

Trace and Pace: Controllable Pedestrian Animation via Guided Trajectory Diffusion

Davis Rempe^{*,1,2}

Zhengyi Luo^{*,1,3}

Xue Bin Peng^{1,4}

Ye Yuan¹

Kris Kitani³

Karsten Kreis¹

Sanja Fidler^{1,5,6}

Or Litany¹

¹NVIDIA

²Stanford University

³Carnegie Mellon University

⁴Simon Fraser University

⁵University of Toronto

⁶Vector Institute

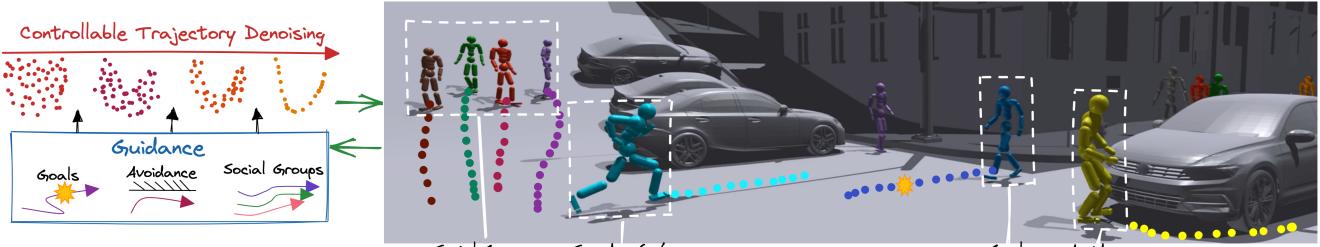


Figure 1. (Left) We propose TRACE, a trajectory diffusion model that enables user control through test-time guidance. (Right) Generated trajectories are passed to a novel physics-based humanoid controller (PACER), forming a closed-loop pedestrian animation system.

Abstract

We introduce a method for generating realistic pedestrian trajectories and full-body animations that can be controlled to meet user-defined goals. We draw on recent advances in guided diffusion modeling to achieve test-time controllability of trajectories, which is normally only associated with rule-based systems. Our guided diffusion model allows users to constrain trajectories through target waypoints, speed, and specified social groups while accounting for the surrounding environment context. This trajectory diffusion model is integrated with a novel physics-based humanoid controller to form a closed-loop, full-body pedestrian animation system capable of placing large crowds in a simulated environment with varying terrains. We further propose utilizing the value function learned during RL training of the animation controller to guide diffusion to produce trajectories better suited for particular scenarios such as collision avoidance and traversing uneven terrain. Video results are available on the [project page](#).

1. Introduction

Synthesizing high-level human behavior, in the form of 2D positional trajectories, is at the core of modeling pedestrians for applications like autonomous vehicles, urban planning, and architectural design. An important feature of such synthesis is *controllability* – generating tra-

jectories that meet user-defined objectives, edits, or constraints. For example, a user may place specific waypoints for characters to follow, specify social groups for pedestrians to travel in, or define a social distance to maintain.

Attaining controllability is straightforward for algorithmic or rule-based models of human behavior, since they have built-in objectives. In the simplest case, human trajectories can be determined by the shortest paths between control points [12], but more sophisticated heuristics have also been developed for pedestrians [3, 16], crowds [24, 50], and traffic [32, 60]. Unfortunately, algorithmically generated trajectories often appear unnatural. Learning-based approaches, on the other hand, can improve naturalness by mimicking real-world data. These methods often focus on short-term trajectory prediction using a single forward pass of a neural network [2, 11, 55, 69]. However, the ability to control these models is limited to sampling from an output trajectory distribution [37, 66] or using an expensive latent space traversal [49]. As a result, learning-based methods can predict implausible motions such as collisions with obstacles or between pedestrians. This motivates another notion of *controllability* – maintaining realistic trajectories during agent-agent and agent-environment interactions.

In this work, we are particularly interested in using controllable pedestrian trajectory models for character animation. We envision a simple interface where a user provides high-level objectives, such as waypoints and social groups, and a system converts them to *physics-based* full-body human motion. Compared to existing kinematic motion mod-

*Equal contribution

els [21, 30, 46], physics-based methods have the potential to produce high-quality motion with realistic subtle behaviors during transitions, obstacle avoidance, traversing uneven terrains, *etc.* Although there exist physics-based animation models [13, 30, 43–45, 65], controlling their behavior requires using task-specific planners that need to be re-trained for new tasks, terrains, and character body shapes.

We develop a generative model of trajectories that is data driven, controllable, and tightly integrated with a physics-based animation system for full-body pedestrian simulation (Fig. 1). Our method enables generating pedestrian trajectories that are realistic and amenable to user-defined objectives at test time. We use this trajectory generator as a planner for a physics-based pedestrian controller, resulting in a closed-loop controllable pedestrian animation system.

For trajectory generation, we introduce a **TRA**jectory Diffusion Model for Controllable PEdestrians (TRACE). Inspired by recent successes in generating trajectories through denoising [10, 22, 72], TRACE generates the future trajectory for each pedestrian in a scene and accounts for the surrounding context through a spatial grid of learned map features that is queried locally during denoising. We leverage classifier-free sampling [19] to allow training on mixed annotations (*e.g.*, with and without a semantic map), which improves controllability at test time by trading off sample diversity with compliance to conditioning. User-controlled sampling from TRACE is achieved through test-time *guidance* [8, 19, 20], which perturbs the output at each step of denoising towards the desired objective. We extend prior work [22] by introducing several analytical loss functions for pedestrians and re-formulating trajectory guidance to operate on clean trajectory outputs from the model [20], improving sample quality and adherence to user objectives.

For character animation, we develop a general-purpose **Pedestrian Animation Controller** (PACER) capable of driving physics-simulated humanoids with diverse body types to follow trajectories from a high-level planner. We focus on (1) motion quality: PACER learns from a small motion database to create natural and realistic locomotion through adversarial motion learning [44, 45]; (2) terrain and social awareness: trained in diverse terrains with other humanoids, PACER learns to move through stairs, slopes, uneven surfaces, and to avoid obstacles and other pedestrians; (3) diverse body shapes: by training on different body types, PACER draws on years of simulation experience to control a wide range of characters; (4) compatibility with high-level planners: PACER accepts 2D waypoints and can be a plugin in model for any 2D trajectory planner.

We demonstrate a controllable pedestrian animation system using TRACE as a high-level planner for PACER, the low-level animator. The planner and controller operate in a closed loop through frequent re-planning according to simulation results. We deepen their connection by guiding

TRACE with the value function learned during RL training of PACER to improve animation quality in varying tasks. We evaluate TRACE on synthetic [3] and real-world pedestrian data [4, 29, 41], demonstrating its flexibility to user-specified and plausibility objectives while synthesizing realistic motion. Furthermore, we show that our animation system is capable and robust with a variety of tasks, terrains, and characters. In summary, we contribute (1) a diffusion model for pedestrian trajectories that is readily controlled at test time through guidance, (2) a general-purpose pedestrian animation controller for diverse body types and terrains, and (3) a pedestrian animation system that integrates the two to drive simulated characters in a controllable way.

2. Related Work

Pedestrian Trajectory Prediction. Modeling high-level pedestrian behavior has been extensively studied in the context of motion prediction (forecasting). Approaches range from physics and planning-based [15, 16, 62] to recent learned methods [2, 6, 27, 55, 69]. We refer the reader to the thorough survey by Rudenko *et al.* [52] for an overview, and focus this discussion on controllability. Most forecasting work is motivated by planning for autonomous vehicles (AVs) or social robots [11] rather than *controllability* or longer-term synthesis. Rule-based models for pedestrians [3, 24, 50] and vehicle traffic [32, 60] can easily incorporate user constraints [28] making them amenable to control. However, the trajectories of these approaches are not always human-like; methods have even been developed to choose the best simulation method and tune parameters to make crowd scenarios more realistic [23].

Data-driven methods produce human-like motions, but neural network-based approaches are difficult to explicitly *control*. Some works decompose forecasting into goal prediction followed by trajectory prediction based on goals [7, 37]. These models offer limited control by selecting goal locations near a target or that minimizes an objective (*e.g.* collisions) [66]. Synthesized pedestrian behavior can also be controlled by strategically choosing a starting location [47]. STRIVE [49] showed that a VAE trajectory model can be controlled through test-time optimization in the learned latent space. Reinforcement learning (RL) agents can be controlled in crowd simulations by incorporating tasks into reward functions for training [26]. By varying the weights of different rewards, the characters can be controlled to exhibit one of several behaviors at test time [40]. Our method, TRACE, trains to mimic trajectories from data and is agnostic to any task: all controls are defined at test time, allowing flexibility to new controls after training. Instead of lengthy test-time optimization, we use guidance for control.

Controllable Character Animation. Full-body pedestrian animation typically involves a high-level task (*e.g.* trajectory following, obstacle avoidance) and low-level body con-

trol. Some methods solve both with a single network that implicitly uses high-level planning and low-level animation. GAMMA [71] trains a kinematic model to go to waypoints, while PFNN [21] follows gamepad inputs. Physics-based humanoid controllers such as AMP [45] train different models for each task, limiting their general applicability.

Two-stage methods split the task into separate high-level planning and low-level character control, where task information is only used by the planner. Planning can be done with traditional A* [12], using learned trajectory prediction [5], searching in a pre-trained latent space [30, 44, 48, 65], or using hierarchical RL [13, 43, 44, 46, 65]. DeepLoco [43], Haworth *et al.* [13], and ASE [44] utilize hierarchical RL to achieve impressive dynamic control for various tasks. They require lengthy training for both low-level and high-level controllers and often jointly train as a final step. They must also train different planners for different tasks.

Our approach follows the two-stage paradigm, with the distinction that both our high-level (TRACE) and low-level (PACER) models consume task information for pedestrian navigation: through test-time guidance and map-conditioned path following, respectively. TRACE and PACER are unaware of each other at training time, yet can be tightly integrated in a closed loop: trace-pace-retrace.

