

# Adaptive Human Trajectory Prediction via Latent Corridors

Neerja Thakkar  
UC Berkeley

Karttikeya Mangalam  
UC Berkeley

Andrea Bajcsy  
Carnegie Mellon University

Jitendra Malik  
UC Berkeley

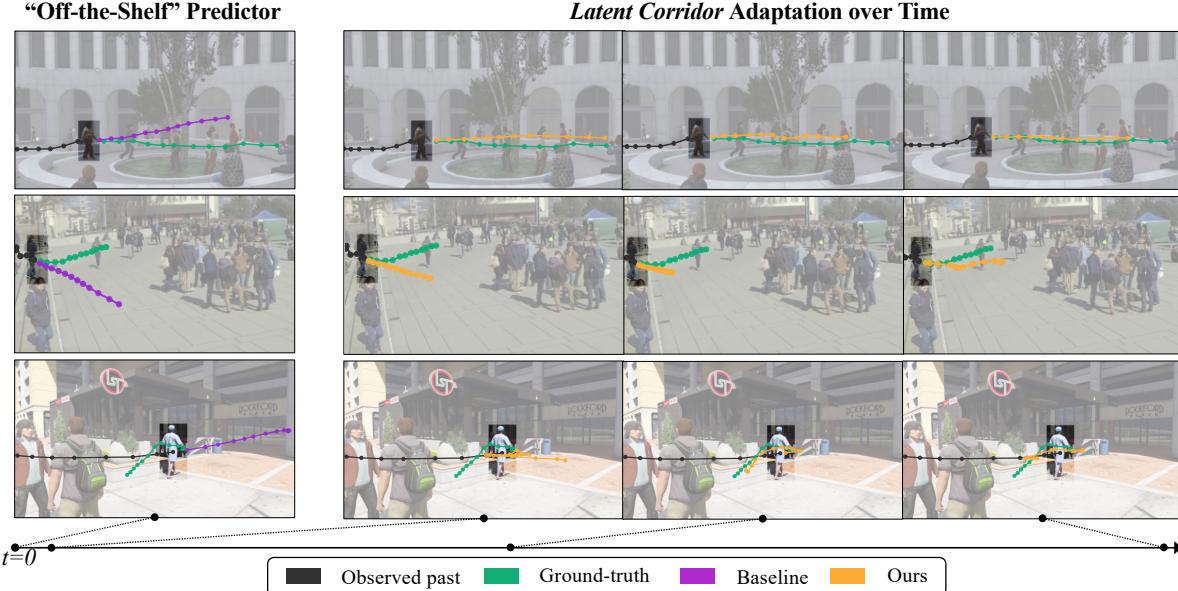


Figure 1. **Adaptive trajectory prediction.** (left) Given a history of human behavior (shown in black), the pre-trained predictor  $\mathcal{P}$  is unable to understand scene-specific behavior trends, like people entering a subterranean subway entrance (bottom row). (right) When adapting, the number of people and amount of time determine the total number of trajectories observed, and we denote this time-dependent quantity human-seconds. Here, the three columns correspond to our method trained for a very small (left), medium (middle), and large amount of human-seconds. Our adaptive *latent corridors* approach enables  $\mathcal{P}$  to quickly learn scene-specific trends, improving predictions with even small amounts of data, and closing the gap between the ground-truth (green) and predicted behavior (orange). For example, in the middle row, we see that  $\mathcal{P}$  predicts the person will move towards the camera, but as our method sees more human-seconds of data, it adapts to the trend that people in this plaza scene tend to avoid the center of the plaza and instead move diagonally across it.

## Abstract

*Human trajectory prediction is typically posed as a zero-shot generalization problem: a predictor is learnt on a dataset of human motion in training scenes, and then deployed on unseen test scenes. While this paradigm has yielded tremendous progress, it fundamentally assumes that trends in human behavior within the deployment scene are constant over time. As such, current prediction models are unable to adapt to scene-specific transient human behaviors, such as crowds temporarily gathering to see buskers, pedestrians hurrying through the rain and avoiding puddles, or a protest breaking out. We formalize the problem of scene-specific adaptive trajectory prediction and propose a new adaptation approach inspired by prompt tuning called latent corridors. By augmenting the input of any pre-trained human trajectory predictor with learnable image prompts,*

*the predictor can improve in the deployment scene by inferring trends from extremely small amounts of new data (e.g., 2 humans observed for 30 seconds). With less than 0.1% additional model parameters, we see up to 23.9% ADE improvement in MOTSynth simulated data and 16.4% ADE in MOT and Wildtrack real pedestrian data. Qualitatively, we observe that latent corridors imbue predictors with an awareness of scene geometry and scene-specific human behaviors that non-adaptive predictors struggle to capture.*

## 1. Introduction

Human motion prediction is a fundamental skill for intelligent systems to effectively navigate the world, assist end-users, and visually survey a scene. To date, learning-based human trajectory prediction has been extensively

studied, making huge strides in predicting multimodal future behavior [27], modeling multi-agent social interactions [2, 17, 35], and accounting for scene context [28, 49].

However, a fundamental assumption underlies the current trajectory prediction paradigm: predictors are assumed to be deployed within an *unchanging* context. Nevertheless, in the real world, the context—and thus patterns of human behavior—will inevitably change over time, as observed in [44]. For example, data from a security camera will capture a new subway entrance being built, and people descending and exiting in new patterns (bottom row, Fig. 1). Humans rushing to work during the day will largely ignore each other while nighttime party-goers will walk in cliques. Even if the trajectory predictor observes data captured from exactly the same location, the human behavior patterns can change in minutes or hours. When faced with such changes, the performance of existing trajectory predictors immediately degrades.

Here, we study how to efficiently *adapt* pre-trained trajectory predictors to observed human behavior within a deployment scene (right, Fig. 1). We focus on adaptation that can be done in a data-efficient way and can preserve any useful structure within the base predictor in the new context (e.g., even if a new subway entrance lets people descend into the ground, other people may still move towards the existing above-ground coffee shops). Our technical approach takes inspiration from the recently popularized paradigm of *prompt-tuning* in large language models [23], but instantiates this idea within the adaptive trajectory prediction problem. Specifically, we augment the input to a human trajectory predictor with a set of latent image prompts, one per each instance of a scene we wish to adapt to. We leave all or most of the predictor frozen, and tune each input prompt with the same trajectory prediction loss that is used to train the base predictor. We find that even extremely small amounts of new human trajectories (e.g., < 1 minute of human seconds, corresponding to 2 humans observed over 30 seconds) can leak sufficient information about how humans interact with the scene (e.g., new subway entrance) or with each other (e.g., nighttime social cliques) to learn this prompt and improve future trajectory prediction. We call our learnable prompt a *latent corridor*: a non-physical corridor that guides human behavior in this scene.

In our experiments, learning latent corridors for prediction adaptation enables up to 23.9% prediction accuracy improvement on simulated data from MOTSynth [12] compared to a non-adaptive predictor that takes as input the scene and has seen all the same data. On real data from MOT and WildTrack [7], we see a 16.4% improvement on non-adaptive predictors and a 11.2% improvement with adapting via our method as opposed to just fine-tuning. We find that our latent corridor adaptation can improve prediction performance even with a very small amount of data

(e.g. a couple of humans for a few seconds), and continually improves as more data (e.g., many humans for a few minutes) is observed over time. Qualitatively, we observe that learning latent corridors can inject 3D ground plane awareness into prediction performed from an arbitrary camera view angle, can help the predictor adapt to obstacles or occlusions in the scene, and learn scene-specific patterns of human behavior.

To summarize, our contributions are as follows:

1. We formalize the adaptive trajectory prediction problem.
2. We propose a novel method for adaptive trajectory prediction by learning lightweight latent corridor prompts in image space, outperforming non-adaptive trajectory predictors with less than a 0.1% parameter increase.
3. We demonstrate qualitative and quantitative improvements of up to 23.9% ADE improvement in MOTSynth simulated data and 16.4% ADE in MOT and Wildtrack real pedestrian data.

