a separate heatmaps of spatial size $H \times W$ corresponding to the position of the agent in the scene over the next $n_f$.

## 3.2. Non-parametric Distribution Sampling

Given a distribution $\mathbb{P}$ of the position of the agent in a future frame represented non-parametrically as a matrix of probabilities $X$, we aim to sample a two-dimensional point as our estimate for the position of the agent at that time step. This is difficult to achieve reliably in practice since the estimated distribution $\mathbb{P}$ is noisy. Hence, taking a naive $\texttt{argmax}$ is not robust. Instead we propose to use the $\texttt{softargmax}$ operation [12] to approximate the most likely position in a robust & stable fashion. Mathematically,

$$\texttt{softargmax}(X) = \left( \sum_i i \frac{\sum_j e^{X_{ij}}}{\sum_{i,j} e^{X_{ij}}}, \sum_j j \frac{\sum_i e^{X_{ij}}}{\sum_{i,j} e^{X_{ij}}} \right)$$

### 3.2.1 Test-Time Sampling Trick (TTST)

In situations where multiple samples are required, such as during testing, sampling from $\mathbb{P}$ can be carried in a straightforward fashion by considering $X$ as a categorical distribution with given probabilities $X(i, j)$ for position $(i, j)$. However, this approach doesn't take into account the number of samples required from $\mathbb{P}$ all of which jointly will represent the quality of the estimated $\mathbb{P}$. For example, if only a few samples are required, sampling indiscriminately spatially is sub-optimal since samples are likely to be wasted if drawn from low probability regions or sampled from adjacent regions.

We propose a 'Test-Time Sampling Trick' (TTST) that is cognizant of the number of goal samples, $K_e$, needed for evaluation. During testing, we propose to first sample a large number of points (10,000 in our experiments) from the estimated distribution $\mathbb{P}$ in a $K_e$-agnostic manner.

To eliminate outliers, we suppress samples from pixels $(i, j)$ with probability $X_{ij}$ below a threshold $thr_{rel}$. It is set adaptively for each probability matrix $X$ separately to a fraction of the highest occurring value in that matrix: $thr_{rel} = \max(X) * 0.01$.

Further, to control the tradeoff between diversity and precision, we use the hyper parameter temperature $T$. Before the pixel-wise sigmoid operation, we divide the predicted logit probability map through $T$. Lower temperature values results in probability maps $X$ with low entropy, *i.e.* the probability mass is concentrated in a smaller number of pixels and samples are increasingly drawn from more likely regions. Higher temperature values increase the diversity of samples. For the short term setting, we use $T = 1.0$, while for our proposed long term setting, we increase the diversity by setting $T = 1.8$.

Then, we propose to run the fast clustering algorithm K-means on these sampled points with the number of clusters set to $K_e - 1$. The cluster centers obtained from the K-means algorithm, along with the $\texttt{softargmax}$ sampled point, form the final set of $K$ samples to be used for evaluation. Note that while in spirit this is similar to the 'truncation trick' proposed in PECNet [28], the 'truncation trick' is K-agnostic and requires a well suited $\sigma_T$ to be chosen experimentally beforehand for a given $K$. Further, their 'truncation trick' operates in the latent variable space with no direct control over the final generated samples since in [28] multi-modality is introduced through implicit approaches like Variational Auto-encoders. Alleviating these limitations, TTST provides direct control on the sampled points, is cognizant of the number of samples needed $K$ and hence, does not require any $K$-specific tuning because of the design choice of using explicit probability heatmaps.

### 3.2.2 Conditioned Waypoint sampling

Goal and waypoints are dependent to each other. Figure 3 shows an example, where the road forks into three different paths. If the sampled goal lies on the bottom path, the agent most likely won't pass through a waypoint on the upper or middle path, and will prefer a waypoint on the bottom road.

Following this intuition, we introduce a hierarchical prior to condition the waypoint sampling on the already sampled goal and waypoints.

We first sample the goal. By fixing the goal the possible locations of the next waypoint is constrained. We assume that the waypoint lies on a straight line segment between the sampled goal and the past trajectory at $\frac{w_i - n_p}{n_f}$ of the

Figure 3: **Conditioned Waypoint sampling**: The first image shows the scene with the past trajectory in blue. The yellow star indicates the sampled goal. The unconditioned waypoint distribution can be seen in the second image; the distribution is goal-agnostic and therefore probability mass is distributed over all three roads. The third image is the resulting waypoint distribution, by multiplying the multivariate Gaussian prior with the unconditioned prediction. The final image shows the trajectory towards the sampled goal while crossing the waypoint.

line segment's length, *e.g.* with $N^w = 1$ waypoint, it would lie in the middle of the segment. This assumption is too hard and can lead to unrealistic paths not complying with the environment constraints. To relax this assumption, we use a multivariate Gaussian prior with mean at the assumed location. The variance is chosen adaptively by considering the distance between the agent's last observed position and the sampled goal as $\sigma_\perp = ||u_{n_p+n_f} - u_{n_p}||/\alpha$ where $\alpha$ is a scaling hyper parameter. We set $\alpha = 6$ in our experiments. Intuitively, the greater the distance between the current position and the sampled goal, the more uncertain is the waypoint position. $\sigma_\perp$ is the variance perpendicular to the line segment, and we set the variance parallel to the line segment to $\sigma_\parallel = \beta * \sigma_\perp$ with $\beta = 0.5$ in our experiments. This constrains the possible waypoint position more in the direction of travel and leaves more room for uncertainty in the perpendicular direction.