Diffusion Models and Guidance. Diffusion models have shown success in generating images [18, 39, 61], videos [17], and point clouds [70]. Guidance has been used for test-time control in several ways: classifier [8] and classifier-free [19] guidance reinforce input conditioning, while reconstruction guidance [20] has been used for coherent video generation. Gu *et al.* [10] adapt the diffusion framework for short-term pedestrian trajectory forecasting conditioned on past trajectories. Diffuser [22] generates trajectories for planning and control in robotics applications with test-time guidance. Closest to ours is the concurrent work of CTG [72], which builds on Diffuser to develop a controllable vehicle traffic model, focusing on following formalized traffic rules like speed limits. Our method contains several key differences: we encode map conditioning into an expressive feature grid queried in denoising, we use classifier-free sampling to enable multi-dataset training and test-time flexibility, we re-formulate guidance to operate on clean model outputs, and we link with a low-level animation model using value function guidance.

3. Method

To model high-level pedestrian behavior, we first introduce the controllable trajectory diffusion model (TRACE). In Sec. 3.2, we detail our low-level physics-based pedestrian controller, PACER, and in Sec. 3.3 how they can be combined into an end-to-end animation system.

3.1. Controllable Trajectory Diffusion

Problem Setting. Our goal is to learn high-level pedestrian behavior in a way that can be *controlled* at test time. For pedestrian animation, we focus on two types of control: (1) user specification, *e.g.*, goal waypoints, social distance, and social groups, and (2) physical plausibility, *e.g.*, avoiding collisions with obstacles or between pedestrians.

We formulate synthesizing pedestrian behavior as an agent-centric trajectory forecasting problem. At each time step, the model outputs a future trajectory plan for a target *ego* agent conditioned on that agent’s past, the past trajectories of all neighboring agents, and the semantic map context. Formally, at timestep t we want the future state trajectory $\tau_s = [\mathbf{s}_{t+1} \ \mathbf{s}_{t+2} \ \dots \ \mathbf{s}_{t+T_f}]$ over the next T_f steps where the state $\mathbf{s} = [x \ y \ \theta \ v]^T$ includes the 2D position (x, y) , heading angle θ , and speed v . We assume this state trajectory is actually the result of a sequence of actions [72] defined as $\tau_a = [\mathbf{a}_{t+1} \ \mathbf{a}_{t+2} \ \dots \ \mathbf{a}_{t+T_f}]$ where each action $\mathbf{a} = [\dot{v} \ \dot{\theta}]^T$ contains the acceleration \dot{v} and yaw rate $\dot{\theta}$. The state trajectory can be recovered from the initial state and action trajectory as $\tau_s = f(\mathbf{s}_t, \tau_a)$ using a given dynamics model f . The full state-action trajectory is then denoted as $\tau = [\tau_s; \tau_a]$. To predict the future trajectory, the model receives as input the past state trajectory of the ego pedestrian $\mathbf{x}^{\text{ego}} = [\mathbf{s}_{t-T_p} \ \dots \ \mathbf{s}_t]$ along with the past trajectories of N neighboring pedestrians $X^{\text{neigh}} = \{\mathbf{x}^i\}_{i=1}^N$. It also gets a crop of the rasterized semantic map $\mathcal{M} \in \mathbb{R}^{H \times W \times C}$ in the local frame of the ego pedestrian at time t . These inputs are summarized as the *conditioning* context $C = \{\mathbf{x}^{\text{ego}}, X^{\text{neigh}}, \mathcal{M}\}$.

Our key idea is to train a diffusion model to conditionally generate trajectories, which can be *guided* at test time to enable controllability. For simplicity, the following formulation uses the full trajectory notation τ , but in practice, the state trajectory is always a result of actions, *i.e.*, diffusion/denoising are on τ_a which determines the states through f . Next, we summarize our diffusion framework, leaving the details to the supplementary material.

3.1.1 Trajectory Diffusion Model

We build on Diffuser [22] and generate trajectories through iterative denoising, which is learned as the reverse of a pre-defined *diffusion* process [18, 57]. Starting from a clean future trajectory $\tau^0 \sim q(\tau^0)$ sampled from the data distribution, the forward noising process produces a sequence of progressively noisier trajectories $(\tau^1, \dots, \tau^k, \dots, \tau^K)$ by adding Gaussian noise at each process step k :

$$q(\tau^{1:K} \mid \tau^0) := \prod_{k=1}^K q(\tau^k \mid \tau^{k-1}) \quad (1)$$

$$q(\tau^k \mid \tau^{k-1}) := \mathcal{N}(\tau^k; \sqrt{1 - \beta_k} \tau^{k-1}, \beta_k \mathbf{I})$$

where β_k is the variance at each step of a fixed schedule, and with a large enough K we get $q(\tau^K) \approx \mathcal{N}(\tau^K; \mathbf{0}, \mathbf{I})$.

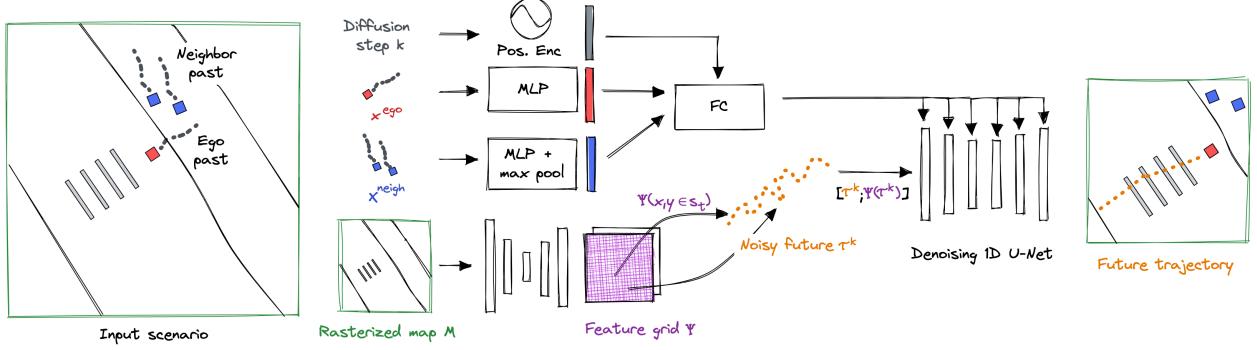


Figure 2. Trajectory diffusion model (TRACE). Future trajectory denoising is conditioned on past and neighbor motion by adding processed features to intermediate U-Net features. Map conditioning is provided through a feature grid queried along the noisy input trajectory.

TRACE learns the reverse of this process so that the sampled noise can be *denoised* into plausible trajectories. Each step of this reverse process is conditioned on C :

$$p_\phi(\tau^{k-1} \mid \tau^k, C) := \mathcal{N}(\tau^{k-1}; \mu_\phi(\tau^k, k, C), \Sigma_k) \quad (2)$$

where ϕ are model parameters and Σ_k is from a fixed schedule. TRACE learns to parameterize the mean of the Gaussian distribution at each step of the denoising process.

Training and Classifier-Free Sampling. Importantly for guidance, the network does *not* directly output μ . Instead, at every step it learns to predict the final clean trajectory τ^0 , which is then used to compute μ [39]. Training supervises this network output $\hat{\tau}^0$ with ground truth future trajectories (*i.e.* denoising score matching [18, 58, 63]):

$$L = \mathbb{E}_{\epsilon, k, \tau^0, C} [| | | \tau^0 - \hat{\tau}^0 | | |^2] \quad (3)$$

where τ^0 and C are sampled from the training dataset, $k \sim \mathcal{U}\{1, 2, \dots, K\}$ is the step index, and $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is used to corrupt τ^0 to give the noisy input trajectory τ^k .

Our training procedure allows the use of *classifier-free sampling*¹ at test time, which has been shown to improve compliance to conditioning in diffusion models [19]. We simultaneously train both a conditional model $\mu_\phi(\tau^k, k, C)$ and unconditional model $\mu_\phi(\tau^k, k)$ by randomly dropping out conditioning during training. At test time, predictions from both models are combined with weight w as:

$$\tilde{\epsilon}_\phi = \epsilon_\phi(\tau^k, k, C) + w (\epsilon_\phi(\tau^k, k, C) - \epsilon_\phi(\tau^k, k)) \quad (4)$$

where ϵ_ϕ is the model’s prediction of how much noise was added to the clean trajectory to produce the input τ^k ; it is straightforward to compute from μ_ϕ [39].

Note that $w > 0$ and $w < 0$ increase and decrease the effect of conditioning, respectively, while $w = 0$ and $w = -1$ result in the purely conditional or unconditional model, respectively. This flexibility allows a user to trade off respecting

¹we refer to it as “sampling” instead of the common term “guidance” [19] to avoid confusion with the guidance introduced in Sec. 3.1.2

conditioning with trajectory diversity, which benefits controllability (see Sec. 4.2). This approach also enables training on multiple distinct datasets with varying annotations: conditioning is already being dropped out randomly, so it is easy to use mixed data with subsets of the full conditioning. Since there are pedestrian datasets with diverse motions but no semantic maps [29, 41], and others with limited motions but detailed maps [4], we find mixed training is beneficial to boost diversity and controllability (see Sec. 4.2).

Architecture. As shown in Fig. 2, TRACE uses a U-Net similar to [22] that has proven effective for trajectories. The input trajectory τ_k at step k is processed by a sequence of 1D temporal convolutional blocks that progressively down and upsample the sequence in time, leveraging skip connections. A key challenge is how to condition the U-Net on C to predict trajectories that comply with the map and other pedestrians. To incorporate step k , ego past x^{ego} , and neighbor past X^{neigh} , we use a common approach [20, 22] that extracts a single conditioning feature and adds it to the intermediate trajectory features within each convolutional block. For the map M , we encode with a 2D convolutional network into a feature grid, where each pixel contains a high-dimensional feature. At step k of denoising, each 2D position $(x, y) \in \tau^k$ is queried by interpolating into the grid to give a feature trajectory, which is concatenated to τ^k and becomes the U-Net input. Intuitively, this allows learning a localized representation that can benefit subtle map interactions such as obstacle avoidance.