## 2. Related Work

**Human Trajectory Prediction** is a well-studied problem with a long and rich history [33]. Prediction methods started from simple models such as social forces [5, 15], and later added multimodality through Gaussian processes [38]. Recently, modern learning-based predictors have made significant progress in incorporating social interactions between humans as modelled via RNNs [16, 35, 40], GANs [13], or conditional VAEs [27], all while generating multimodal future trajectories. Another line of works incorporate the scene context in predictions, getting neural network features from a scene map [20, 22, 34, 39, 46], and more recently using transformers [30, 36]. A few works have incorporated the scene context by projecting the trajectories into heatmaps so that the network can reason in image space [6, 28], allowing for longer-term predictions. Our work builds upon these advances by adapting an RGB scene-aware pre-trained trajectory predictor [28] over time.

**Adapting Human Trajectory Predictors.** In the past few years, there has been a growing interest in lightweight ways to modify pre-trained predictors to new data. The key differences between these works lie in 1) *what* they are adapting to, and 2) *how* they adapt. Several works focus on cross-domain transfer, wherein a predictor trained for trajectory prediction in domain A (e.g., New York) is adapted to work in domain B (e.g., London). These works leverage architectures that partition generic trajectory prediction from domain-specific features [42, 45], leverage adaptive meta-learning via Kalman filtering [18], or simply finetune the prediction network on new data [25]. Other works focus on online adaptation over time, for example by adapting to different agent dynamics using recursive least squares [1, 10]. [47] carries out continual learning on 2D birds-eye view scenes for pedestrian trajectory prediction, using a memory

bank to store embedded pairs of trajectory pasts and futures. In this work, we focus on adaptation over time within a specific deployment domain, but propose a novel prompting-based adaptation method which learns scene-specific information in a scalable way from in-the-wild data.

**Prompt Tuning in Language & Vision.** With the rise of large pre-trained models, efficient adaptation for downstream tasks has gained increasing interest. In the large language model domain, prompt tuning—wherein a few tunable tokens (i.e., prompt) are prepended to input text and tuned per downstream task—has been shown to be remarkably effective [23, 24, 26]. In vision models, image prompts [41] have been introduced to adapt a vision model to new tasks [3] or instruct an inpainting model to carry out various computer vision tasks [4]. Recently, there have also been attempts at visual prompt tuning, showing that large vision transformers can benefit from utilizing tuned prompts for continual learning or transfer learning [19, 31, 37, 43]. While our task is different, we were inspired by the prompting scheme of [19] that learned classification on a new dataset by introducing a trainable input space prompt and tuning the prompt along with the transformer predictor head. We use visual prompts as a lightweight way of adapting a trajectory predictor. To our knowledge, we are the first work to apply prompt tuning to human trajectory prediction.

### 3. Problem Formulation

Here, we present a formalization of our proposed adaptive trajectory prediction problem. We seek to generate accurate future trajectories of  $N$  agents in a scene. At any time  $\tau$ , let  $\mathbf{o}_\tau := o_{\tau-H:\tau} \in \mathcal{O}$  be the  $H$ -step history of observations input into the predictor,  $\mathbf{y}_\tau := y_{\tau+1:\tau+T} \in \mathcal{Y}$  be the ground-truth  $T$ -step future trajectory, and  $\hat{\mathbf{y}}_\tau := \hat{y}_{\tau+1:\tau+T} \in \hat{\mathcal{Y}}$  be the prediction outputs. We assume access to a base predictor  $\mathcal{P}$  pre-trained on a human motion dataset consisting of past and ground-truth future trajectories,  $\mathcal{D} = \{\mathbf{o}_i, \mathbf{y}_i\}_{i=1}^K$  to minimize a loss function on trajectories,  $\ell$  (e.g., mean squared error).

**Adaptive Trajectory Prediction (ATP).** We aim to adapt  $\mathcal{P}$  to a deployment scene by observing new human trajectory data over a time interval. Let the dataset of *new* human trajectory data be  $\mathcal{D}'_{0:T}$  collected over the  $T$  step time interval. The deployment scenes being adapted to in  $\mathcal{D}'$  can be real or synthetic, can be of a new physical environment or same as the pre-training dataset, and can be obtained by a new or previously seen camera configuration. However, we assume that human pedestrians are present in the scene and they are navigating an outdoor environment. We adapt to scenes captured for a similar amount of time, between  $T = 45$  seconds and 5 minutes. We aim to adapt to scene-specific characteristics and fluctuating events beyond those that can be exploited by conditioning only on a segmenta-

tion map.

In the adaptive trajectory prediction problem, a subset of the deployment dataset  $\mathcal{D}'_{0:t}$ ,  $t < T$  is observed. The goal is to create an adapted model  $\mathcal{A}[\mathcal{P}]$  such that the prediction error decreases over the future observed deployment data  $\mathcal{D}'_{t:T}$ :

$$\ell(\mathcal{A}[\mathcal{P}(\mathbf{o})], \mathbf{y}) < \ell(\mathcal{P}(\mathbf{o}), \mathbf{y}), \quad (\mathbf{o}, \mathbf{y}) \in \mathcal{D}'_{t:T}. \quad (1)$$

Intuitively, as the value of  $t$  increases and the adapted predictor sees more data, the performance of  $\mathcal{A}[\mathcal{P}]$  should improve over the pre-trained predictor  $\mathcal{P}$ . At deployment,  $\mathcal{A}[\mathcal{P}]$  should have learned scene context-specific characteristics of human behavior.

## 4. Adaptation via Learned Latent Corridors

To adapt our trajectory predictor efficiently, we take inspiration from language-based prompt-tuning [23] and instantiate it within our adaptive trajectory prediction problem. Specifically, we augment a frozen base predictor with a learnable latent prompt we call a *latent corridor*, a non-physical corridor that informs the predictor of human behavior patterns in a specific deployment scene. In this section, we detail our predictor architecture, prompt representation and training, and investigate a suite of prompting configurations that are compatible with latent corridors.

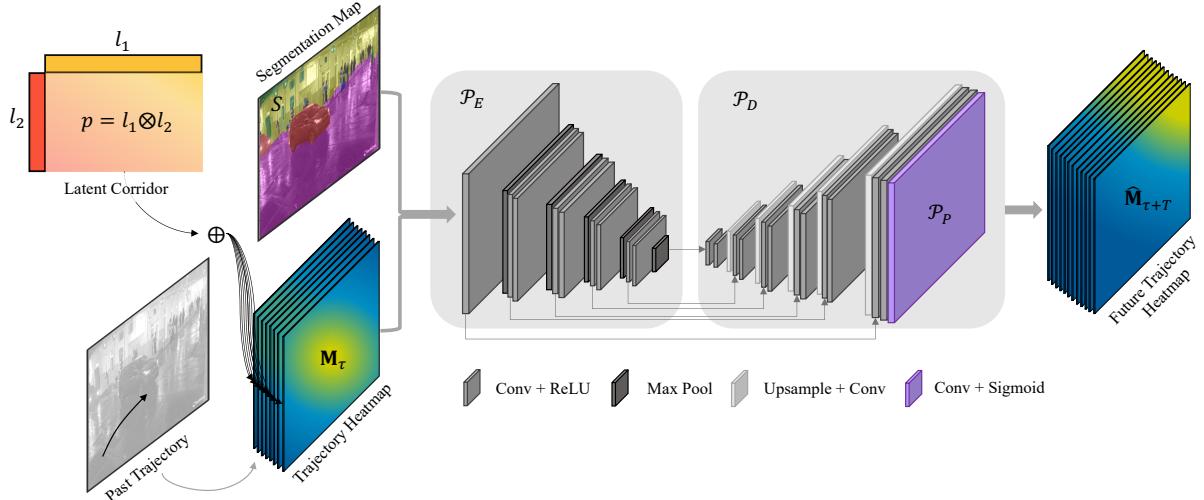
### 4.1. Base Trajectory Predictor Architecture

Our base predictor model,  $\mathcal{P}$ , consists of three modules: an encoder  $\mathcal{P}_E$ , decoder  $\mathcal{P}_D$ , and predictor head  $\mathcal{P}_P$  (see Figure 2). We use the YNet encoder for  $\mathcal{P}_E$ , and trajectory decoder architecture for  $\mathcal{P}_D$  and  $\mathcal{P}_P$  [28].