We multiply the described multivariate Gaussian prior pixel-wise to the predicted waypoint distribution. Fusing prior and predicted distribution leads to scene-compliant waypoints in the correct direction. As seen in the example, it suppresses probability mass on the upper and middle road. From the resulting distribution we use `softargmax` for the first waypoint and sample the remaining $K_a - 1$ waypoints randomly to get two-dimensional points from the distribution.

If there is more than one waypoint, *i.e.* $N^w > 1$, we repeat the above process for the next waypoint at $w_i$ and condition it to the previously sampled waypoint at $w_{i+1}$.

### 3.3. Loss Function

Since we use the trajectory on scene representation (Section 3.1) we impose losses directly on the estimated distribution $\hat{\mathbb{P}}$ rather than on the samples. The ground truth future is represented as a Gaussian heatmap $\mathbb{P}$ centered at the observed points with a predetermined variance $\sigma_H$. All three networks, $\mathbf{U}_e$, $\mathbf{U}_g$ and $\mathbf{U}_t$ are trained end to end jointly using a weighted combination of binary cross entropy losses on the predicted goal, waypoint and trajectory distributions.

$$\mathcal{L}_{\text{goal}} = \text{BCE}(\mathbb{P}(\mathbf{u}_{n_p+n_f}), \hat{\mathbb{P}}(\mathbf{u}_{n_p+n_f}))$$

$$\mathcal{L}_{\text{waypoint}} = \sum_{i=1}^{N^w} \text{BCE}(\mathbb{P}(\mathbf{u}_{w_i}), \hat{\mathbb{P}}(\mathbf{u}_{w_i}))$$

$$\mathcal{L}_{\text{trajectory}} = \sum_{i=n_p}^{n_p+n_f} \text{BCE}(\mathbb{P}(\mathbf{u}_i), \hat{\mathbb{P}}(\mathbf{u}_i))$$

$$\mathcal{L} = \mathcal{L}_{\text{goal}} + \lambda_1 \mathcal{L}_{\text{waypoint}} + \lambda_2 \mathcal{L}_{\text{trajectory}}$$

## 4. Results

### 4.1. Datasets

We use a total of three datasets to study Y-net's performance – the Stanford Drone Dataset (SDD) [32] (for both short & long term), the Intersection Drone Dataset (InD) [5] (for long term only) and the ETH [30] / UCY [23] forecasting benchmark (for short term only).

**Stanford Drone Dataset**: We benchmark our proposed model on the popular Stanford drone dataset [32] where several recently proposed methods have improved state-

6

of-the-art performance significantly in the past few years [45]. The dataset is comprised of more than $11,000$ unique pedestrians across 20 top-down scenes captured on the Stanford university campus in bird's eye view using a flying drone. It has over $40,000$ agent-scene interactions and has enjoyed very popular usage in trajectory prediction literature in the short temporal horizon setting. For short term prediction, we follow the well-established preprocessed data and splits from the TrajNet benchmark [34] setup used by [14, 35, 28, 10], sampling at $\texttt{FPS} = 2.5$ yielding an input sequence of length $n_p = 8$ (3.2 seconds) and output of length $n_f = 12$ (4.8 seconds).

In our proposed long term setting, we split the raw data of Stanford Drone Dataset (SDD) in the same fashion as proposed in TrajNet benchmark [34] evaluating on the same scenes, all of which are not seen during training. The raw data is recorded in $\texttt{FPS} = 30$ and we first downsample the data to our proposed $\texttt{FPS} = 1$, thus yielding a $n_p = 5$ for $t_f = 5$ seconds in the past and predicting upto one minute into the future. We use the middle point of the raw bounding boxes to get the same coordinate representation as the preprocessed short term setting. The data contains various types of agent beyond pedestrians (bicyclists, skateboarders, cars, buses, and golf carts), we filter out all non-pedestrians and short trajectories below $n_p + n_f$ out. As the raw data is noisy and contains temporal discontinuities, we split the trajectories at those discontinuities. We use a sliding window approach without overlap to split up long trajectories, resulting in our final dataset. After those steps, the dataset contains 1502 trajectories for $n_p = 5$ and $n_f = 30$. Further, we label the scenes with semantic segmentation maps consisting of the following $N_c = 5$ "stuff" classes [7] depending on the affordability of class to pedestrians: pavement, terrain, structure, tree and road.

**Intersection Drone Dataset**: We propose to use the Intersection drone dataset [5] for benchmarking trajectory forecasting in long horizon settings. The dataset comprises over 10 hours of measurements over 4 distinct intersection in an urban environment. We use similar steps for the Intersection Drone Dataset as for SDD. To evaluate Y-net and the baselines performance on unseen scenes during training, we only use location ID 4 during testing. The raw video and detection are in $\texttt{FPS} = 25$, and again, we downsample the data to $\texttt{FPS} = 1$. We then filter out non-pedestrians and short trajectories and use a sliding window approach without overlap to split long trajectories. Since the data lies in world coordinates, we convert it into pixel coordinates by scaling with the provided scale factors from the authors. After the preprocessing steps, inD contains 1,396 long term trajectories with $n_p = 5$ and $n_f = 30$. To evaluate performance on unseen environments, we are using location ID 4 only during testing time. The scene is labeled with the same $N_c = 5$ classes as in SDD.