3.1.2 Controllability through Clean Guidance

After training TRACE to generate realistic trajectories, controllability is implemented through test-time *guidance*. Intuitively, guidance nudges the sampled trajectory at each step of denoising towards a desired outcome. Let $\mathcal{J}(\tau)$ be a guidance loss function measuring how much a trajectory τ violates a user objective. This may be learned [22] or an analytical differentiable function [72]. Guidance uses the gradient of \mathcal{J} to perturb the predicted mean from the model at each denoising step such that the right side of Eq. (2)

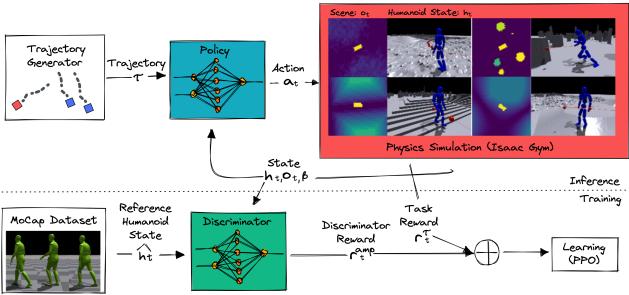


Figure 3. Pipeline: Pedestrian Animation Controller (PACER).

becomes $\mathcal{N}(\boldsymbol{\tau}^{k-1}; \tilde{\boldsymbol{\mu}}_\phi(\boldsymbol{\tau}^k, k, C), \boldsymbol{\Sigma}_k)$ where $\tilde{\boldsymbol{\mu}}$ is the perturbed (guided) mean. Prior work [22, 72] directly perturbs the *noisy* network-predicted mean with

$$\tilde{\boldsymbol{\mu}} = \boldsymbol{\mu} - \alpha \boldsymbol{\Sigma}_k \nabla_{\boldsymbol{\mu}} \mathcal{J}(\boldsymbol{\mu}) \quad (5)$$

where α determines the guidance strength. Note that Eq. (5) evaluates \mathcal{J} at the noisy mean, so learned loss functions must be trained at varying noise levels and analytic loss functions may suffer from numerical issues.

To avoid this, we build upon “reconstruction guidance”, which operates on the *clean* model prediction $\hat{\boldsymbol{\tau}}^0$ [20]. We extend the guidance formulation introduced in [20] for temporal video upsampling to work with arbitrary loss functions. At each denoising step with input $\boldsymbol{\tau}^k$, we first perturb the clean trajectory predicted from the network $\hat{\boldsymbol{\tau}}^0$ with

$$\tilde{\boldsymbol{\tau}}^0 = \hat{\boldsymbol{\tau}}^0 - \alpha \boldsymbol{\Sigma}_k \nabla_{\boldsymbol{\tau}^k} \mathcal{J}(\hat{\boldsymbol{\tau}}^0), \quad (6)$$

then compute $\tilde{\boldsymbol{\mu}}$ in the same way as we would in Eq. (2), *i.e.*, as if $\tilde{\boldsymbol{\tau}}^0$ were the output of the network. Note that the gradient is evaluated wrt the noisy input trajectory $\boldsymbol{\tau}^k$ rather than the clean $\hat{\boldsymbol{\tau}}^0$, requiring backpropagation through the denoising model. We formulate several analytical guidance objectives like waypoint reaching, obstacle avoidance, collision avoidance, and social groups (see Sec. 4.1, 4.2). A learned RL value function can also be used (Sec. 4.3).

3.2. Physics-Based Pedestrian Animation

To enable full-body pedestrian simulation, we design the Pedestrian Animation ContROLLER (PACER) to execute the 2D trajectories generated by TRACE in a physics simulator.

Background: Goal-Conditioned RL. Our framework (Fig. 3) follows the general goal-conditioned reinforcement learning framework, where a goal-conditioned policy π_{PACER} is trained to follow 2D target trajectories specified by $\boldsymbol{\tau}_s$. The task is formulated as a Markov Decision Process (MDP) defined by a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, R, \gamma)$ of states, actions, transition dynamics, reward function, and discount factor. The state \mathcal{S} , transition dynamics \mathcal{T} , and reward R are calculated by the environment based on the current simulation and goal, while the action \mathcal{A} is computed by the policy π_{PACER} . The policy’s objective is to maximize

the discounted return $\mathbb{E} \left[\sum_{t=1}^T \gamma^{t-1} r_t \right]$ where r_t is the reward per timestep. We utilize Proximal Policy Optimization (PPO) [56] to find the optimal control policy π_{PACER} .

Terrain, Social, and Body Awareness. To create a controller that can simulate crowds in realistic 3D scenes (*e.g.* scans, neural reconstructions, or artist-created meshes (Fig. 1)), our humanoid must be terrain aware, socially aware of other agents, and support diverse body types. We use a humanoid model that conforms to the kinematic structure of SMPL [31], and is automatically generated using a procedure similar to [33, 34, 68]. Our control policy $\pi_{\text{PACER}}(\mathbf{a}_t | \mathbf{h}_t, \mathbf{o}_t, \beta, \boldsymbol{\tau}_s)$ is conditioned on the state of the simulated character \mathbf{h}_t , environmental features \mathbf{o}_t , body type β , and goal trajectory $\boldsymbol{\tau}_s$. The environment input is a rasterized local height and velocity map of size $\mathbf{o}_t \in \mathbb{R}^{64 \times 64 \times 3}$, which gives agents crucial information about their surroundings. To allow for social awareness, nearby humanoids are represented as a cuboid and rendered on the global height map. In this way, each humanoid views other people as dynamic obstacles to avoid. Obstacle and interpersonal avoidance are learned by using obstacle collision as a termination condition. By conditioning and training with different body parameters β our policy learns to adapt to characters with diverse morphologies.

Realistic Motion through Adversarial Learning. To learn the optimal control policy π_{PACER} that (1) follows a 2D trajectory closely and (2) creates realistic pedestrian motions, we follow Adversarial Motion Prior (AMP) [45]. AMP uses a motion discriminator to encourage the policy to generate motions that are similar to the movement patterns contained in a dataset of motion clips recorded by human actors. The discriminator $D(\mathbf{h}_{t-10:t}, \mathbf{a}_t)$ is then used to specify a motion style reward r_t^{amp} for training the policy. The style reward is combined with a trajectory following reward r_t^τ and an energy penalty r_t^{energy} [9] to produce the total reward $r_t = r_t^{\text{amp}} + r_t^\tau + r_t^{\text{energy}}$. To mitigate artifacts arising from asymmetric gaits, such as limping, we utilize the motion-symmetry loss proposed by [67]:

$$L_{\text{sym}}(\theta) = \| \pi_{\text{PACER}}(\mathbf{h}_t, \mathbf{o}_t, \beta, \boldsymbol{\tau}_s) - \Phi_a(\pi_{\text{PACER}}(\Phi_s(\mathbf{h}_t, \mathbf{o}_t, \beta, \boldsymbol{\tau}_s))) \|^2, \quad (7)$$

where Φ_s and Φ_a mirror the state and action along the character’s sagittal plane. This loss encourages the policy to produce more symmetric motions, leading to natural gaits. During training, random terrains are generated following the procedure used in [53]. We create stairs, slopes, uneven terrains, and obstacles consisting of random polygons. Character morphology is also randomized by sampling a gender and body shape β parameters. More details are available in the supplementary material.

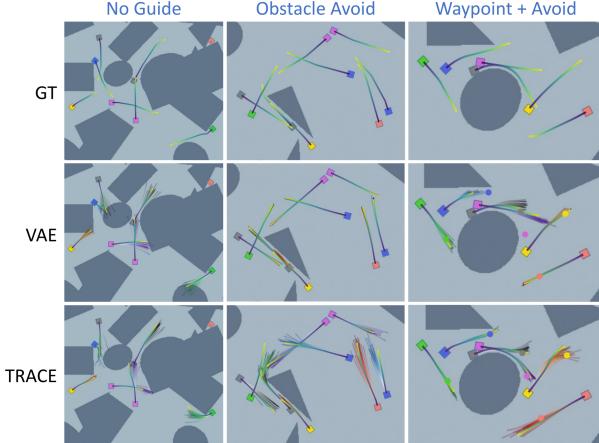


Figure 4. Guidance results on *ORCA-Maps*. For VAE and TRACE, 20 samples are visualized for each pedestrian (the boxes) along with the final trajectory chosen via filtering which is bolded.

3.3. Controllable Pedestrian Animation System

The high-level trajectory planning from TRACE is combined with the low-level character control from PACER to create an end-to-end pedestrian animation system. The two components are trained independently, but at runtime they operate in a closed feedback loop: PACER follows planned trajectories for 2s before TRACE re-planning, taking past character motion from PACER as input. By combining terrain and social awareness of PACER with collision avoidance guidance, both high and low-level systems are task-aware and work in tandem to prevent collisions and falls.

Value Function as Guidance. To enable tighter two-way coupling between TRACE and PACER, in Sec. 4.3 we explore using the value function learned during RL training of PACER to guide trajectory diffusion. The value function predicts expected future rewards and is aware of body pose and surrounding terrain and agents. Using the value function to guide denoising encourages TRACE to produce trajectories that are easier to follow and better suited to the current terrain (which TRACE is unaware of otherwise). Unlike Diffuser [22], which requires training a reward function with samples from the diffusion model at varying noise levels, our guidance (Eq. (6)) operates on clean trajectories so we can use the value function directly from RL training.

4. Experiments

We first demonstrate the controllability of TRACE when trained on synthetic (Sec. 4.1) and real-world (Sec. 4.2) pedestrian data. Sec. 4.3 evaluates our full animation system on several tasks and terrains. **Video results** are provided in the supplementary material.

Implementation Details. TRACE is trained to predict 5s of future motion from 3s of past motion (both at 10Hz), and uses $K=100$ diffusion steps. During training, map and neighbor conditioning inputs are independently dropped

with 10% probability. At test time, we sample (and guide) multiple future trajectories for each pedestrian in a scene and choose one with the lowest guidance loss, which we refer to as *filtering*. PACER operates at 30Hz; we randomly sample terrain, body type, and procedural 2D trajectories during training and use a dataset of locomotion sequences from AMASS [35]. All physics simulations are performed using NVIDIA’s Isaac Gym simulator [36].

Datasets. The *ORCA* dataset (Sec. 4.1) contains synthetic trajectory data from 10s scenes generated using the ORCA crowd simulator [3]. Up to 20 agents are placed in a $15m \times 15m$ environment with ≤ 20 static primitive obstacles. Agent placement and goal velocity, along with obstacle placement and scale, are randomized per scene. The dataset contains two distinct subsets: *ORCA-Maps* has many obstacles but few agents, while *ORCA-Interact* has no obstacles (*i.e.* no map annotations) but many agents.