**Scene and trajectory representation.** To capture the scene, the RGB image  $I$  of the first video frame is processed with the Mask2Former semantic segmentation model [8, 9]. The final semantic segmentation map,  $S$ , is obtained by downsampling the segmentation classes to  $C = 12$  meaningful classes for pedestrians in outdoor environments. We follow [28], and convert the history of observed agent positions  $x_{\tau-H:\tau} \in \mathbb{R}^{H \times 2}$  into a trajectory of heatmaps,  $\mathbf{M}_{\tau-H:\tau}$ , each of the same spatial size as  $I$ , for a total size of  $H \times h \times w$ . Heatmaps are concatenated with the semantic segmentation map  $S$ ; the final input into  $\mathcal{P}$  is  $\mathbf{o}_\tau := ([\mathbf{M}_{\tau-H:\tau}, S])$  of size  $(C + H) \times h \times w$ . The predictor also outputs a trajectory of heatmaps,  $\hat{\mathbf{M}}_{\tau+1:\tau+T}$ . The corresponding ground truth  $T$ -step future agent trajectory,  $\mathbf{y}_\tau \in \mathbb{R}^{T \times 2}$ , is also converted into a trajectory of heatmaps,  $\mathbf{M}_{\tau+1:\tau+T}$ , for loss computation.

**Loss function.** We train the predictor with a binary cross entropy loss on the trajectory heatmaps,

$$\ell := \sum_{\tau} \text{BCE}(\mathcal{P}([\mathbf{M}_{\tau-H:\tau}, S]), \mathbf{M}_{\tau+1:\tau+T}). \quad (2)$$



**Figure 2. Adapting a predictor  $\mathcal{P}$  with latent corridors.**  $\mathcal{P}_E$ ,  $\mathcal{P}_D$  and  $\mathcal{P}_P$  are pre-trained on the task of human trajectory prediction, taking as input trajectory heatmaps  $M_{\tau-H:\tau}$  and segmentation  $S$ , and outputting predicted trajectory heatmaps  $\hat{M}_{\tau+1:\tau+T}$ . We augment  $\mathcal{P}$  with a per-scene *latent corridor*  $p$  which is summed element-wise to the input trajectory heatmaps. The latent corridors are trained with  $\mathcal{P}_E$  and  $\mathcal{P}_D$  frozen. The predictor head  $\mathcal{P}_P$  can be frozen, tuned on a single deployment scene, or tuned jointly across multiple scenes.

For evaluation, a `soft argmax` operation is used to sample 2D points from the predicted heatmap  $\hat{M}_{\tau+1:\tau+T}$  as in [28]. This yields predictions in x-y pixel space,  $\hat{\mathbf{y}}_\tau = \hat{y}_{\tau+1:\tau+T}$ , for computing displacement error metrics in Section 6 with respect to the ground-truth future positions,  $\mathbf{y}_\tau$ .

#### 4.2. Representing & Learning Latent Corridors

Our key idea for adapting trajectory predictors is to augment the frozen base model  $\mathcal{P}$  with a set of trainable prompts, called *latent corridors*, which learn new trends in how humans interact with the scene or with each other. For effective adaptation, we seek two key properties for our latent corridors: parameter-efficient (i.e., we want to minimize the number of parameters that are tuned from new deployment data) and spatially scene-grounded (i.e., the latent should have pixel-wise alignment with the scene). We first outline our prompting approach, and then discuss a compact but spatially-grounded representation of the prompt that is amenable to parameter-efficient learning.

**Latent corridor prompt.** For each of the  $K$  deployment scenes, we introduce a unique trainable prompt,  $p_k \in \mathbb{R}^{h \times w}$ , of the same spatial size as the image  $I$  that is input into the predictor (left, Fig. 2). For a network that has adapted to  $K$  scenes, we have a set  $K$  of prompts  $\mathcal{P} := \{p_0, p_1, \dots, p_K\}$ . Thus, our adaptation rule is

$$\mathcal{A} := \mathbf{M}_t \oplus p \quad \forall t \in \{\tau - H, \dots, \tau\}. \quad (3)$$

The prompt is summed element-wise to each of the input heatmaps corresponding to the observed trajectories; let  $\mathbf{M} = \mathbf{M} \oplus p$  denote this for any original  $\mathbf{M}$ . See Sec. 8 for an ablation on prompt location. The adapted predictor takes as input  $\mathcal{P}([\mathbf{M}_{\tau-H:\tau}, S])$ .

**A compact but spatially-grounded representation.** In Equation (3), the prompt is assumed to be the same size as the image,  $h \times w$ , which is nice for spatial alignment between the prompt and scene. However, in our experiments, where all images are of size  $h = 288$  and  $w = 480$ , this naive image-based representation requires learning an additional 138K parameters. Considering that our base predictor  $\mathcal{P}$  has  $\sim 900K$  trainable parameters, this prompt increases the model parameter size by over 15% and makes it infeasible to learn meaningful latent corridors from small amounts of new human trajectories (see Sec. 8). Instead, we propose to represent the prompt as a *low-rank representation* with rank 1. We initialize our latent corridor as a vector of dimension  $h + w$  using Kaiming initialization [14], and parameterize the full prompt as the outer product of the  $h$  and  $w$  dimensional vectors. This lightweight representation preserves the spatial relationship between the prompt, scene, and trajectory heatmaps and is significantly more compact: it increases the model parameter size by less than 0.1%.

**Training.** We learn the prompt  $p$  while the predictor encoder  $\mathcal{P}_E$  and decoder  $\mathcal{P}_D$  are frozen. The predictor head  $\mathcal{P}_P$  is optionally tuned and the latent corridors can be trained individually on one scene at a time, or simultaneously amongst many scenes (see Section 4.3). We use the same trajectory loss  $\ell$  from Equation (2) that the base predictor was pre-trained with.

#### 4.3. Prompting configurations

One of the strengths of our latent corridors approach is that it is compatible with a suite of deployment desiderata that can inform how the prompts are incorporated into the predictor. We identify and study three settings. In each setting,

$K$  unique prompts are trained on  $K$  deployment scenes individually, but the treatment of the predictor head differs.

1. **Latent corridor adaptation (LC):** The simplest use of latent corridors is to keep the entire base predictor frozen, including the original predictor head  $\mathcal{P}_P$ . This design is the fastest to train since it has the smallest number of trainable parameters, so it is desirable when rapid adaptation to short-term transient events occurs. It also enables easy recovery of the base predictor model and its original performance by simply not inputting a prompt.
2. **Multi-scene finetuning (LC + Joint FT):** In this setting, each of the  $K$  deployment scenes has a unique latent, but one single predictor head  $\mathcal{P}_P$  is jointly tuned across the  $K$  scenes. The latents retain unique scene information, but  $\mathcal{P}$  is better adapted to in conjunction with the per-scene latents. Similarly to the LC approach above, this is a more compact configuration since one prediction model is used for all scenes, so it is desirable if deployment hardware has limited space. It is also faster to train than per-scene finetuning especially as  $K$  grows, since multiple predictor heads do not have to be tuned.
3. **Per-scene finetuning (LC + Per-Scene FT):** If one seeks to maximize performance within a *specific* deployment scene, one can *jointly* finetune  $K$  predictor heads  $\mathcal{P}_{P_k}$  with  $K$  latent corridors per scene. While this results in the need for a unique predictor for each deployment scene, we find empirically that this method achieves best in-scene adaptation performance.

## 5. Experimental Setup

We study our approach on synthetic and real datasets. Here, we detail these datasets, describe how we evaluate adaptation quality over time, and outline our trajectory predictor baselines.

### 5.1. Datasets

**MOTSynth.** We start in simulation with MOTSynth [12], a synthetic pedestrian detection and tracking dataset of over 700 scenes with varying camera viewpoints and outdoor environments, with each video 90 seconds long. Pedestrians carry out simple actions such as walking, standing, or running, and follow manually pre-planned flows as well as a collision avoidance algorithm. MOTSynth has over 17 hours of video; we select a subset of approximately 11 hours of video corresponding to the 437 scenes with a static camera. Due to the large size of the dataset and perfect ground truth detections, MOTSynth was a good starting point.

**MOT & WildTrack.** We also evaluated our approach on the real-world pedestrian datasets MOT and WildTrack. WildTrack [7] consists of 7 static camera viewpoint videos of pedestrians walking through a plaza. We randomly selected 3 videos and took the first 5 minutes of each video.