**ETH & UCY datasets**: The ETH/UCY benchmarks have been widely used for benchmarking trajectory forecasting models in the short horizon setting in the recent years [44]. Forecasting performance has improved by over $\sim 64\%$ on average, within the last two years itself [14].

It comprises of five different scenes all of which report position in world coordinates (in meters). We follow the leave one out validation strategy as outlines in prior works [14, 35, 10, 36]. We use the preprocessed data from [14] [1], also used by state-of-the-art methods [28, 36]. Similar to SDD, the frames are sampled at $\texttt{FPS} = 2.5$ predicting $n_f = 12$ frames, $t_f = 4.8$ seconds into the future given last $n_p = 8$ frames comprising of $t_p = 3.2$ seconds of the past.

Our model represents trajectories as heatmaps and hence needs the coordinates in pixel space. The ETH/UCY data lie in world coordinates. To project the data from meter into pixel space we use the provided homography matrices from the original dataset for the ETH [30] scenes ETH and HOTEL and create our own homography matrices for the UCY[23] scenes UNIV, ZARA1 and ZARA2. To enable fair comparisons, we convert our predictions back to world coordinates using the inverse homography matrices and calculate our errors with the untouched raw data in world coordinates, to avoid any errors from our projection.

For all ETH and & UCY datasets, since the classes of affordances furnished by the surfaces present is small, we use $N_c = 2$, identifying each pixel as either belonging to class 'road' or 'not road'.

## 4.2. Implementation Details

### 4.2.1 Segmentation Model Implementation Details

To incorporate constraints and interactions of the agents with the scene, we pretrain a semantic segmentation model to efficiently use the sparse scene image data. Stanford Drone Dataset contains 60 scene images in total, while inD only contains images from four recording locations. We use the U-net model [33] with ResNet101 [15] backbone. The ResNet101 encoder's weights are pretrained on ImageNet, while the weights for the U-net decoder and segmentation head are randomly initialized. The images are downsampled by a factor of four (SDD) and three (inD), padded to be divisible by 32 as required for U-net and cropped to $256 \times 256$. The data is augmented spatially by rotation, flipping, scaling and perspective transformation and we introduce Gaussian noise, blurring as well as color, brightness and contrast shifts. The semantic maps are manually labeled into the $N_c = 5$ classes mentioned above, as well as a dummy class for padding and black areas in the inD dataset. We only use the corresponding images from the trajectory train scenes for training to evaluate the performance of Y-net on unseen environments for both SDD and inD.

---

[1] https://github.com/agrimgupta92/sgan

| | S-GAN | CF-VAE | P2TIRL | SimAug | PECNet* | Y-net (Ours) | DESIRE | TNT* | PECNet | Y-net (Ours) |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | $K = 20$ | | | | | $K = 5$ | | |
| ADE | 27.23 | 12.60 | 12.58 | 10.27 | 9.96 | **7.85** | 19.25 | 12.23 | 12.79 | **11.49** |
| FDE | 41.44 | 22.30 | 22.07 | 19.71 | 15.88 | **11.85** | 34.05 | 21.16 | 29.58 | **20.23** |

Table 1: **Short temporal horizon forecasting results on SDD**: Our method significantly outperforms previous state-of-the-art methods (indicated by *) on the Stanford Drone Dataset [32] on both the ADE & FDE metrics for both settings of $K$, where $K$ represents the number of multimodal samples . Reported errors are in pixels with $t_p = 3.2$ sec, $t_f = 4.8$ sec, $n_p = 8, n_f = 12$. Lower is better.

| | S-GAN | | Sophie | | CGNS | | PECNet | | Trajectron++* | | Y-net (Ours) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ADE | FDE | ADE | FDE | ADE | FDE | ADE | FDE | ADE | FDE | ADE | FDE |
| ETH | 0.81 | 1.52 | 0.70 | 1.43 | 0.62 | 1.40 | 0.54 | 0.87 | 0.39 | 0.83 | **0.28** | **0.33** |
| HOTEL | 0.72 | 1.61 | 0.76 | 1.67 | 0.70 | 0.93 | 0.18 | 0.24 | 0.12 | 0.21 | **0.10** | **0.14** |
| UNIV | 0.60 | 1.26 | 0.54 | 1.24 | 0.48 | 1.22 | 0.35 | 0.60 | **0.20** | 0.44 | 0.24 | **0.41** |
| ZARA1 | 0.34 | 0.69 | 0.30 | 0.63 | 0.32 | 0.59 | 0.22 | 0.39 | **0.15** | 0.33 | 0.17 | **0.27** |
| ZARA2 | 0.42 | 0.84 | 0.38 | 0.78 | 0.35 | 0.71 | 0.17 | 0.30 | **0.11** | 0.25 | 0.13 | **0.22** |
| AVG | 0.58 | 1.18 | 0.54 | 1.15 | 0.49 | 0.97 | 0.29 | 0.48 | 0.19 | 0.41 | **0.18** | **0.27** |

Table 2: **Short term forecasting results on ETH/UCY benchmark**: Our proposed method establishes new state-of-the-art results (previous results denoted by *) on both the ADE & the FDE metrics on the popular ETH-UCY benchmark using standard short-horizon settings (same as SDD) and $K = 20$. Reported errors are in meters. Lower is better.