For real-world data (Sec. 4.2), we use ETH/UCY and nuScenes. ETH/UCY [29, 41] is a common trajectory forecasting benchmark that contains scenes with dense crowds and interesting pedestrian dynamics but does not have semantic maps. nuScenes [4] contains 20s driving scenes in common street settings. We convert the pedestrian bounding-box annotations to 2D trajectories and use them for training and evaluation. Detailed semantic maps are also annotated with layers for roads, crosswalks, and sidewalks.

Metrics. We care about trajectory plausibility and meeting user controls. Controllability is evaluated with a *Guidance Error* that depends on the task: *e.g.*, for avoidance objectives this is collision rate, while the waypoint error measures the minimum distance from the trajectory. *Obstacle and Agent Collision Rates* measure the frequency of collisions. *Realism* is measured at the dataset or trajectory level by (1) computing the Earth Mover’s Distance (EMD) between the generated and ground truth test-set histograms of trajectory statistics (*e.g.* velocity, longitudinal/lateral acceleration) [66], or (2) measuring the mean accelerations of each trajectory, assuming pedestrians generally move smoothly.

4.1. Augmenting Crowd Simulation

We first evaluate TRACE trained on *ORCA-Maps* and *ORCA-Interact*. These provide a clean test bed for comparisons since there is a clear definition of correct pedestrian behavior – no obstacle or agent collisions are present in the data. All methods operate in an open loop by predicting a single 5s future for each pedestrian. This way, compounding errors inherent to closed-loop operation are not a factor.

Results for single and multi-objective guidance on the *ORCA-Maps* test set are shown in Tab. 1. TRACE is compared to a VAE baseline [49] adapted to our setup, which achieves controllability through test-time latent optimization. This is a strong baseline that generally works well, but requires expensive optimization at test time. We also

Guide	Method	Guidance		Collision Rate		Realism (EMD)		
		Error		Obstacle	Agent	Vel	Lon Acc	Lat Acc
None	VAE [49]	–		0.076	0.118	0.038	0.039	0.040
Obstacle Avoid	TRACE	–	0.050	0.132	0.029	0.008	0.009	
	VAE [49]	0.018		0.018	0.116	0.040	0.036	0.039
	TRACE-Filter	0.018		0.018	0.123	0.019	0.011	0.015
	TRACE-Noisy	0.015		0.015	0.125	0.021	0.012	0.017
Agent Avoid	TRACE	0.014	0.014	0.124	0.020	0.011	0.017	
	VAE [49]	0.010		0.075	0.010	0.041	0.038	0.039
	TRACE-Filter	0.049	0.050	0.049	0.031	0.012	0.013	
	TRACE-Noisy	0.000		0.056	0.000	0.035	0.013	0.012
Waypoint	TRACE	0.000		0.058	0.000	0.025	0.010	0.012
	VAE [49]	0.078		0.051	0.092	0.070	0.031	0.033
	TRACE-Filter	0.333		0.046	0.112	0.044	0.013	0.013
	TRACE-Noisy	0.129		0.052	0.110	0.067	0.038	0.033
Waypoint & Obs Avoid	TRACE	0.105		0.048	0.093	0.057	0.013	0.014
	VAE [49]	0.207		0.021	0.015	0.053	0.032	0.032
	TRACE-Filter	0.527		0.023	0.096	0.025	0.014	0.016
	TRACE-Noisy	0.236		0.022	0.017	0.057	0.025	0.022
& Agt Avoid	TRACE	0.211		0.021	0.009	0.036	0.007	0.009

Table 1. Guidance evaluation on *ORCA-Maps* dataset. TRACE using full diffusion guidance improves upon VAE latent optimization and selective sampling (*TRACE-Filter*) in terms of meeting objectives, while maintaining strong realism.

Guide	Method	Train Data	Guidance		Realism (Mean)		
			w	Error	Lon Acc	Lat Acc	
Waypoint	VAE [49]	Mixed	–	0.340	0.193	0.172	
	TRACE	nuScenes	-0.5	0.421	0.177	0.168	
	TRACE	Mixed	0.0	0.551	0.159	0.145	
	TRACE	Mixed	-0.5	0.366	0.140	0.132	
Waypoint perturbed	VAE [49]	Mixed	–	0.962	0.443	0.441	
	TRACE	nuScenes	-0.5	0.977	0.239	0.238	
	TRACE	Mixed	0.0	1.129	0.233	0.218	
	TRACE	Mixed	-0.5	0.802	0.212	0.204	
Social groups	VAE [49]	Mixed	–	0.297	0.109	0.104	
	TRACE	nuScenes	-0.5	0.287	0.155	0.158	
	TRACE	Mixed	0.0	0.244	0.110	0.101	
	TRACE	Mixed	-0.5	0.245	0.094	0.087	

Table 2. Guidance evaluation on nuScenes. Training on mixed data and using $w < 0$ for classifier-free sampling are important to achieve controllability for out-of-distribution objectives.

compare to two ablations: *TRACE-Filter* samples from the diffusion model *without guidance* and chooses the best sample according to the guidance loss (similar to [66]), while *TRACE-Noisy* uses the guidance formulated in Eq. (5) from prior works [22, 72]. Models are trained on the combined dataset of *ORCA-Maps* (with map annotations) and *ORCA-Interact* (no map annotations). The guidance losses are: **None** samples randomly with no guidance; **Obstacle avoid** discourages collisions between map obstacles and pedestrian bounding boxes; **Agent avoid** discourages collisions between pedestrians by denoising all their futures in a scene jointly; **Waypoint** encourages a trajectory to pass through a goal location at any point in the planning horizon. For this experiment, the waypoint is set as the position of each pedestrian at 4s into the future in the ground truth data. These are *in-distribution* objectives, since they reinforce behavior already observed in the ground truth data.

In Tab. 1, TRACE successfully achieves all objectives through the proposed guidance. It is competitive or better than the VAE optimization in terms of guidance, while maintaining velocity and acceleration distributions closer to

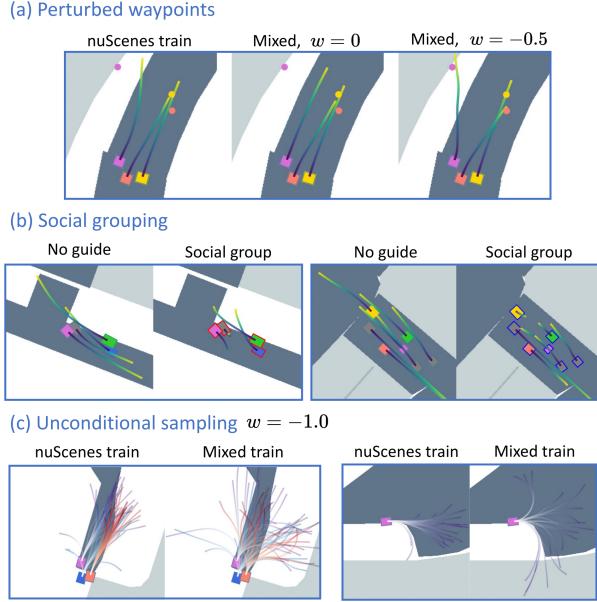


Figure 5. nuScenes results demonstrating flexibility of TRACE. (a) Using mixed training and $w = -0.5$ is best for noisy waypoints. (b) Social group guidance encourages sets of pedestrians to stay close. (c) Mixed training (ETH/UCY+nuScenes) learns a more diverse distribution as demonstrated by unconditional sampling.

ground truth as indicated by *Realism*. Fig. 4 shows that random samples from the VAE contain collisions, and using latent optimization for controllability gives similar local minima across samples thereby limiting diversity compared to TRACE. Finally, using our proposed clean guidance (Eq. (6)) instead of the noisy version produces consistently better results in guidance and realism.

4.2. Real-world Data Evaluation

We next evaluate controllability when trained on real-world data, and focus on *out-of-distribution* (OOD) guidance objectives to emphasize the flexibility of our approach. In this experiment, methods operate in a *closed loop*: pedestrians are rolled out for 10s and re-plan at 1Hz. Results on a held out nuScenes split are shown in Tab. 2. We compare TRACE trained on mixed data (ETH/UCY+nuScenes), after training on nuScenes only, and using two different classifier-free sampling weights w . Along with in-distribution **Waypoint** (now at 9s into the future), two additional objectives are evaluated: **Waypoint perturbed** uses a noisily perturbed ground truth future position (at 9s), requiring pedestrians to go off sidewalks or into streets to reach the goal; **Social groups** specifies groups of agents to stay close and travel together. Groups are set heuristically based on spatial proximity and velocity at initialization.

In Tab. 2, we observe that OOD flexibility requires (1) training on mixed data, and (2) classifier-free sampling. Since nuScenes data is less diverse (people tend to follow the sidewalk), TRACE trained on just nuScenes struggles

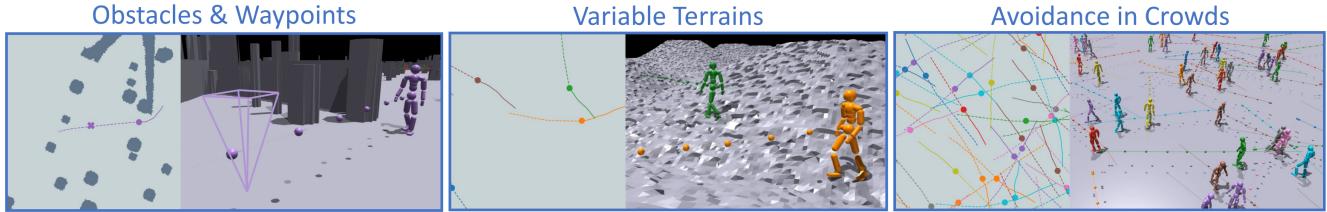


Figure 6. Our animation system enables avoiding obstacles, meeting goals, traversing variable terrains, and large crowds.