MOT, the multiple object tracking benchmark, consists of several video datasets of pedestrians navigating outdoor environments. We combined the datasets MOT15 [21], MOT16 [29], and MOT20 [11], and removed videos with dynamic cameras, pedestrian density of less than 10, and length less than 45s, leaving 4 scenes: MOT16-03, MOT20-02, AVG-TownCentre, and PETS09-S2L2.

**EarthCam.** We evaluate our approach in-the-wild on data from <https://www.earthcam.com/>, which has livestream webcam data from around the world available, and from which we scrape 5 minute segments of data.

### 5.2. Train-Test Split Over Time in Human Seconds

Given a video of pedestrians in motion, we prepare our dataset to adapt to the scene as follows. We use provided annotations or run ByteTrack [48] to get tracklets of all  $N$  detectable identities in the video. Tracklets are downsampled to 2 timesteps per second, then windowed by 20 timesteps with an observed length  $H = 8$  and future length  $T = 12$ . Then, we sort the  $N$  identities chronologically by their first appearance in the scene. We select the first 80% of identities that appear in the scene for our training dataset, and hold out the last 20% for the testing set. When we conduct experiments over time, the testing set is always the same 20% of identities, whereas the training set consists of the first  $m\%$  of identities,  $m \leq 80\%$ . We report adaptation time in person-seconds of observation, where 1 person for 1 second is 1 person-second, so for instance, 30 people observed for 10 seconds each would be 300 person-seconds.

### 5.3. Base Predictor Pre-Training Dataset

We first pre-train our base predictor  $\mathcal{P}$  described in Sec 4.1 on dataset  $\mathcal{D}$  which consists of simulated human pedestrian agents moving around various outdoor environments, captured from a single static RGB camera with any point of view (not necessarily birds eye view). This dataset comes from 90-second MOTSynth videos with a static camera. We pre-train our predictor on a train-test split over time of  $\mathcal{D}$  as described in the previous section.

### 5.4. Baselines

We evaluate our method with the following baselines.

1. *Constant velocity:* As in [5].
  2. *Learned trajectory:* Simplified PECNet [27], which is a learned predictor with position history but no scene.
  3. *Scene-aware:* Simplified YNet [28], which is a learned predictor with position and scene input.
  4. *ATP Finetune:* Adaptation baseline which finetunes scene-aware predictor head  $\mathcal{P}_P$  without latent corridors.
- We implement the learned trajectory baseline by taking the encoder, decoder and trajectory loss from PECNet [27]. This removes multimodality for simplicity, but uses the an



Figure 3. **Qualitative results** on MotSynth (top; synthetic) and MOT and WildTrack (bottom; real). These examples show scenarios where our LC + per-scene finetune ATP method (orange) outperforms the scene-aware baseline (purple). In several MOTSynth examples, the baseline predicts the pedestrian floats into the air (top row), while our method has gained awareness of where the 3D ground plane lies in the 2D image. We also note that patterns of behaviour such as walking on the sidewalk instead of into the road (second row left) and walking up the traversable portion of stairs (second row right) is captured. On real data, we observe similar awareness of the ground plane and obstacles, as well as a better understanding of nuanced human behavior patterns such as crossing diagonally across a plaza.

architecture proven to be successful on the trajectory prediction problem. Similarly, YNet’s [28] encoder and trajectory decoder are used for the scene-aware baseline.

### 5.5. Metrics

Our experiments are evaluated using the Average Displacement Error (ADE) [32] and Final Displacement Error (FDE) [2] metrics. ADE is the average  $l_2$  error between the entire predicted and ground truth trajectory, while FDE is the  $l_2$  error between just the final point of the predicted trajectory and the final point of the ground truth trajectory.

## 6. Results

Here, we detail quantitative and qualitative results on MOT-Synth, MOT, and WildTrack, and EarthCam datasets.

### 6.1. MOTSynth Results

We first evaluate if latent corridors enable predictor improvement when the deployment scene is in the pre-training dataset,  $\mathcal{D}$ . We train the baselines and two variants of our ATP method (LC and LC + Joint Finetune, described in Sec. 4.3) on the MOTSynth pre-training dataset  $\mathcal{D}$ . All models see the same training data, and both the scene-aware baseline and our adaptive method are conditioned on each scene semantic segmentation map. Results are in the first column of Table 1. The learned trajectory baseline significantly improves over constant velocity, and adding in scene awareness results in a 7.6% performance gain on ADE and 3.5% gain on FDE. Our latent corridor approach results in a 5.7% and 9.9% improvement over the scene-aware baseline for ADE, without and with joint finetuning of the predictor

Method	MOTSynth		MOT		Wildtrack		NoLA	
	ADE	FDE	ADE	FDE	ADE	FDE	ADE	FDE
Constant velocity	78.5	160.4	47.7	99.3	44.9	90.1	42.4	82.8
Learned trajectory (PECNet-Ours)	51.2	100.0	49.7	103.4	43.8	83.1	36.9	65.9
Scene-aware (YNet-Ours)	47.3	96.5	44.2	91.0	33.6	67.7	34.0	65.1
ATP (LC)	44.6	90.2	43.0	89.2	32.9	65.9	33.5	63.1
ATP (LC + Joint Finetune)	<b>42.6</b>	<b>85.3</b>	-	-	-	-	-	-
ATP (LC + Per-Scene Finetune)	-	-	<b>37.4</b>	<b>74.3</b>	<b>27.4</b>	<b>54.7</b>	<b>26.1</b>	<b>47.4</b>

Table 1. **Our method vs baselines on real-world datasets.** The average ADE and FDE are in pixel space of (left-to-right) 437 MOTSynth scenes, 4 MOT scenes, 3 WildTrack scenes, and 2 NoLA webcam scenes. Across these synthetic and real-world deployment scenes, our latent corridor adaptation method is comparable to finetuning the predictor head and outperforms non-adaptive baselines. Our method with latent corridors and per-scene finetuning consistently outperforms all other approaches.

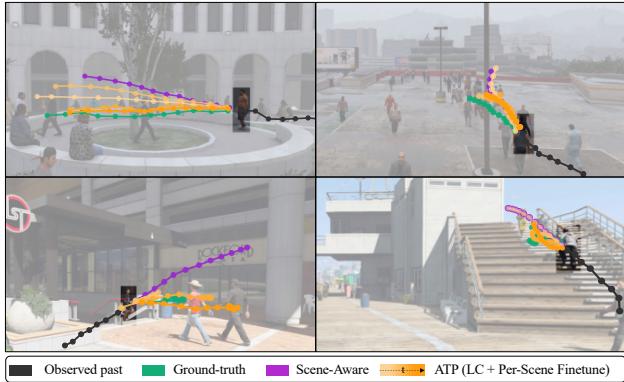


Figure 4. **Qualitative results over time.** Our method’s predictions that saw a short number of human seconds in training (2%) are shown in light orange, with a human seconds training time of up to 80%. With the latent corridor trained on a tiny amount of data, the predictions sometimes significantly improve, but other times are close to the baseline. When more human seconds of data have been seen, the adaptation results consistently improve.

layer respectively, and a 6.5% and 11.6% improvement on FDE. This indicates that our latent corridor approach is able to more effectively learn scene-conditioned information that is useful for the trajectory prediction task.