The SDD segmentation model is trained using ADAM optimizer to reduce the Dice Loss [38] with an initial learning rate of $1 \times 10^{-4}$ and batch size of 4. The learning rate is decreased to $1 \times 10^{-5}$ after 1500 epochs. Further we freeze the ResNet101 backbone for the first 200 epochs.

As inD only contains images from four locations, we use the pretrained SDD model and freeze the encoder for the first 1000 epochs to avoid catastrophic forgetting. All other hyper parameters are the same as for training SDD.

### 4.2.2 Y-net Implementation Details

We train the entire network end to end with ADAM optimizer with a learning rate of $1 \times 10^{-4}$ and batch size of 8. We scale the overall loss by a factor of 1000. Since the scene images $\mathcal{I}$ have different heights and widths in all datasets, we ensure that each batch only contains image and trajectories from the same scene. Y-net does not use fully-connected layers and therefore can handle images of different sizes, without cropping or padding to the same shape. The RGB scene image $\mathcal{I}$ and trajectory heatmaps $\mathbf{H}$ are downsampled by 4 for SDD, 3 for inD and 1.5 for

ETH/UCY to save memory and padded to be divisible by 32. For fair comparisons with previous methods we upsample the predicted trajectories back to its original size and compare with the ground-truth data in original scale. All scene images and trajectories are augmented by spatial flipping and rotation in 90° steps, increasing the number of training data by a factor of eight. The encoder blocks in $\mathbf{U}_e$ have output channel dimensions $[32, 32, 64, 64, 64]$ and both $\mathbf{U}_g$ and $\mathbf{U}_t$ start with two convolutional layers of output channels 128, followed by blocks of output channel dimensions $[64, 64, 64, 32, 32]$. We use $\lambda_1 = \lambda_2 = 1$ to weight the binary cross entropy loss. The ground-truth trajectory heatmaps has a variance of $\sigma_H = 4$ pixels.

During training $\mathbf{U}_g$ predicts the goal and waypoint distribution for all $n_p$ time steps as an auxiliary task. This helps to let the sub-network learn the dynamics of pedestrian trajectories better. During inference, we only use the goal and $N^w$ waypoint distributions as needed.

The trajectory sub-network $\mathbf{U}_t$ is trained using the ground-truth goal and waypoints. Those are represented as trajectory heatmaps as described in Section 3.1.1 and downsampled spatially to fit the corresponding feature map

| | Stanford Drone Dataset | | | | | | Intersection Drone Dataset | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | S-GAN | PECNet | R-PECNet | Y-net (Ours) | | | S-GAN | PECNet | R-PECNet | Y-net (Ours) | | |
| $K_a$ | 1 | 1 | 1 | 1 | 2 | 5 | 1 | 1 | 1 | 1 | 2 | 5 |
| ADE | 155.32 | 72.22 | 261.27 | 47.94 | 44.94 | 39.49 | 38.57 | 20.25 | 341.80 | 14.99 | 14.02 | 12.67 |
| FDE | 307.88 | 118.13 | 750.42 | 66.71 | 66.71 | 66.71 | 84.61 | 32.95 | 1702.64 | 21.13 | 21.13 | 21.13 |

Table 3: **Long term trajectory forecasting Results**: We benchmark performance on our proposed long horizon forecasting setting predicting $t_f = 30$ seconds into the future given $t_p = 5$ seconds past motion history. All reported error are in pixels (lower is better) for $K_e = 20$ with additional results for varying $K_a$ with a fixed $K_e$.

shapes of $\mathbf{U}_t$ blocks. By using the ground-truth, $\mathbf{U}_t$ learns to predict trajectories leading towards the goal, while passing the waypoints. During inference, we use the (TTST) sampled goals and waypoints predicted by $\mathbf{U}_g$.

On ETH/UCY, we further experiment with deformable convolutional layers as proposed in [9].

We will release all our data, both raw & processed, code for reproducing experiments across all the datasets & labeled semantic maps for reproducibility and future work.

### 4.3. Metrics

We use the established Average Displacement Error (ADE) and Final Displacement Error (FDE) metrics for measuring performance of future predictions. ADE is calculated as the $\ell_2$ error between the predicted future and the ground truth averaged over the entire trajectory while FDE is the $\ell_2$ error between the predicted future and ground truth for the final predicted point [2]. Following prior works [14], in the case of multiple future predictions, the final error is reported as the `min` error over all predicted futures.

### 4.4. Baseline models

We benchmarks against several state-of-the-art methods across both short and long term trajectory forecasting settings which we describe briefly.

- Social GAN [14] proposes a Generative Adversarial Network (GAN) to predict multi-modal trajectory autoregressively.

- SoPhie [35] also uses a GAN and extends it by attention modules to incorporate other agents and scene context.

- Conditional Flow VAE (CF-VAE) [4] proposes a conditional normalizing flow based Variational Auto-Encoder that models future uncertainty without disentangling underlying factors.

- Conditional Generative Neural System (CGNS) [24] proposes variational divergence minimization in latent space to learn feasible regions for future trajectories.

- P2TIRL [10] proposes a grid based trajectory forecasting method learnt using maximum entropy inverse reinforcement learning.

- DESIRE [22] also proposes an inverse reinforcement learning approach for prediction by planning.