Terrain	Guide	Fail Rate	Traj Follow Error	Discrim Reward
Random	Procedural	0.133	0.680	1.950
	None	0.093	0.104	1.887
	Waypoint	0.107	0.111	2.113
Obstacles	Procedural	0.307	0.948	2.278
	None	0.125	0.093	2.512
	Obs Avoid	0.063	0.089	2.521
Flat (Crowd)	Procedural	0.127	0.371	2.320
	None	0.087	0.082	2.374
	Agt Avoid	0.013	0.071	2.402

Table 3. Closed-loop animation results. Our system successfully follows waypoints and avoids collisions in a variety of terrains, and additional guidance improves performance.

to hit perturbed waypoints. Though the VAE is trained on mixed data, it struggles to produce diverse dynamics on the nuScenes maps to achieve OOD objectives, even though it uses 200 optimization steps ($2\times$ more than the diffusion steps $K=100$ in TRACE). TRACE reaches OOD objectives using classifier-free sampling with $w=-0.5$ to down-weight the conditioning of the semantic map and leverage diverse trajectories learned from ETH/UCY. The flexibility of TRACE is further highlighted in Fig. 5.

4.3. Controllable Pedestrian Animation

Finally, we demonstrate our full controllable pedestrian animation system. TRACE is trained on *ORCA* and used as a planner for the pre-trained PACER without any fine-tuning. We evaluate the animations by: *Fail Rate* measures the fraction of agents that fall down or collide with an obstacle or other agent, *Trajectory Following Error* measures the average deviation of the character from TRACE’s plan, and *Discriminator Reward* is the mean reward returned by the adversarial motion prior used to train PACER, which measures how human-like a generated motion appears.

Tab. 3 evaluates the animations from our system using TRACE with and without guidance in various settings: *Random* is an assortment of smooth and rough slopes and stairs with varying difficulties, *Obstacles* is a flat terrain with large obstacles, and *Flat* is a flat terrain with pedestrians spawned in a crowd of 30. For each setting, 600 rollouts of 10s are simulated across 30 characters with random bodies from AMASS [35]. To put the difficulty of environments and discriminator rewards in context, we also include metrics when using the (terrain and obstacle unaware) *Procedural*

Terrain	Guide	Waypoint Value	Waypoint Error	Fail Rate	Traj Follow Error	Discrim Reward
Random	✓		0.541	0.107	0.111	2.113
	✓	✓	0.481	0.100	0.112	2.162
Obstacles	✓		1.065	0.220	0.138	2.552
	✓	✓	0.929	0.178	0.113	2.609
Flat (Crowd)	✓		0.248	0.063	0.084	2.555
	✓	✓	0.175	0.053	0.084	2.607

Table 4. Using the value function learned in RL training as guidance improves quality of trajectory following and robustness to varying terrains, obstacles, and other agents.

dural trajectory generation method used to train PACER.

Our combined system performs well in the physically-simulated environment with TRACE providing easy-to-follow trajectories resulting in high-quality animations from PACER, as evaluated by the discriminator. Diffusion guidance can further improve failure rates, especially for avoiding agent collisions in dense crowds. Fig. 6 shows some qualitative applications of our animation system and we highly encourage viewing the supplementary **video results** to qualitatively evaluate the motion quality. Tab. 4 shows the effect of using the learned value function from training PACER as a guidance loss for TRACE. In each setting, adding value guidance in addition to waypoint guidance makes trajectories easier to follow, reduces failures, and improves the discriminator reward. As a result, waypoint guidance error also improves.

5. Discussion

We have introduced a controllable trajectory diffusion model, a robust physics-based humanoid controller, and an end-to-end animation system that combines the two. This represents an exciting step in being able to control the high-level behavior of learned pedestrian models, and opens several directions for future work. First is improving the efficiency of sampling from trajectory diffusion models to make them real-time: TRACE currently takes 1-3s to sample for a single character, depending on the guidance used (see the supplement for full analysis). Recent work in diffusion model distillation [38] offers a potential solution. In addition to high-level motion controllability, exploring how diffusion models can be extended to low-level full-body character control is an interesting next step.

Acknowledgments. Davis Rempe was supported by an NVIDIA Graduate Fellowship. The authors thank Ziyuan Zhong and the AV Research Group for helpful discussions on trajectory diffusion.

References

- [1] Principles of robot autonomy i, lecture 1 course notes: Mobile robot kinematics. <https://stanfordasl.github.io/aa274a/pdfs/notes/lecture1.pdf>. Accessed: 2022-11-15. 12
- [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. 1, 2
- [3] Jur van den Berg, Stephen J Guy, Ming Lin, and Dinesh Manocha. Reciprocal n-body collision avoidance. In *Robotics Research: The 14th International Symposium (ISRR)*, pages 3–19. Springer, 2011. 1, 2, 6
- [4] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liang, Qiang Xu, Anush Krishnan, Yu Pan, Giacarlo Baldan, and Oscar Beijbom. muscenes: A multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. 2, 4, 6, 17
- [5] Zhe Cao, Hang Gao, Karttikeya Mangalam, Qi-Zhi Cai, Minh Vo, and Jitendra Malik. Long-term human motion prediction with scene context. In *European Conference on Computer Vision*, pages 387–404. Springer, 2020. 3
- [6] Yuning Chai, Benjamin Sapp, Mayank Bansal, and Dragomir Anguelov. Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction. In *Conference on Robot Learning (CoRL)*, pages 86–99. PMLR, 2020. 2
- [7] Patrick Dendorfer, Aljosa Osep, and Laura Leal-Taixé. Goalgan: Multimodal trajectory prediction based on goal position estimation. In *Proceedings of the Asian Conference on Computer Vision*, 2020. 2
- [8] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 2, 3
- [9] Zipeng Fu, Xuxin Cheng, and Deepak Pathak. Deep whole-body control: Learning a unified policy for manipulation and locomotion. *ArXiv*, abs/2210.10044, 2022. 5
- [10] Tianpei Gu, Guangyi Chen, Junlong Li, Chunze Lin, Yongming Rao, Jie Zhou, and Jiwen Lu. Stochastic trajectory prediction via motion indeterminacy diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17113–17122, 2022. 2, 3
- [11] 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. 1, 2
- [12] Mohamed Hassan, Duygu Ceylan, Ruben Villegas, Jun Saito, Jimei Yang, Yi Zhou, and Michael J Black. Stochastic scene-aware motion prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11374–11384, 2021. 1, 3
- [13] M. Brandon Haworth, Glen Berseth, Seonghyeon Moon, Petros Faloutsos, and Mubbasis Kapadia. Deep integration of physical humanoid control and crowd navigation. *Proceedings of the 13th ACM SIGGRAPH Conference on Motion, Interaction and Games*, 2020. 2, 3
- [14] 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. 13
- [15] Dirk Helbing, Illés Farkas, and Tamas Vicsek. Simulating dynamical features of escape panic. *Nature*, 407(6803):487–490, 2000. 2
- [16] Dirk Helbing and Peter Molnar. Social force model for pedestrian dynamics. *Physical review E*, 51(5):4282, 1995. 1, 2
- [17] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 3
- [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 3, 4
- [19] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 2, 3, 4
- [20] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022. 2, 3, 4, 5
- [21] Daniel Holden, Taku Komura, and Jun Saito. Phase-functioned neural networks for character control. *ACM Transactions on Graphics (TOG)*, 36(4):1–13, 2017. 2, 3
- [22] Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. *International Conference on Machine Learning (ICML)*, 2022. 2, 3, 4, 5, 6, 7, 12
- [23] Ioannis Karamouzas, Nick Sohre, Ran Hu, and Stephen J Guy. Crowd space: a predictive crowd analysis technique. *ACM Transactions on Graphics (TOG)*, 37(6):1–14, 2018. 2
- [24] Jongmin Kim, Yeongho Seol, Taesoo Kwon, and Jehee Lee. Interactive manipulation of large-scale crowd animation. *ACM Transactions on Graphics (TOG)*, 33(4):1–10, 2014. 1, 2
- [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 13, 18
- [26] Jaedong Lee, Jungdam Won, and Jehee Lee. Crowd simulation by deep reinforcement learning. In *Proceedings of the 11th Annual International Conference on Motion, Interaction, and Games*, pages 1–7, 2018. 2
- [27] 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
- [28] Marilena Lemonari, Rafael Blanco, Panayiotis Charalambous, Nuria Pelechano, Marios Avraamides, Julien Pettré, and Yiorgos Chrysanthou. Authoring virtual crowds: A survey. In *Computer Graphics Forum*, volume 41, pages 677–701. Wiley Online Library, 2022. 2