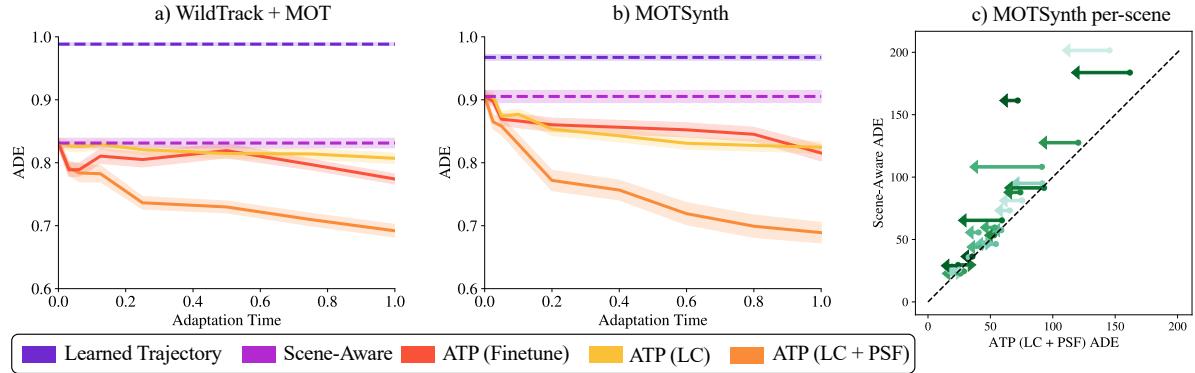
Next, for a random subset of 25 of the MOTSynth scenes, we additionally train latent corridors with per-scene finetuning (LC + Per-Scene FT, from Sec 4.3). The datasets  $\mathcal{D}'_{0:t}$  correspond to increasing human-second lengths. The test set trajectories are the last 20% of agents in the deployment scene, and the training sets consist of the first 2, 4, 8, 16, 32, 48, 64 and 80% of agents to enter the scene. Results normalized per-scene and averaged are shown in Fig. 5b. Using only latent corridors has comparable results to only finetuning the last layer for each deployment scene, while latent corridors with per-scene finetuning yields significant performance gains. We visualize the ADE results per-scene with our LC + finetuning method trained with 8% and then trained with 80% of the data in Fig. 5c. We see that while the effectiveness of the adaptation for varying

time horizons differs for each deployment scene, there are many scenes where adaptation yields significant gains, and some where the improvement on ADE error is up to 63.4%. We hypothesize that the scenes where our method does not help exhibit behaviors and environment geometries that the prior  $\mathcal{P}$  is sufficient for.

**Qualitative results.** Visualizations from the per-scene models can be seen in Fig. 3 and Fig. 4. We observe many examples where the scene-aware baseline seems to lack awareness of the ground plane. There are multiple instances where the scene-aware baseline predicts a pedestrian floats into the air (Fig. 3 top row and Fig. 4), whereas our predictions lie closely in line with the ground truth trajectory in terms of distance from the ground plane. This holds even for non-planar ground planes: When a pedestrian walks up stairs, our model seems to better understand their 3D structure, for example in the bottom right of Fig. 4. Additionally, our approach captures trends of behaviour in scenes such as walking on the sidewalk instead of into the road, walking around the rotunda, and walking up the section of smaller, easily traversable steps (see Fig. 3 second row, left to right).

## 6.2. MOT and WildTrack Results

We next investigate if latent corridors help predictor adaptation when the deployment scene is *outside* of the pre-training dataset, and if latent corridors outperform ATP via direct finetuning on deployment data. Specifically, we use the base predictor  $\mathcal{P}$  which only saw MOTSynth data and then directly do sim-to-real adaptation via latent prompting, finetuning, or a combination thereof on seven *real human pedestrian scenes* from MOT and Wildtrack. Plots of the ADE over time averaged over these scenes is shown in Figure 5a. Regardless of the ATP setting, adaptation over time results in a consistent error reduction compared to the baselines. While only finetuning seems to be slightly more effective on real-world scenes than prompting alone, our latent corridors + per-scene finetuning approach is significantly more effective (11.2%) than finetuning alone. Quantitative results for these experiments can be seen in Table 1.



**Figure 5. Learning latent corridors over time.** (a and b): The x-axis represents adaptation time in person-seconds, or the amount trajectories use to train  $\mathcal{A}(\mathcal{P})$ , and the y-axis represents the ADE. Results are normalized per-scene and averaged over models trained on 25 MOTSynth scenes (a) and 7 from MOT and WildTrack (b), with shaded area  $\sigma/10$ . On all data, even with a short adaptation time, our methods improve on the baselines, and as adaptation time increases, performance improves. Latent corridors + per-scene finetuning has the best performance. c) Comparison to baseline over many MOTSynth scenes for models trained with 8% (point) and 80% (arrowhead) h-s datasets. Each arrow represents one scene, with the ADE using our ATP method plotted against the scene-aware baseline ADE. For some deployment scenes, the scene-aware baseline suffices, whereas other scenes see much more significant benefits from our method.

**Qualitative results.** Visualizations of baselines and our LC + PSF approach on MOT and WildTrack are in Fig. 3. Similar trends from experiments on synthetic data carry over. Our approach enables the predictor to better ground future behavior in the scene geometry: for example, pedestrians are no longer predicted to float upwards (third row, left and middle), and future behavior is guided by a finer-grained awareness of obstacles (third row right, bottom row left and middle). We also observe our adapted predictor learning trends in human behavior: in the WildTrack plaza environment shown in Fig. 1 middle row and Fig. 3 bottom right, pedestrians tend to avoid the middle of the plaza and instead cross it diagonally. Our method learns this subtle, scene-specific behavior pattern and thus predicts more accurately.

### 6.3. EarthCam NoLA Results

Finally, we evaluate the performance of our approach on truly “in-the-wild” data by scraping two 5-minute videos of pedestrian data captured from Bourbon Street in New Orleans: one video during daytime activity (top row, Fig. 6) and one video during nighttime activity (bottom row, Fig. 6). Additional results on different scenes of webcam data are in Sec 10. Quantitatively, the results on this in-the-wild data align with our results on real and synthetic data: our ATP model with latent corridors and per-scene finetuning outperforms the non-adaptive baselines and pure finetuning. Qualitatively, our method adapts to daytime patterns where pedestrians navigate around carts and don’t frequently walk diagonally through streets because of through-traffic. When the adaptive predictor observes nighttime behavior, it no longer has to respect these daytime patterns and learns to predict pedestrians as crossing the diagonally. Similarly to the prior datasets we explored, our ATP model



**Figure 6. Results on webcam data of New Orleans.** We train our method on one intersection during the daytime, when pedestrians must obey vehicles moving through the street, and at nighttime, when pedestrians take over. During the day, the model learns to account for humans walking through crosswalks instead of diagonally into the intersection, and to account for transient push-carts in the area. We also see awareness of the 3D ground plane.

again learns to ground pedestrian future behavior in the 3D ground plane.

## 7. Conclusion

In this work, we formalize and study the problem of adaptive trajectory prediction: the ability of human predictors to adapt to changing deployment conditions and environments. To this end, we proposed a lightweight adaptation approach grounded in image-based prompt tuning called latent corridors. Through extensive experiments on both simulated and real world pedestrian datasets, we observed that latent corridors enable a data-efficient way to adapt pre-trained predictors to new scene-specific human behavior.