- SimAug [25] is a recently proposed method that uses additional adversarially generated 3D multi-view data for adapating to novel viewpoints in forecasting and improve the Multiverse model [26].

- PECNet [28] is the prior state-of-the art method on short-term trajectory prediction on the Stanford Drone Dataset. They propose to use goal-conditioning but does not account for multi-modality in the path to the goal.

- TNT [47] closely improves upon PECNet's performance for $K = 5$ samples on SDD and is the prior state-of-the-art in that setting.

- Trajectron++ [36] proposes a recurrent graph based forecasting model incorporating dynamic constrains such as other moving agents and scene information. This work also hold the prior state-of-the art on ETH/UCY short-term trajectory prediction benchmark.

### 4.5. Short Term Forecasting Results

#### 4.5.1 Stanford Drone Results

Table 1 presents results on SDD in the short term setting *i.e.* $t_p = 3.2$ seconds, $t_f = 4.8$ seconds. We follow the standard split from [34] and report results with $K_e = 5$ & 20. Since there's limited aleatoric multimodality in short term settings, we use $K_a = 1$ thus being comparable to prior

| | SDD | | | inD | |
|---|---|---|---|---|---|
| TTST | ✗ | ✗ | ✓ | ✗ | ✓ |
| CWS | ✗ | ✓ | ✓ | ✓ | ✓ |
| ADE | 65.00 | 52.31 | 47.94 | 17.77 | 14.99 |
| FDE | 86.98 | 86.98 | 66.71 | 28.52 | 21.13 |

Table 4: **Ablation results for Conditioned Waypoint Sampling (`CWS`) and `TTST`**: We benchmark the performance of Y-net with and without our proposed `CWS` and `TTST` on our long horizon forecasting setting, predicting $t_f = 30$ seconds into the future given $t_p = 5$ seconds of past motion history. All reported errors are in pixels (lower is better) for $N^w = 1$, $K_e = 20$ and $K_a = 1$.

works using $K = 20$ trajectory samples for final evaluation. Table 1 shows our proposed model achieving an ADE of 7.85 and FDE of 11.85 at $K_e = 20$ which outperforms the previous state-of-the-art performance of PECNet [28] by 26.8% on ADE and 34.0% on FDE. Further, at $K = 5$ it achieves an ADE of 11.49 & FDE of 20.23 outperforming previous state-of-the-art performance of TNT.

### 4.5.2 ETH & UCY Results

We also report results on the ETH/UCY benchmark in Table 2. Similar to SDD, we set $K_e = 20, K_a = 1$. Here as well, we observe that our proposed model achieves an ADE of 0.18 & FDE of 0.27 improving the previously state-of-the-art performance from Trajectron++ [36] of 0.19 ADE & 0.41 FDE by about 5.6% ADE and 51.9% FDE.

### 4.6. Long Term Forecasting Results

To study the effect of *epistemic* & *aleatoric* uncertainty, we propose a long term trajectory forecasting setting with a prediction horizon upto 10 times longer than prior works, extending to a minute. To benchmark, we retrain PECNet [28], the previous state-of-the-art method from short term forecasting & Social GAN [14] for each $t_f$ in the long horizon setting. We also train a recurrent short term baseline using the PECNet model (R-PECNet) where the model is trained only for $t_f = 5$ seconds and is fed its own predictions recurrently for predicting longer temporal horizons.

### 4.6.1 Forecasting Results

Table 3 reports the baseline and our results on the Stanford Drone (SDD) and Intersection Drone Datasets (InD) for a time horizon of $t_f = 30$ seconds in the future given the
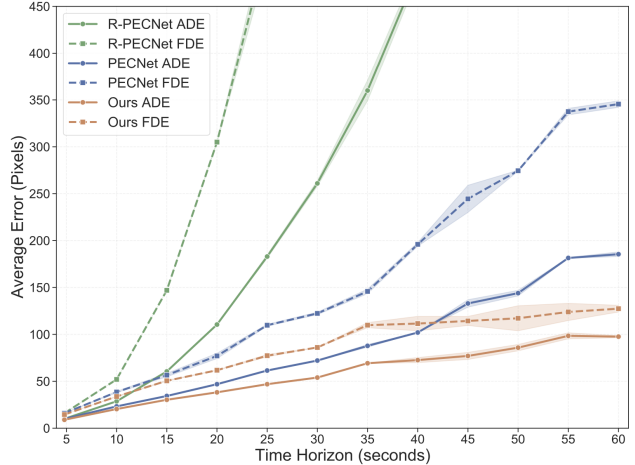


Figure 4: **Benchmarking Performance against Time Horizons**: On prediction horizons upto a minute, we observe a consistently growing difference in ADE between Y-net and PECNet, highlighting the important of factorized goal & path modeling in long term forecasting.
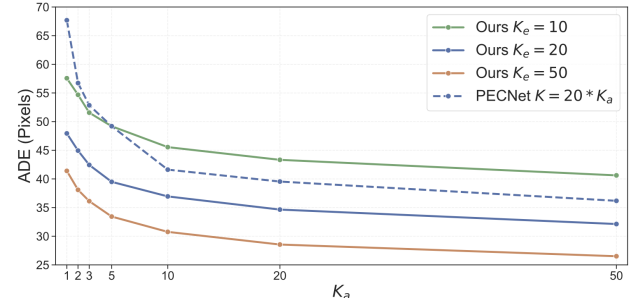