- [29] Alon Lerner, Yiorgos Chrysanthou, and Dani Lischinski. Crowds by example. In *Computer graphics forum*, volume 26, pages 655–664. Wiley Online Library, 2007. [2](#), [4](#), [6](#), [17](#)
- [30] Hung Yu Ling, Fabio Zinno, George Cheng, and Michiel Van De Panne. Character controllers using motion vaes. *ACM Transactions on Graphics (TOG)*, 39(4):40–1, 2020. [2](#), [3](#)
- [31] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Smpl: a skinned multi-person linear model. *ACM Trans. Graph.*, 34:248:1–248:16, 2015. [5](#)
- [32] Pablo Alvarez Lopez, Michael Behrisch, Laura Bieker-Walz, Jakob Erdmann, Yun-Pang Flötteröd, Robert Hilbrich, Leonhard Lücken, Johannes Rummel, Peter Wagner, and Eva-marie Wießner. Microscopic traffic simulation using sumo. In *2018 21st international conference on intelligent transportation systems (ITSC)*, pages 2575–2582. IEEE, 2018. [1](#), [2](#)
- [33] Zhengyi Luo, Ryo Hachiuma, Ye Yuan, and Kris Kitani. Dynamics-regulated kinematic policy for egocentric pose estimation. In *Advances in Neural Information Processing Systems*, 2021. [5](#)
- [34] Zhengyi Luo, Shun Iwase, Ye Yuan, and Kris Kitani. Embodied scene-aware human pose estimation. In *Advances in Neural Information Processing Systems*, 2022. [5](#)
- [35] Naureen Mahmood, N. Ghorbani, N. Troje, Gerard Pons-Moll, and Michael J. Black. Amass: Archive of motion capture as surface shapes. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5441–5450, 2019. [5](#), [6](#), [8](#), [16](#)
- [36] Viktor Makovychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, and Gavriel State. Isaac gym: High performance gpu-based physics simulation for robot learning, 2021. [6](#)
- [37] Karttikayaa 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. [1](#), [2](#)
- [38] Chenlin Meng, Ruiqi Gao, Diederik P Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. *arXiv preprint arXiv:2210.03142*, 2022. [8](#), [21](#)
- [39] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. [3](#), [4](#), [12](#)
- [40] Andreas Panayiotou, Theodoros Kyriakou, Marilena Lemonari, Yiorgos Chrysanthou, and Panayiotis Charalambous. Ccp: Configurable crowd profiles. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022. [2](#)
- [41] 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. [2](#), [4](#), [6](#), [17](#)
- [42] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Trans. Graph.*, 37(4):143:1–143:14, July 2018. [16](#), [17](#)
- [43] Xue Bin Peng, Glen Berseth, KangKang Yin, and Michiel van de Panne. Deeploco: Dynamic locomotion skills using hierarchical deep reinforcement learning. *ACM Transactions on Graphics (Proc. SIGGRAPH 2017)*, 36(4), 2017. [2](#), [3](#)
- [44] Xue Bin Peng, Yunrong Guo, Lina Halper, Sergey Levine, and Sanja Fidler. Ase: Large-scale reusable adversarial skill embeddings for physically simulated characters. *ACM Trans. Graph.*, 41(4), July 2022. [2](#), [3](#)
- [45] Xue Bin Peng, Ze Ma, Pieter Abbeel, Sergey Levine, and Angjoo Kanazawa. Amp: Adversarial motion priors for stylized physics-based character control. *ACM Trans. Graph.*, 40(4), July 2021. [2](#), [3](#), [5](#), [16](#)
- [46] Maria Priisalu, Ciprian Paduraru, Aleksi Pirinen, and Cris-tian Sminchisescu. Semantic synthesis of pedestrian locomotion. In *Proceedings of the Asian Conference on Computer Vision*, 2020. [2](#), [3](#)
- [47] Maria Priisalu, Aleksi Pirinen, Ciprian Paduraru, and Cris-tian Sminchisescu. Generating scenarios with diverse pedestrian behaviors for autonomous vehicle testing. In *Conference on Robot Learning*, pages 1247–1258. PMLR, 2022. [2](#)
- [48] Davis Rempe, Tolga Birdal, Aaron Hertzmann, Jimei Yang, Srinath Sridhar, and Leonidas J. Guibas. Humor: 3d human motion model for robust pose estimation. In *International Conference on Computer Vision (ICCV)*, 2021. [3](#)
- [49] Davis Rempe, Jonah Philion, Leonidas J. Guibas, Sanja Fidler, and Or Litany. Generating useful accident-prone driving scenarios via a learned traffic prior. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [1](#), [2](#), [6](#), [7](#), [14](#), [18](#), [20](#)
- [50] Zhiguo Ren, Panayiotis Charalambous, Julien Bruneau, Qunsheng Peng, and Julien Pettré. Group modeling: A unified velocity-based approach. In *Computer Graphics Forum*, volume 36, pages 45–56. Wiley Online Library, 2017. [1](#), [2](#)
- [51] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. [13](#)
- [52] 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](#)
- [53] Nikita Rudin, David Hoeller, Philipp Reist, and Marco Hutter. Learning to walk in minutes using massively parallel deep reinforcement learning, 2021. [5](#), [16](#)
- [54] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamvar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. [21](#)
- [55] Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and Marco Pavone. Trajectron++: Dynamically-feasible trajec-

- tory forecasting with heterogeneous data. In *European Conference on Computer Vision*, pages 683–700. Springer, 2020. [1](#), [2](#)
- [56] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *ArXiv*, abs/1707.06347, 2017. [5](#)
- [57] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015. [3](#)
- [58] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019. [4](#)
- [59] Simon Suo, Sebastian Regalado, Sergio Casas, and Raquel Urtasun. Trafficsim: Learning to simulate realistic multi-agent behaviors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10400–10409, 2021. [14](#)
- [60] Martin Treiber, Ansgar Hennecke, and Dirk Helbing. Congested traffic states in empirical observations and microscopic simulations. *Physical review E*, 62(2):1805, 2000. [1](#), [2](#)
- [61] Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. *Advances in Neural Information Processing Systems*, 34:11287–11302, 2021. [3](#)
- [62] Jur Van den Berg, Ming Lin, and Dinesh Manocha. Reciprocal velocity obstacles for real-time multi-agent navigation. In *2008 IEEE international conference on robotics and automation*, pages 1928–1935. Ieee, 2008. [2](#)
- [63] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011. [4](#)
- [64] Jingkang Wang, Ava Pun, James Tu, Sivabalan Manivasagam, Abbas Sadat, Sergio Casas, Mengye Ren, and Raquel Urtasun. Advsim: Generating safety-critical scenarios for self-driving vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9909–9918, 2021. [18](#)
- [65] Jungdam Won, Deepak Gopinath, and Jessica Hodgins. Physics-based character controllers using conditional vaes. *ACM Transactions on Graphics (TOG)*, 41(4):1–12, 2022. [2](#), [3](#)
- [66] Danfei Xu, Yuxiao Chen, Boris Ivanovic, and Marco Pavone. Bits: Bi-level imitation for traffic simulation. *arXiv preprint arXiv:2208.12403*, 2022. [1](#), [2](#), [6](#), [7](#)
- [67] Wenhao Yu, Greg Turk, and C. Karen Liu. Learning symmetric and low-energy locomotion. *ACM Transactions on Graphics (TOG)*, 37:1 – 12, 2018. [5](#), [16](#)
- [68] Ye Yuan, Shih-En Wei, Tomas Simon, Kris Kitani, and Jason Saragih. Simpoe: Simulated character control for 3d human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. [5](#)
- [69] Ye Yuan, Xinshuo Weng, Yanglan Ou, and Kris M Kitani. Agentformer: Agent-aware transformers for socio-temporal multi-agent forecasting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9813–9823, 2021. [1](#), [2](#), [19](#)
- [70] Xiaohui Zeng, Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, and Karsten Kreis. Lion: Latent point diffusion models for 3d shape generation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. [3](#)
- [71] Yan Zhang and Siyu Tang. The wanderings of odysseus in 3d scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20481–20491, 2022. [3](#)
- [72] Ziyuan Zhong, Davis Rempe, Danfei Xu, Yuxiao Chen, Sushant Veer, Tong Che, Baishakhi Ray, and Marco Pavone. Guided conditional diffusion for controllable traffic simulation. *International Conference on Robotics and Automation (ICRA)*, 2023. [2](#), [3](#), [4](#), [5](#), [7](#), [12](#), [19](#)

Appendices

Appendix A and Appendix B go over details of the TRACE and PACER models, respectively. Appendix C provides additional details of the experiments presented in the main paper, while Appendix D gives additional results to supplement those in the main paper. Appendix E discusses limitations and future work in more detail.

Video Results. Extensive video results are included on the [project page](#). We highly encourage readers to view them to better understand our method’s capabilities.

A. TRACE Details

In this section, we provide details on the TRAjectory Diffusion Model for Controllable PEdestrians (TRACE) presented in Sec 3.1 of the main paper.

A.1. Model Details

A.1.1 Denoising-Diffusion Formulation

Input Representations. In practice, the history trajectories of the ego pedestrian $\mathbf{x}^{\text{ego}} = [\mathbf{s}_{t-T_p}, \dots, \mathbf{s}_t]$ and N neighboring pedestrians $X^{\text{neigh}} = \{\mathbf{x}^i\}_{i=1}^N$ given as input to the diffusion model include more than just positions, heading, and speed. In particular, each past state is

$$\mathbf{s} = [x \ y \ h_x \ h_y \ v \ l \ w \ p] \in \mathbb{R}^8$$

where (x, y) is the 2D position, (h_x, h_y) is the 2D heading vector computed from the heading angle θ , v is the speed, (l, w) is the 2D bounding box dimensions of the person, and $p \in \{0, 1\}$ indicates whether the person is present (visible) at that timestep or not (*e.g.*, due to occlusions in real-world data). If a person is not visible at some step, *i.e.*, $p = 0$, then the full state vector is zeroed out before being given to the diffusion model. All trajectories are transformed into the local frame of the ego pedestrian at the current time step.

The rasterized map input $\mathcal{M} \in \mathbb{R}^{H \times W \times C}$ is in bird’s-eye view, and is cropped around the ego pedestrian and transformed into their local frame. For all experiments $H = W = 224$ px at a resolution of 12 px/m. The map is cropped such that 14 m to the front, left, and right of the ego are visible, and ~ 4.6 m behind. Each channel of \mathcal{M} is a binary map with 1 indicating the presence of some semantic property. For example, in the ORCA dataset, there are only two layers – one for walkable area and one for obstacles. In nuScenes, there are seven layers representing lane, road segment, drivable area, road divider, lane divider, crosswalk, and sidewalk. Notice that the map for TRACE does not contain fine-grained height information as in the map for PACER. As such, TRACE is in charge of high-level obstacle avoidance while PACER factors in both obstacles and terrain.

Denoising with Dynamics. As discussed in the main paper, during denoising the future state trajectory is always a result of actions, *i.e.* diffusion/denoising are on τ_a , similar to [72]. In detail, given an input noisy action sequence τ_a^k the denoising process is as follows: (1) compute the input state sequence $\tau_s^k = f(\mathbf{s}_t, \tau_a^k)$ using the dynamics model f , (2) pass the full input trajectory $\tau^k = [\tau_s^k; \tau_a^k]$ to the denoising model to predict the clean action trajectory $\hat{\tau}_a^0$, (3) compute the output state trajectory $\hat{\tau}_s^0 = f(\mathbf{s}_t, \hat{\tau}_a^0)$, (4) if training, compute the loss in Eqn 3 of the main paper on the full output clean trajectory $\hat{\tau}^0 = [\hat{\tau}_s^0; \hat{\tau}_a^0]$.

We use a unicycle dynamics model for f [1]. Though humans are in theory more agile than the unicycle model, we find it regularizes predictions to be generally smooth, which is how pedestrians usually move and is amenable to being followed by an animation model. Since our model requires actions as input, we compute these from the state-only input data through a simple inverse dynamics procedure.