## References

- [1] Abulkemu Abuduweili, Siyan Li, and Changliu Liu. Adaptable human intention and trajectory prediction for human-robot collaboration. *arXiv preprint arXiv:1909.05089*, 2019. 2
- [2] 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, 6
- [3] Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. Exploring visual prompts for adapting large-scale models. *arXiv preprint arXiv:2203.17274*, 2022. 3
- [4] Amir Bar, Yossi Gandelsman, Trevor Darrell, Amir Globerson, and Alexei A. Efros. Visual prompting via image inpainting. *arXiv preprint arXiv:2209.00647*, 2022. 3
- [5] R.A. Best and J.P. Norton. A new model and efficient tracker for a target with curvilinear motion. *IEEE Transactions on Aerospace and Electronic Systems*, 33(3):1030–1037, 1997. 2, 5
- [6] Zhe Cao, Hang Gao, Karttikeya Mangalam, Qi-Zhi Cai, Minh Vo, and Jitendra Malik. Long-term human motion prediction with scene context. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 387–404. Springer, 2020. 2
- [7] Tatjana Chavdarova, Pierre Baqué, Stéphane Bouquet, Andrii Maksai, Cijo Jose, Timur Bagautdinov, Louis Lettry, Pascal Fua, Luc Van Gool, and François Fleuret. Wildtrack: A multi-camera hd dataset for dense unscripted pedestrian detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5030–5039, 2018. 2, 5
- [8] Bowen Cheng, Alexander G. Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. 2021. 3
- [9] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. 2022. 3
- [10] Yujiao Cheng, Weiye Zhao, Changliu Liu, and Masayoshi Tomizuka. Human motion prediction using semi-adaptable neural networks. In *2019 American Control Conference (ACC)*, pages 4884–4890. IEEE, 2019. 2
- [11] Patrick Dendorfer, Hamid Rezatofighi, Anton Milan, Javen Shi, Daniel Cremers, Ian Reid, Stefan Roth, Konrad Schindler, and Laura Leal-Taixé. Mot20: A benchmark for multi object tracking in crowded scenes. *arXiv preprint arXiv:2003.09003*, 2020. 5
- [12] Matteo Fabbri, Guillem Brasó, Gianluca Maugeri, Orcun Cetintas, Riccardo Gasparini, Aljoša Ošep, Simone Calderara, Laura Leal-Taixé, and Rita Cucchiara. Motsynth: How can synthetic data help pedestrian detection and tracking? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10849–10859, 2021. 2, 5
- [13] 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. 2
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015. 4
- [15] Dirk Helbing and Peter Molnar. Social force model for pedestrian dynamics. *Physical review E*, 51(5):4282, 1995. 2
- [16] Boris Ivanovic and Marco Pavone. The trajectron: Probabilistic multi-agent trajectory modeling with dynamic spatiotemporal graphs. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2375–2384, 2019. 2
- [17] Boris Ivanovic, Edward Schmerling, Karen Leung, and Marco Pavone. Generative modeling of multimodal multi-human behavior. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3088–3095. IEEE, 2018. 2
- [18] Boris Ivanovic, James Harrison, and Marco Pavone. Expanding the deployment envelope of behavior prediction via adaptive meta-learning. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7786–7793. IEEE, 2023. 2
- [19] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *European Conference on Computer Vision*, pages 709–727. Springer, 2022. 3
- [20] Kris M Kitani, Brian D Ziebart, James Andrew Bagnell, and Martial Hebert. Activity forecasting. In *Computer Vision–ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7–13, 2012, Proceedings, Part IV 12*, pages 201–214. Springer, 2012. 2
- [21] Laura Leal-Taixé, Anton Milan, Ian Reid, Stefan Roth, and Konrad Schindler. Motchallenge 2015: Towards a benchmark for multi-target tracking. *arXiv preprint arXiv:1504.01942*, 2015. 5
- [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. 2
- [23] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021. 2, 3
- [24] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021. 3
- [25] Yiheng Li, Seth Z Zhao, Chenfeng Xu, Chen Tang, Chenran Li, Mingyu Ding, Masayoshi Tomizuka, and Wei Zhan. Pre-training on synthetic driving data for trajectory prediction. *arXiv preprint arXiv:2309.10121*, 2023. 2
- [26] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroyuki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023. 3

- [27] 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. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 759–776. Springer, 2020. [2](#) [5](#)
- [28] Karttikeya Mangalam, Yang An, Harshayu Girase, and Jitendra Malik. From goals, waypoints & paths to long term human trajectory forecasting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15233–15242, 2021. [2](#) [3](#) [4](#) [5](#) [6](#)
- [29] Anton Milan, Laura Leal-Taixé, Ian Reid, Stefan Roth, and Konrad Schindler. Mot16: A benchmark for multi-object tracking. *arXiv preprint arXiv:1603.00831*, 2016. [5](#)
- [30] Jiquan Ngiam, Benjamin Caine, Vijay Vasudevan, Zhengdong Zhang, Hao-Tien Lewis Chiang, Jeffrey Ling, Rebecca Roelofs, Alex Bewley, Chenxi Liu, Ashish Venugopal, et al. Scene transformer: A unified architecture for predicting multiple agent trajectories. *arXiv preprint arXiv:2106.08417*, 2021. [2](#)
- [31] Xing Nie, Bolin Ni, Jianlong Chang, Gaofeng Meng, Chunlei Huo, Shiming Xiang, and Qi Tian. Pro-tuning: Unified prompt tuning for vision tasks. *IEEE Transactions on Circuits and Systems for Video Technology*, 2023. [3](#)
- [32] Stefano Pellegrini, Andreas Ess, Konrad Schindler, and Luc Van Gool. You’ll never walk alone: Modeling social behavior for multi-target tracking. In *2009 IEEE 12th international conference on computer vision*, pages 261–268. IEEE, 2009. [6](#)
- [33] Andrey Rudenko, Luigi Palmieri, Michael Herman, Kris M Kitani, Dariu M Gavrila, and Kai O Arras. Human motion trajectory prediction: A survey. *The International Journal of Robotics Research*, 39(8):895–935, 2020. [2](#)
- [34] 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/CVF conference on computer vision and pattern recognition*, pages 1349–1358, 2019. [2](#)
- [35] Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and Marco Pavone. Trajectron++: Dynamically-feasible trajectory forecasting with heterogeneous data. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16*, pages 683–700. Springer, 2020. [2](#)
- [36] Tim Salzmann, Lewis Chiang, Markus Ryll, Dorsa Sadigh, Carolina Parada, and Alex Bewley. Robots that can see: Leveraging human pose for trajectory prediction. *IEEE Robotics and Automation Letters*, 8(11):7090–7097, 2023. [2](#)
- [37] Mark Sandler, Andrey Zhmoginov, Max Vladymyrov, and Andrew Jackson. Fine-tuning image transformers using learnable memory. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12155–12164, 2022. [3](#)
- [38] Meng Keat Christopher Tay and Christian Laugier. Modelling smooth paths using gaussian processes. In *Field and Service Robotics: Results of the 6th International Conference*, pages 381–390. Springer, 2008. [2](#)
- [39] Daksh Varshneya and G Srinivasaraghavan. Human trajectory prediction using spatially aware deep attention models. *arXiv preprint arXiv:1705.09436*, 2017. [2](#)
- [40] Anirudh Vemula, Katharina Muelling, and Jean Oh. Social attention: Modeling attention in human crowds. In *2018 IEEE international Conference on Robotics and Automation (ICRA)*, pages 4601–4607. IEEE, 2018. [2](#)
- [41] Jiaqi Wang, Zhengliang Liu, Lin Zhao, Zihao Wu, Chong Ma, Sigang Yu, Haixing Dai, Qiushi Yang, Yiheng Liu, Songyao Zhang, et al. Review of large vision models and visual prompt engineering. *arXiv preprint arXiv:2307.00855*, 2023. [3](#)
- [42] Letian Wang, Yeping Hu, Liting Sun, Wei Zhan, Masayoshi Tomizuka, and Changliu Liu. Transferable and adaptable driving behavior prediction. *arXiv preprint arXiv:2202.05140*, 2022. [2](#)
- [43] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 139–149, 2022. [3](#)
- [44] William H. Whyte. *The social life of small urban spaces*. Washington, D.C.:Conservation Foundation, 1980. [2](#)
- [45] Yi Xu, Lichen Wang, Yizhou Wang, and Yun Fu. Adaptive trajectory prediction via transferable gnn. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6520–6531, 2022. [2](#)
- [46] Hao Xue, Du Q. Huynh, and Mark Reynolds. Ss-lstm: A hierarchical lstm model for pedestrian trajectory prediction. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1186–1194, 2018. [2](#)
- [47] Biao Yang, Fucheng Fan, Rongrong Ni, Jie Li, Loochu Kiong, and Xiaofeng Liu. Continual learning-based trajectory prediction with memory augmented networks. *Knowledge-Based Systems*, 258:110022, 2022. [2](#)
- [48] Yifu Zhang, Peize Sun, Yi Jiang, Dongdong Yu, Fucheng Weng, Zehuan Yuan, Ping Luo, Wenyu Liu, and Xinggang Wang. Bytetrack: Multi-object tracking by associating every detection box. 2022. [5](#)
- [49] Tianyang Zhao, Yifei Xu, Mathew Monfort, Wongun Choi, Chris Baker, Yibiao Zhao, Yizhou Wang, and Ying Nian Wu. Multi-agent tensor fusion for contextual trajectory prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12126–12134, 2019. [2](#)

# Adaptive Human Trajectory Prediction via Latent Corridors

## Supplementary Material

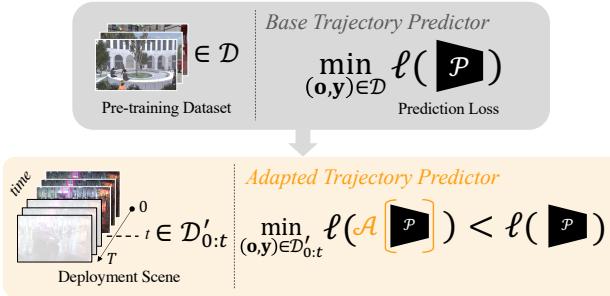


Figure 7. **Adaptive trajectory prediction.** ATP, formulated in Sec. 3, allows a pre-trained predictor  $\mathcal{P}$  to adapt to a new deployment scene by learning over time on the deployment scene. Once adaptation has occurred, the adapted predictor  $\mathcal{A}[\mathcal{P}]$  should perform better on the deployment scene.