Figure 5: **Benchmarking performance against aleatoric uncertainty** ($K_a$): Fixing the goal multimodality ($K_e$) we vary $K_a$ to observe the effect of path multimodality. Also, we benchmark against PECNet by allowing it 20 times more samples for each $K_a$ for a fair compare against the $K_e = 20$ Y-net curve.

past $t_p = 5$ seconds input. All reported results are with $K_e = 20$ for Y-net conditioned on $N_w = 1$ intermediate waypoint at $w_i = 20$, *i.e.* midway temporally between the observed inputs and the estimated goal. All reported baseline results are at $K = 20$ for fair comparisons with our $K_e = 20, K_a = 1$ setting. On SDD, we observe that our proposed model outperforms the state-of-the-art short term baseline on the long horizon setting as well, achieving an ADE of 47.94 and FDE 66.72 improving upon PECNet's performance by over 50%. Similarly, Y-net outperforms PECNet on InD improving ADE performance from 20.25 to 14.99 and FDE from 32.95 to 21.13.
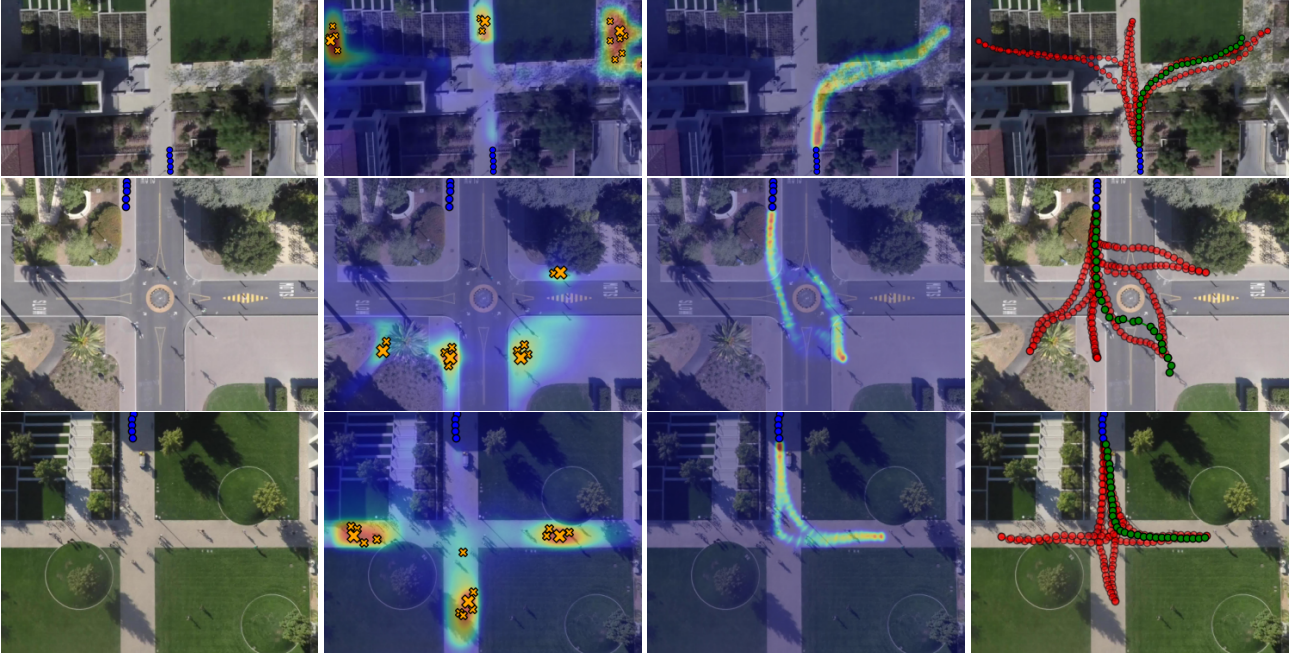
Figure 6: **Qualitative Long Term Trajectory Forecasting Results**: We show various heatmaps and visualizations for three different scenes (rows) in SDD testset. The first column shows the past observed trajectory for last $t_p = 5$ seconds in blue. The second column shows the heatmap from $\mathbf{U}_g$ for $t_f = 30$ seconds in the future (goal multimodality) and some sampled goals from the estimated distribution. The third column shows trajectory heatmaps from $\mathbf{U}_t$ conditioned on a sampled goal from column three (path multimodality). The last column shows the predicted trajectories, green indicating the ground-truth trajectories & red our multimodal predictions.

### 4.6.2 Conditoned Waypoint Sampling

Table 4 shows an ablation of the Conditioned Waypoint Sampling `CWS`. While it doesn't affect the goal sampling, *i.e.* the FDE, the ADE decreases by 24.3%.

### 4.6.3 Test-Time Sampling Trick

Table 4 shows the effectiveness of our proposed `TTST`. `TTST` reduces the error on SDD and inD by 9.1% and 18.5% in ADE, respectively, and 30.4% and 35.0% in FDE.

### 4.6.4 Varying Prediction Horizon

We also compare Y-net with two flavours of PECNet, one retrained separately for each prediction horizon $t_f$, and another trained only for $t_f = 5$ seconds but evaluated recurrently for long horizons (R-PECNet). We observe that the difference in ADE between Y-net & PECNet grows as prediction horizon increases from 5 to 60 seconds. This shows Y-net's adaptability for long prediction horizons owing to factorized mulitmodality modeling. We also observe that for PECNet, training a separate model for different time horizons is significantly better than using a short temporal horizon model recurrently. This motivates our proposal

for studying long term forecasting since short term models behave very poorly when applied out of box recurrently to longer term settings.