Parameterization. At each denoising step k , TRACE must predict the mean of the distribution used to sample the slightly less noisy trajectory for step $k - 1$:

$$\mu_\phi(\tau^{k-1} | \tau^k, C) := \mathcal{N}(\tau^{k-1}; \mu_\phi(\tau^k, k, C), \Sigma_k). \quad (8)$$

There are three common ways to parameterize this prediction (we recommend [39] for a full background on these formulations): (1) directly output μ from the network, (2) output the denoised clean trajectory τ^0 , or (3) output the noise ϵ used to corrupt the clean trajectory. TRACE uses (2), but the formulations are equivalent. In particular, we can compute μ from τ^k and τ^0 using

$$\mu(\tau^0, \tau^k) := \frac{\sqrt{\bar{\alpha}_{k-1}}\beta_k}{1 - \bar{\alpha}_k}\tau^0 + \frac{\sqrt{\alpha_k}(1 - \bar{\alpha}_{k-1})}{1 - \bar{\alpha}_k}\tau^k \quad (9)$$

where β_k is the variance from the schedule (we follow [22, 39] and use a cosine schedule), $\alpha_k := 1 - \beta_k$, and $\bar{\alpha}_k := \prod_{j=0}^k \alpha_j$. Therefore, we can plug the output from TRACE $\hat{\tau}^0$ along with the noisy input τ^k into Eq. (9) to get the desired next step mean μ_ϕ . We can also use the fact that τ^0 is corrupted by

$$\tau^k = \sqrt{\bar{\alpha}_k}\tau^0 + \sqrt{1 - \bar{\alpha}_k}\epsilon \quad (10)$$

where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ to compute ϵ from the output of TRACE:

$$\epsilon = \frac{\tau^k - \sqrt{\bar{\alpha}_k}\tau^0}{\sqrt{1 - \bar{\alpha}_k}}. \quad (11)$$

This allows the use of the classifier-free sampling strategy defined in Eqn 4 of the main paper, which requires mixing ϵ outputs from the conditional and unconditional models.

A.1.2 Architecture

The denoising architecture is shown in Fig. 2 of the main paper. At each denoising step k , the step index is processed

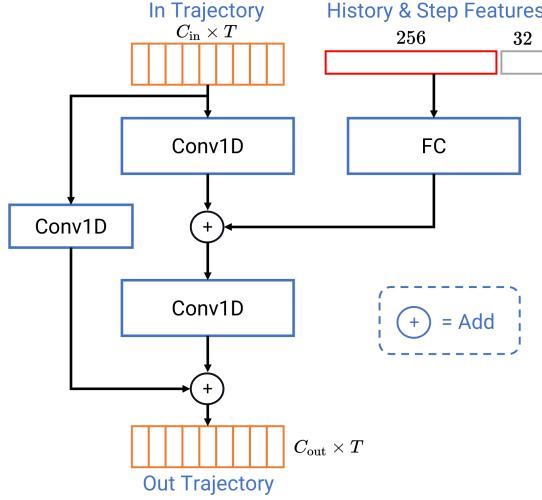


Figure 7. Architecture of a single layer of denoising 1D U-Net.

with a positional embedding followed by a small MLP that gives a 32-dim feature. The map feature extractor uses the convolutional backbone of ResNet-18 [14] as the encoder followed by a 2D U-Net [51] decoder that leverages skip connections to layers in the encoder. For all experiments, the output feature grid Ψ is then $56 \times 56 \times 32$.

The ego x^{ego} and neighbor X^{neigh} history encoders operate on past trajectories $\in \mathbb{R}^{T_p \times 8}$ that are flattened to be a single input vector. The ego and all neighbor trajectories are processed by an MLP with 4 hidden layers giving a feature vector of size 128. A different MLP is learned for ego and neighbor trajectories (*i.e.* all neighbors are processed by the same MLP, which is different from the ego MLP). Neighbor trajectory features are max-pooled to get a single *interaction* feature. The resulting ego and interaction features are finally jointly processed by another MLP with 4 hidden layers, giving a 256-dim feature summarizing the past trajectory context.

Note that the processing of input conditioning described thus far is only necessary to do *once* before starting the denoising process. Only the denoising 1D U-Net needs to be run at every step. At step k of denoising, the 2D position at each timestep $t+i$ of the current noisy input trajectory τ^k is queried in the map feature grid to obtain a feature $\mathbf{g}_{t+i} = \Psi(x_{t+i}, y_{t+i}) \in \mathbb{R}^{32}$. This query is done through bilinear interpolation of map features at the corresponding point. Over all timesteps, these form a feature trajectory $\mathbf{G} = [\mathbf{g}_{t+1} \dots \mathbf{g}_{t+T_f}]$ that is concatenated along the channel dimension with $\tau^k \in \mathbb{R}^{T_f \times 6}$ (containing both actions and states) to get the full trajectory input to the denoising U-Net $[\tau^k; \mathbf{G}] \in \mathbb{R}^{T_f \times 38}$. Each layer of the U-Net also receives the concatenation of the past trajectory context and denoising step feature.

The architecture of a single U-Net layer is shown in Fig. 7. The input trajectory at each layer is first processed by

a 1D convolution. The input trajectory history and step index feature are projected to the same feature size, broadcast over the temporal dimension, and then added to the intermediate trajectory features. Another convolution is performed before adding to the input trajectory in a residual fashion. In the encoding part of the U-Net shown in Fig. 2 of the main paper, a $2 \times$ downsampling over the temporal dimension is performed between layers, while a $2 \times$ upsampling is done in the decoding part. The encoder is three layers with output channels being 64, 128, and 256-dim.

A.1.3 Training Details

TRACE uses $K = 100$ denoising steps in both training and testing. Training uses a fixed learning rate of 2e-4 with the Adam optimizer [25] and runs for 40k iterations. TRACE trains on a 32 GB NVIDIA V100 GPU and takes ~ 2 days on the ORCA dataset and ~ 3 days on the mixed nuScenes+ETH/UCY data.

During training, the neighbor history and map conditioning are randomly dropped out with a 10% probability. Note that these are dropped independently and that the ego history is never dropped. In practice, to drop map conditioning, all pixels are filled with a 0.5 value; this means that the model is aware that it “does not know” about the map context, it *is not* simply fed an empty map with no obstacles (all zeros). To drop neighbor conditioning, the neighbor history feature is zeroed out. The same mechanism is used to train on “mixed” data with varying annotations, *e.g.*, some data samples have no maps. In this case, a map is still given to the model but filled with 0.5 value pixels.

A.2 Guidance Details

A.2.1 Scene-Level Guidance

Some guidance objectives are based on multi-agent interactions, *e.g.*, agent collision avoidance and social groups. In this case, we assume that all pedestrians in a scene can be denoised simultaneously in a batched fashion. At each denoising step, the loss function is evaluated at the current trajectory prediction of all pedestrians and gradients are propagated back to each one for guidance. This can be seen as sampling a scene-level future rather than a single agent future. If we want to sample M possible scene futures, we can draw M samples from each agent and assume that the m th sample from each agent corresponds to the same scene sample. In other words, we compute the scene-level guidance by considering only the m th sample from each agent.

This multi-agent guidance slightly complicates the *filtering* procedure described in Sec. 4 of the main paper whereby trajectory samples are strategically chosen to minimize the guidance loss. In the case of a multi-agent objective, the trajectory that minimizes the guidance loss for one

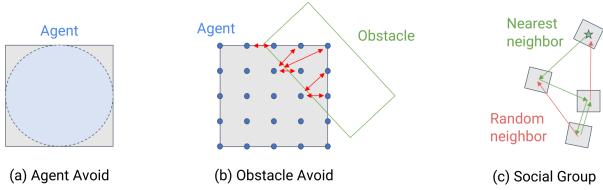


Figure 8. Illustrations of various guidance objectives. See text for details.

agent may not be globally optimal, so it is undesirable to filter the agents independently. We instead do filtering at the scene level, similar to how guidance is computed: we compute the summed guidance loss across the m th sample from all agents and choose the *scene-level sample* that minimizes this aggregate loss.

A.2.2 Guidance Objectives

Next, we describe the different test-time guidance objectives (losses) that we have implemented for TRACE. Objectives operate on the future state trajectory τ_s that starts at timestep $t + 1$ and contains state s_j at time j .

Agent Avoid and Social Distance. We use the same agent collision penalty as in TrafficSim [59] and STRIVE [49] which approximates each agent with disks to efficiently and differentiably compute a collision loss. For pedestrians, it is sufficient to use a single disk for each agent (see Fig. 8(a)). With this approximation, collision detection is fast and the collision loss is computed based on the extent of disk overlap between agents. It is easy to artificially inflate the size of the disk in order to implement a desired social distance between pedestrians. Note that this is a multi-agent objective, so guidance is enforced at the scene level, as previously discussed.

Obstacle Avoid. We extend the differentiable environment collision penalty introduced in STRIVE [49] to more robustly handle collision avoidance and provide more useful gradients. The core idea is illustrated in Fig. 8(b); for each timestep where the pedestrian’s bounding box is overlapping with an obstacle, we query a grid of points on the agent (in our experiments, this is 10×10) and define a loss with respect to points that are embedded in the obstacle. For each embedded point, we compute the minimum distance d_{\min} to a non-embedded point on the agent and define the loss at that point as $\mathcal{L} = 1 - (d_{\min}/b)$ where b is the bounding box diagonal of the agent. Summing the loss at all embedded points gives the total loss.

A subtle, but very important, implementation detail here is that the embedded points must be detached (*i.e.* stop grad) before computing the loss. Intuitively, embedded points are treated as points on the obstacle, *not* the agent. So when the loss is computed with respect to these points, it gives a

meaningful gradient back to the non-embedded agent points which propagates back to the agent position and heading.

Local Waypoint at Specific Time. This loss encourages an agent to be at a specific 2D goal waypoint $\mathbf{p}_g = (x, y)$ at a specific time step j that falls *within* the planning horizon T_f of the model. It simply encourages the trajectory to be at that location at the timestep with the loss $\mathcal{L} = \|\mathbf{s}_j - \mathbf{p}_g\|_2$.