In the supplementary material, we provide an ablation on the implementation of our prompting method and a visual overview of our ATP problem formulation. We also show additional qualitative results on MOTSynth and EarthCam data, as well as more detailed quantitative results.

## 8. Details on Latent Corridors for Adaptive Trajectory Prediction

### 8.1. Adaptive Trajectory Prediction Visualization

We provide an illustrative overview of adaptive trajectory prediction problem, formulated in Sec. 3, in Fig. 7.

### 8.2. Ablation on Prompt Location/Size

We ablated the prompt location and method of combining with the input. For the input location, we experimented with combining the prompt with several parts of the input to  $\mathcal{P}$ ,  $[\mathbf{M}_{\tau-H:\tau}, S]$ : all of the input heatmaps  $\mathbf{M}_{\tau-H:\tau}$ , just the first input heatmap  $\mathbf{M}_0$ , just the segmentation map  $S$ , and all of the inputs. For method of combination, we experimented with element-wise summing and element-wise multiplication. We ran this ablation on the latent corridors-only approach to training on MOTSynth. Results can be seen in Table 2. While summing seems to lead to better performance than multiplying, and summing to the heatmaps seems to be more helpful than summing to the segmentation, generally, the prompts are effective at improving performance on a variety of input locations. This shows promise for training latent corridors in a variety of settings.

Prompt Method	ADE	FDE
Sum to all heatmaps $\mathbf{M}_{\tau-H:\tau}$	<b>44.6</b>	<b>90.2</b>
Sum to $\mathbf{M}_0$	45.1	92.1
Sum to $S$	46.4	97.3
Sum to $\mathbf{M}_{\tau-H:\tau}$ and $S$	46.4	96.5
Multiply to all heatmaps $\mathbf{M}_{\tau-H:\tau}$	46.5	96.1
Multiply to $S$	45.8	93.6
Multiply to $\mathbf{M}_0$	47.6	101.6
Multiply to $\mathbf{M}_{\tau-H:\tau}$ and $S$	46.5	96.1

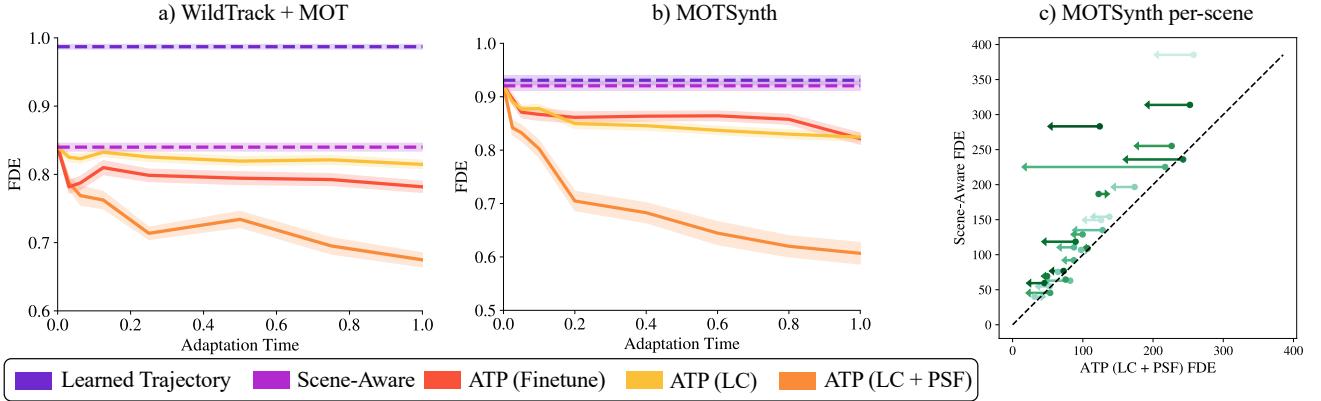
Table 2. Ablation on 437 MOTSynth scenes in the ATP latent corridor adaptation configuration. We experiment with different locations and methods of combining the prompt with the input, and find that summing the prompt to all input heatmaps yields the best result, but most combinations result in an improvement on the scene-aware baseline (see Table 1).

## 9. Additional Experimental Results on MOT-Synth, MOT and WildTrack

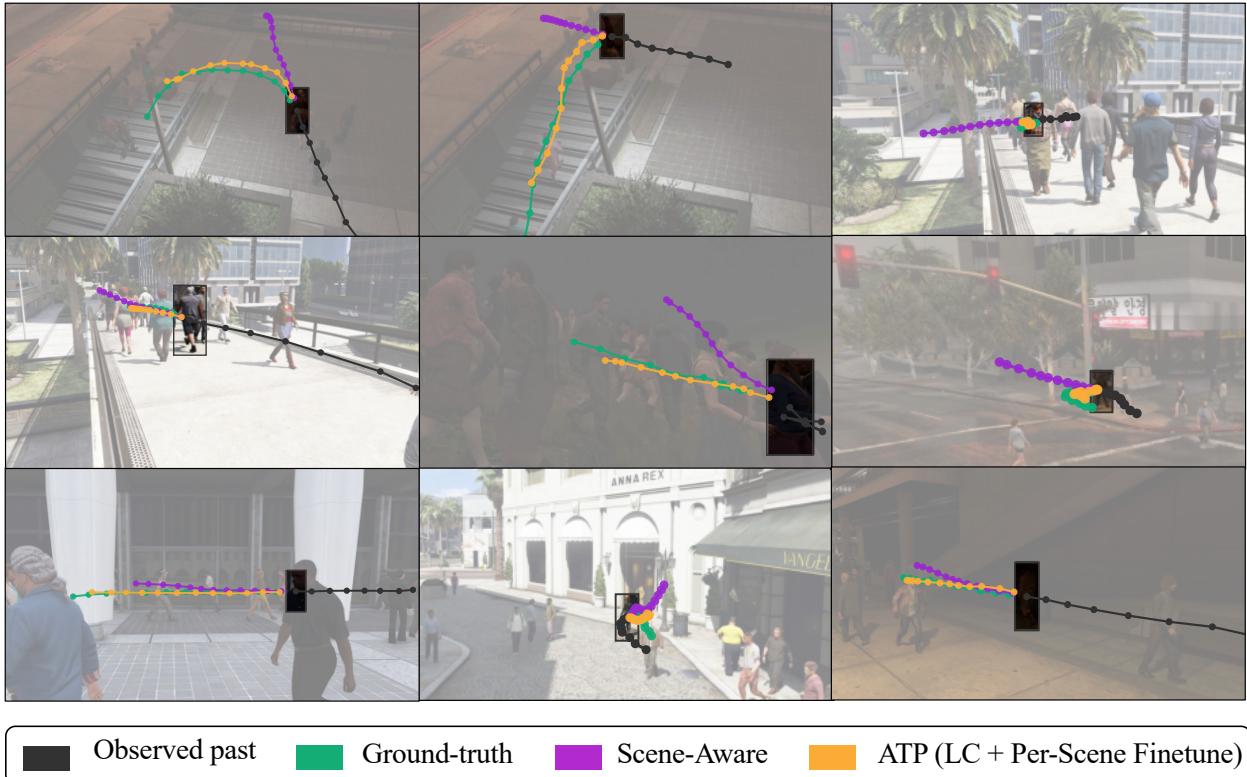
In Fig. 8, we see the FDE results for MOTSynth, WildTrack and MOT over time. Similar to results with the ADE metric, we see that on real data (Fig. 8a), the adaptive finetuning baseline is slightly better than just latent corridors, but ATP via both latent corridor prompting and per-scene finetuning largely outperforms both of those, and all adaptive methods outperform the non-adaptive baselines. On the MOTSynth data, as with ADE, the ATP via just per-scene finetuning or just latent corridors are comparable (Fig. 8b). With FDE, while some scenes have minimal gains, we see even more significant error reduction for some scenes than with ADE, of up to 91.4%, and an average FDE improvement of 33.9% from ATP via latent corridors and per-scene finetuning on 25 MOTSynth scenes (Fig. 8c).