### 4.6.5 Varying $K_a$

We also report results with $K_a = 2 \& 5$ for studying the improvement in performance from aleatoric multimodality in Table 3. We observe a consistent improvement in ADE on both datasets, thus indicating the diversity in predicted paths given the same estimated final goal $\mathbf{u}_{n_p+n_f}$. We also report extensive results for varying the path multimodality $K_a$ with a fixed $K_e$ for various choice of $K_e \& K_a$ in Figure 5. Additionally for baselining, we benchmark against PECNet [27] evaluated with $K_e$ times more samples than the corresponding Y-net model while varying $K_a$. We show consistent ADE improvements for various $K_e$ while increasing $K_a$, indicating effective use of multimodality. Further, even with $K = K_e * 20$ samples, *i.e.* 20 times more samples for each $K_e$, PECNet's performance is significantly worse than Y-net at $K_e = 20$ for all $K_a$ highlighting the importance of factorizing goal and path multimodality for diverse & accurate future trajectory modeling.

epistemic and aleatoric uncertainty dichotomy. In this setting, we benchmark on the Stanford Drone & Intersection Drone dataset where Y-net exceeds previous state-of-the-art by over $77.1\%$ and $55.9\%$ respectively thereby highlighting the importance of modeling factorized stochasticity.

## References

[1] Alexandre Alahi, Kratarth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–971, 2016.

[2] Alexandre Alahi, Vignesh Ramanathan, and Li Fei-Fei. Socially-aware large-scale crowd forecasting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2203–2210, 2014.

[3] Maren Bennewitz, Wolfram Burgard, and Sebastian Thrun. Learning motion patterns of persons for mobile service robots. In *Proceedings 2002 IEEE International Conference on Robotics and Automation (Cat. No. 02CH37292)*, volume 4, pages 3601–3606. IEEE, 2002.

[4] Apratim Bhattacharyya, Michael Hanselmann, Mario Fritz, Bernt Schiele, and Christoph-Nikolas Straehle. Conditional flow variational autoencoders for structured sequence prediction. *arXiv preprint arXiv:1908.09008*, 2019.

[5] Julian Bock, Robert Krajewski, Tobias Moers, Steffen Runde, Lennart Vater, and Lutz Eckstein. The ind dataset: A drone dataset of naturalistic road user trajectories at german intersections. 2019.

[6] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.

[7] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Cocostuff: Thing and stuff classes in context. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1209–1218, 2018.

[8] Zhe Cao, Hang Gao, Karttikeya Mangalam, Qizhi Cai, Minh Vo, and Jitendra Malik. Long-term human motion prediction with scene context. In *Eur. Conf. Comput. Vis.* 2020.

[9] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 764–773, 2017.

[10] Nachiket Deo and Mohan M Trivedi. Trajectory forecasts in unknown environments conditioned on grid-based plans. *arXiv preprint arXiv:2001.00735*, 2020.

[11] Armen Der Kiureghian and Ove Ditlevsen. Aleatory or epistemic? does it matter? *Structural safety*, 31(2):105–112, 2009.

Figure 7: **GIF Visualization**: Demonstrating the goal, waypoint and path multimodality for long term human trajectory prediction (30 seconds horizon). Given the past 5 seconds input history (green), we predict diverse future trajectories (current location in orange, past in red). *Best viewed in Adobe Acrobat Reader.*

### 4.6.6 Qualitative Results

We show qualitative results for long term trajectory prediction ($t_f = 30$) on SDD in Figure 6 and through a GIF temporally in Figure 7. We observe that Y-net predicts diverse scene-complaint trajectories, with both future goal and path multimodalities.

## 5. Conclusion

In summary, we present Y-net, a scene-compliant trajectory forecasting network with factorized goal and path multimodalities. Y-net uses the U-net structure [33] for explicitly modeling probability heatmaps for epistemic and aleatoric uncertainties. Overall, Y-net improves previous state-of-the-art performance by $34.0\%$ on the SDD and by $51.9\%$ on ETH/UCY benchmarks in the short term setting. We also propose a new long term trajectory forecasting setting (prediction horizon upto a minute) for exemplifying the

[12] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 6.2. 2.3 softmax units for multinoulli output distributions. In *Deep Learning.*, pages 180–184. MIT Press, 2016.

[13] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014.

[14] Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social gan: Socially acceptable trajectories with generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2255–2264, 2018.

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[16] Dirk Helbing. *Stochastische Methoden, nichtlineare Dynamik und quantitative Modelle sozialer Prozesse*. Shaker, 1993.

[17] Dirk Helbing and Peter Molnar. Social force model for pedestrian dynamics. *Physical review E*, 51(5):4282, 1995.

[18] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[19] Daniel Kahneman. *Thinking, fast and slow*. Macmillan, 2011.

[20] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[21] Kris M Kitani, Brian D Ziebart, James Andrew Bagnell, and Martial Hebert. Activity forecasting. In *European Conference on Computer Vision*, pages 201–214. Springer, 2012.

[22] Namhoon Lee, Wongun Choi, Paul Vernaza, Christopher B Choy, Philip HS Torr, and Manmohan Chandraker. Desire: Distant future prediction in dynamic scenes with interacting agents. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 336–345, 2017.

[23] Alon Lerner, Yiorgos Chrysanthou, and Dani Lischinski. Crowds by example. In *Computer graphics forum*, volume 26, pages 655–664. Wiley Online Library, 2007.

[24] Jiachen Li, Hengbo Ma, and Masayoshi Tomizuka. Conditional generative neural system for probabilistic trajectory prediction. *arXiv preprint arXiv:1905.01631*, 2019.

[25] Junwei Liang, Lu Jiang, and Alexander Hauptmann. Simaug: Learning robust representations from 3d simulation for pedestrian trajectory prediction in unseen cameras. 2020.

[26] Junwei Liang, Lu Jiang, Kevin Murphy, Ting Yu, and Alexander Hauptmann. The garden of forking paths: Towards multi-future trajectory prediction, 2020.

[27] Karttikeya Mangalam, Ehsan Adeli, Kuan-Hui Lee, Adrien Gaidon, and Juan Carlos Niebles. Disentangling human dynamics for pedestrian locomotion forecasting with noisy supervision, 2020.

[28] Karttikeya Mangalam, Harshayu Girase, Shreyas Agarwal, Kuan-Hui Lee, Ehsan Adeli, Jitendra Malik, and Adrien Gaidon. It is not the journey but the destination: Endpoint conditioned trajectory prediction. *arXiv preprint arXiv:2004.02025*, 2020.

[29] A. A. Markov. An example of statistical investigation of the text eugene onegin concerning the connection of samples in chains. *Lecture at the physical-mathematical faculty, Royal Academy of Sciences, St. Petersburg*, 23 January 1913.

[30] Stefano Pellegrini, Andreas Ess, and Luc Van Gool. Improving data association by joint modeling of pedestrian trajectories and groupings. In *European conference on computer vision*, pages 452–465. Springer, 2010.

[31] A Robicquet, A Sadeghian, A Alahi, and S Savarese. Learning social etiquette: Human trajectory prediction in crowded scenes. In *European Conference on Computer Vision (ECCV)*.

[32] Alexandre Robicquet, Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. Learning social etiquette: Human trajectory understanding in crowded scenes. In *European conference on computer vision*, pages 549–565. Springer, 2016.

[33] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In Nassir Navab, Joachim Hornegger, William M. Wells, and Alejandro F. Frangi, editors, *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, pages "234–241", Cham, 2015. Springer International Publishing.

[34] A Sadeghian, V Kosaraju, A Gupta, S Savarese, and A Alahi. Trajnet: Towards a benchmark for human trajectory prediction. *arXiv preprint*, 2018.

[35] Amir Sadeghian, Vineet Kosaraju, Ali Sadeghian, Noriaki Hirose, Hamid Rezatofighi, and Silvio Savarese. Sophie: An attentive gan for predicting paths compliant to social and physical constraints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1349–1358, 2019.

[36] Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and Marco Pavone. Trajectron++: Dynamically-feasible trajectory forecasting with heterogeneous data. *arXiv preprint arXiv:2001.03093*, 2020.

[37] Rajat Sen, Hsiang-Fu Yu, and Inderjit S Dhillon. Think globally, act locally: A deep neural network approach to high-dimensional time series forecasting. In *Advances in Neural Information Processing Systems*, pages 4837–4846, 2019.

[38] Carole H Sudre, Wenqi Li, Tom Vercauteren, Sebastien Ourselin, and M Jorge Cardoso. Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In *Deep learning in medical image analysis and multimodal learning for clinical decision support*, pages 240–248. Springer, 2017.

[39] Sebastian Thrun. Probabilistic robotics. *Communications of the ACM*, 45(3):52–57, 2002.

[40] Michael Tomasello, Malinda Carpenter, Josep Call, Tanya Behne, and Henrike Moll. Understanding and sharing intentions: The origins of cultural cognition. *Behavioral and brain sciences*, 28(5):675–691, 2005.

[41] Pavan Vasishta, Dominique Vaufreydaz, and Anne Spalanzani. Natural vision based method for predicting pedestrian behaviour in urban environments. In *2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)*, pages 1–6. IEEE, 2017.

[42] Anirudh Vemula, Katharina Muelling, and Jean Oh. Social attention: Modeling attention in human crowds, 2018.

[43] Peter Whittle. *Hypothesis testing in time series analysis*, volume 4. Almqvist & Wiksells boktr., 1951.

[44] Paper with code. Eth/ucy trajectory prediction benchmark, `https : / / paperswithcode . com / sota / trajectory-prediction-on-ethucy`, 2020.

[45] Paper with code. Stanford drone trajectory prediction benchmark, `https://paperswithcode.com/sota/ trajectory-prediction-on-stanford-drone`, 2020.

[46] Takuma Yagi, Karttikeya Mangalam, Ryo Yonetani, and Yoichi Sato. Future person localization in first-person videos, 2018.

[47] Hang Zhao, Jiyang Gao, Tian Lan, Chen Sun, Benjamin Sapp, Balakrishnan Varadarajan, Yue Shen, Yi Shen, Yuning Chai, Cordelia Schmid, et al. Tnt: Target-driven trajectory prediction. *arXiv preprint arXiv:2008.08294*, 2020.