We showcase additional qualitative results on MOT-Synth data in Fig. 9. We see that our approach is able to learn that in a scene with a large staircase, pedestrians mostly move towards the staircase, regardless of the direction of their observed history (top left and middle). We also see that our approach learns that pedestrians tend to stay on a walkway (top right). Finally, we see several examples of our approach having awareness of the 3D ground plane and predicting future trajectories that lie within the ground plane (bottom two rows).

Finally, we break down the numbers for the three ATP methods used on Wildtrack and MOT shown in Fig. 5a and Fig. 8a in Table 3. We see that on these real datasets, for both ADE and FDE, ATP via per-scene finetuning alone yields slightly better results than ATP via latent corridor



**Figure 8. Adaptation over time FDE.** As in Fig. 5, the x-axis represents normalized adaptation time in person-seconds. The y-axis represents the FDE (lower is better). Results are normalized per-scene and averaged over models trained on 25 MOTSynth scenes (a) and 7 from MOT and WildTrack (b), with shaded area  $\sigma/10$ . For the FDE metric, our methods improve on the baselines increasingly with adaptation time. Latent corridors + per-scene finetuning has the best performance, as with FDE, and ATP via just finetuning or just latent corridor learning is still comparable. c) Comparison to baseline over many MOTSynth scenes for models trained with 8% (point) and 80% (arrowhead) human-second datasets for FDE. For many deployment scenes, FDE improves significantly more with our method than ADE improved, but still, the per-scene improvements are varied.



**Figure 9. Additional MOTSynth qualitative results.** (top left and middle) From several examples of pedestrians in motion, our approach (orange) is able to learn that in this scene, most pedestrians will turn to go down the stairs, while the scene-aware baseline (purple) struggles to understand this scene-specific feature, and instead assumes that the pedestrian will continue walking in the direction of the observed history. Our model is also able to gain understanding that most pedestrians will stay on a walkway, even if they move in a direction orthogonal to it (top right). We also see many more examples of our approach having better awareness of the 3D nature of the ground plane projected into the 2D image (bottom two rows), even when the ground plane is tilted (middle middle).

Method	MOT		Wildtrack	
	ADE	FDE	ADE	FDE
ATP (Finetune)	41.8	86.9	31.7	63.5
ATP (LC)	43.0	89.2	32.9	65.9
ATP (LC + PSF)	<b>37.4</b>	<b>74.3</b>	<b>27.4</b>	<b>54.7</b>

Table 3. ATP via latent corridors, per-scene finetuning, or a combination thereof on MOT and Wildtrack. We see that on real data, adaptation via per-scene finetuning alone outperforms latent corridors alone. However, a combination of our latent corridor approach with per-scene finetuning yields significantly better results.

adaption only. However, a combination of the two significantly outperforms either alone.

## 10. Additional Results on Webcam Data

We scraped two additional 5-minute videos from EarthCam as described in Sec. 5.1, one from Rick’s Cafe in Jamaica and one from Times Square in New York City. Qualitative results for these scenes are in Fig. 10. In the top row, at the cafe, there is a complicated path through an overlook that twists down stairs towards the water. The scene-aware baseline often predicts that people will jump over a ledge, whereas our method learns the boundaries of the paths that people follow and is able to predict that people will stay within those boundaries. In the bottom row, we see that the scene-aware baseline does not recognize that people on a billboard will stay within the billboard, whereas our method is able to recognize that.

Quantitative results for these scenes, as well as for the NoLA daytime and nighttime scenes, can be seen in Table 4. We see that across the four EarthCam scenes, ATP via per-scene finetuning alone works better than using latent corridor adaptation alone (by a narrow margin on both NoLA scenes, and by a significant amount on the Rick’s Cafe scene), but a combination of the two is significantly better than any other adaptive or non-adaptive approach. Interestingly, for the Rick’s Cafe and Times Square scenes, a constant velocity baseline is better than our non-adaptive learned baselines, but our ATP approach outperforms all baselines.

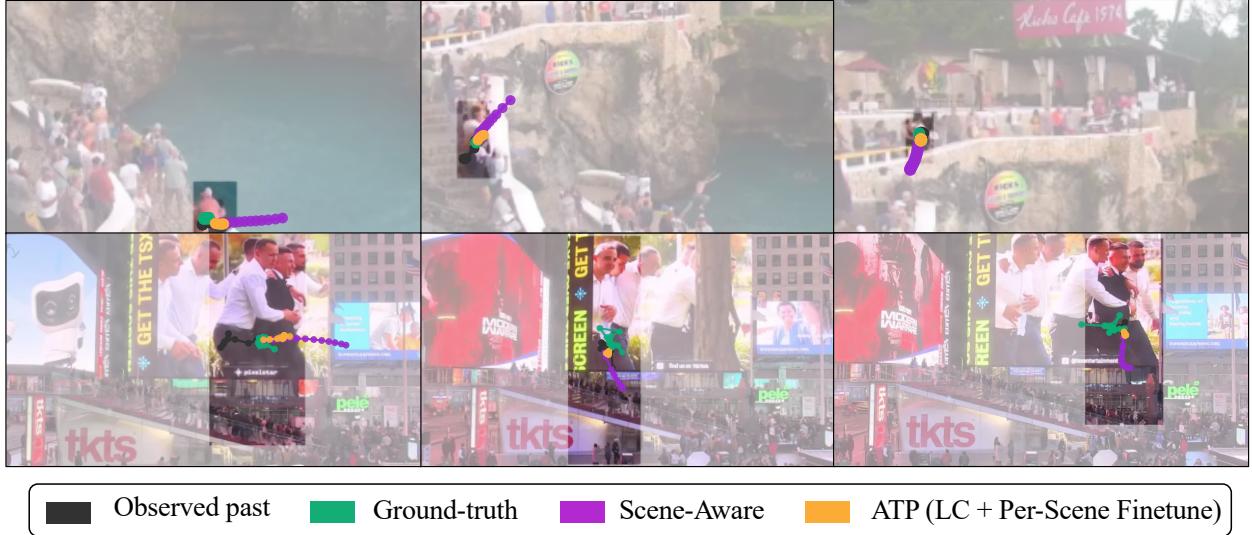


Figure 10. **Qualitative results on additional EarthCam scenes.** (top) At a cafe in Jamaica with an overlook and stairs on the edge of the water, our model (orange) is able to correctly predict that people will stick to moving up and down the path created by the stairs, while the scene-aware baseline (purple) predicts people will walk directly into the water (left) or over the edge of the path (middle, right). (bottom) In Times Square, a billboard depicts human actors in motion. The scene-aware baseline assumes that these actors will move following their observed history, while our model correctly predicts that anyone in the area of the billboard will stay in the billboard.

Method	Rick's Cafe		Times Square		NoLA Daytime		NoLA Nighttime	
	ADE	FDE	ADE	FDE	ADE	FDE	ADE	FDE
Constant velocity	9.6	16.3	15.1	26.9	36.0	70.7	48.4	94.8
Learned Traj (PECNet-Ours)	20.6	29.6	42.7	60.5	35.1	64.1	38.8	67.8
Scene-aware (YNet-Ours)	10.4	16.8	17.3	30.5	30.7	59.2	37.3	71.0
ATP (Finetune)	7.4	10.6	14.7	25.1	30.4	57.1	34.9	62.0
ATP (LC)	10.2	16.5	16.7	29.0	30.3	58.0	36.5	67.8
ATP (LC + Per-Scene FT)	<b>6.2</b>	<b>8.8</b>	<b>11.7</b>	<b>19.5</b>	<b>22.6</b>	<b>43.0</b>	<b>29.5</b>	<b>51.9</b>

Table 4. Results on four EarthCam scenes. Across all of these challenging in-the-wild scenarios, our ATP method using latent corridors and per-scene finetuning outperforms the baselines.