Local Waypoint at Any Time. This loss encourages an agent to be at a specific 2D goal waypoint $\mathbf{p}_g = (x, y)$ at *any* timestep that is *within* the planning horizon T_f of the model. The loss is defined as

$$\mathcal{L} = \sum_{j=t+1}^{t+T_f} \delta_j \cdot \|\mathbf{s}_j - \mathbf{p}_g\|_2^2 \quad (12)$$

$$\delta_j = \frac{\exp(-\|\mathbf{s}_j - \mathbf{p}_g\|)}{\sum_j \exp(-\|\mathbf{s}_j - \mathbf{p}_g\|)}. \quad (13)$$

Intuitively, it tries to minimize the distance from all points in the trajectory to the target location, but each timestep is weighted by δ_j which is the *softmax* over the distances of each step from the waypoint. The δ_j form a distribution over trajectory timesteps where states close to the waypoint will have higher probability and therefore be weighted more in the loss.

Global Waypoint at Specific Time. This loss encourages an agent to be at a specific 2D goal waypoint $\mathbf{p}_g = (x, y)$ at a specific timestep j that falls *beyond* the planning horizon T_f of the model. This is useful during closed-loop operation in which the agent should eventually reach the point, but at the current step t is not within the planning horizon. Intuitively, the loss encourages making enough progress toward the waypoint such that when it becomes in range, we can revert to the Local Waypoint loss and hit the target exactly at the desired time.

To do this, we would like to ensure that the future trajectory ends in a location such that the pedestrian can travel in a straight line at a “preferred” speed v_{pref} and get to the waypoint on time. Formally, the trajectory should be within a target distance defined as:

$$d_{\text{goal}} = (j - t) \cdot dt \cdot v_{\text{pref}} \quad (14)$$

where dt is the step size of TRACE output (0.1 sec in our experiments). Since the pedestrian may not be able to actually travel a straight line path (*e.g.* in environments with obstacles), we incorporate an *urgency* parameter $u \in [0, 1]$ that encourages getting there earlier and modifies the goal distance as

$$\tilde{d}_{\text{goal}} = d_{\text{goal}} \cdot (1 - u). \quad (15)$$

Then the loss with respect to this target distance is defined as

$$\mathcal{L} = \text{ReLU}(\|\mathbf{s}_{T_f} - \mathbf{p}_g\|_2 - \tilde{d}_{\text{goal}}) \quad (16)$$

which penalizes the trajectory if the final state is not within the goal distance.

Global Waypoint at Any Time. This loss encourages an agent to be at a specific 2D goal waypoint $\mathbf{p}_g = (x, y)$ at *any* timestep *beyond* the planning horizon T_f of the model. To determine whether the goal waypoint is outside the current horizon, we check if the agent could progress along a straight line to the goal at a preferred speed v_{pref} and reach the waypoint within the planning horizon. If so, the loss reverts to the Local Waypoint loss.

If the waypoint is indeed beyond the planning horizon, the loss attempts to progress according to some urgency $u \in [0, 1]$. To do this, we first compute the maximum distance that could be covered in the current horizon:

$$d_{\max} = T_f \cdot dt \cdot v_{\text{pref}} \quad (17)$$

and use the urgency to get the goal distance we wish to cover over the horizon

$$d_{\text{goal}} = u \cdot d_{\max}. \quad (18)$$

The loss is computed based on how much progress is made over the horizon:

$$\mathcal{L} = \text{ReLU}(d_{\text{goal}} - d_{\text{progress}}) \quad (19)$$

$$d_{\text{progress}} = \|\mathbf{s}_t - \mathbf{p}_g\|_2 - \|\mathbf{s}_{t+T_f} - \mathbf{p}_g\|_2. \quad (20)$$

Social Groups. This loss encourages groups of agents to travel together. A social group is based on one *leader* pedestrian that is not affected by the social group loss (via detach/stop grad); others in the group will tend to move with the leader. Intuitively, we want each agent in the social group to maintain a specified social distance d_{soc} to the closest agent also in the same social group. Let ψ be a map from one agent index to another, *e.g.* $i = \psi(k)$ means agent k in the social group is mapped to agent i . Then the social group loss for agent i at timestep j in the future trajectory is

$$\mathcal{L} = \left(\|\mathbf{s}_j^i - \mathbf{s}_j^{\psi(i)}\|_2 - d_{\text{soc}} \right)^2. \quad (21)$$

As shown in Fig. 8, most of the time ψ maps each agent to the closest agent in the group, but with some probability based on a cohesion parameter $c \in [0, 1]$ the mapping will be to a random agent in the group. So with a larger cohesion, agents in the group are all encouraged to be equidistant from each other, while with low cohesion agents will not closely follow the leader and connected components in the social group graph may break off and ignore others.

Learned Value Function. Given a learned value function $V(\tau_s)$ that predicts the future rewards over a given trajectory, in general this guidance loss is simply $\mathcal{L} = \exp(-V(\tau_s))$. When TRACE is used with PACER, the value function is $V(v_t | \tau_s, \mathbf{h}_t, \mathbf{o}_t, \beta)$; it takes in the current humanoid state \mathbf{h}_t , environmental feature \mathbf{o}_t , and humanoid shape β , which are fixed throughout denoising.

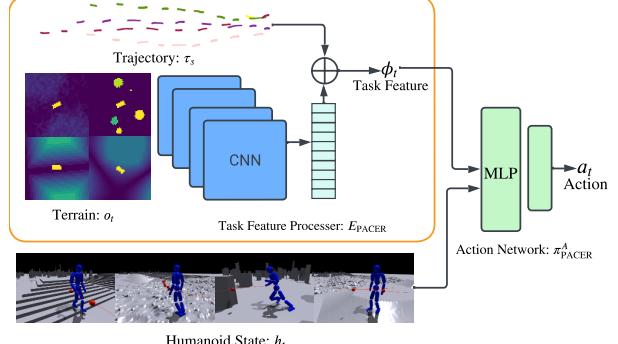


Figure 9. The PACER policy network π_{PACER} consists of a task feature processor $E_{\text{PACER}}(\phi_t | \mathbf{o}_t, \tau_s)$ and an action policy network $\pi_{\text{PACER}}^A(\phi_t, \mathbf{h}_t)$.

B. PACER Details

In this section, we give details on the **Pedestrian ControllIEr** (PACER) presented in Sec 3.2.

B.1. Implementation Details

Humanoid State. The state \mathbf{h}_t holds joint positions $\mathbf{j}^t \in \mathbb{R}^{24 \times 3}$, rotations $\mathbf{q}^t \in \mathbb{R}^{24 \times 6}$, linear velocities $\mathbf{v}^t \in \mathbb{R}^{24 \times 3}$, and angular velocities $\omega^t \in \mathbb{R}^{24 \times 3}$ all normalized w.r.t. the agent’s heading and root position. The rotation is represented in the 6-degree-of-freedom rotation representation. SMPL has 24 body joints with the root (pelvis) as the first joint, which is not actuated, resulting in an action dimension of $\mathbf{a}_t \in \mathbb{R}^{23 \times 3}$. No special root forces/torques are used.

Network Architecture. As mentioned in the main paper, the environmental feature \mathbf{o}_t is a rasterized local height and velocity map of size $\mathbf{o}_t \in \mathbb{R}^{64 \times 64 \times 3}$. The first channel is the terrain height map relative to the current humanoid root height, and the second & third channels are the 2D linear velocities (x and y directions) in the egocentric coordinate system. The map corresponds to a $4\text{m} \times 4\text{m}$ square area centered at the humanoid root, sampled on an evenly spaced grid. The trajectory $\tau_s \in \mathbb{R}^{10 \times 2}$ consists of the 2D waypoints for the next 5 seconds sampled at 0.5 s intervals.

The architecture of π_{PACER} can be found in Fig. 9. Due to the high dimensionality of the environmental features \mathbf{o}_t , we separate the policy network into a task feature processor $E_{\text{PACER}}(\phi_t | \mathbf{o}_t, \tau_s)$ and an action network $\pi_{\text{PACER}}^A(\mathbf{a}_t | \phi_t, \mathbf{h}_t, \beta)$. The task feature processor transforms task-related features, such as environmental features \mathbf{o}_t and trajectory τ_s into a latent vector $\phi_t \in \mathbb{R}^{256}$. Then, π_{PACER}^A computes the action \mathbf{a}_t based on the humanoid state \mathbf{h}_t , body shape β , and ϕ_t . The overall policy network is then $\pi_{\text{PACER}}(\mathbf{a}_t | \mathbf{o}_t, \mathbf{h}_t, \beta, \tau_s) \triangleq \pi_{\text{PACER}}^A(E_{\text{PACER}}(\mathbf{o}_t, \tau_s), \mathbf{h}_t, \beta)$. E_{PACER} is a four-level convolutional neural network with a stride of 2, 16 filters, and a kernel size of 4. π_{PACER}^A is a standard MLP with ReLU activations. It has two layers, each with 2048 and 1024 units.

The policy maps to the Gaussian distribution over actions $\pi_{\text{PACER}}(\mathbf{a}_t | \mathbf{o}_t, \mathbf{h}_t, \boldsymbol{\beta}, \boldsymbol{\tau}_s) = \mathcal{N}(\mu(\mathbf{o}_t, \mathbf{h}_t, \boldsymbol{\beta}, \boldsymbol{\tau}_s), \Sigma)$ with a fixed covariance matrix Σ . Each action vector $\mathbf{a}_t \in \mathbb{R}^{23 \times 3}$ corresponds to the PD targets for the 23 actuated joints on the SMPL human body. The discriminator $D(\mathbf{h}_{t-10:t}, \mathbf{a}_t)$ shares the same architecture as π_{PACER}^A , while the value function $V(v_t | \mathbf{o}_t, \mathbf{h}_t, \boldsymbol{\beta}, \boldsymbol{\tau}_s)$ shares the same architecture as the policy π_{PACER} .

B.2. Reward and Loss

Reward. Following AMP [45], our policy π_{PACER} is learned through goal-conditioned RL where the reward contains a task reward r_t^τ , a style reward r_t^{amp} , and an energy penalty r_t^{energy} . The style reward is computed by the discriminator $D(\mathbf{h}_{t-10:t}, \mathbf{a}_t)$ based on 10 steps of aggregated humanoid state. We use the same set of observations, loss formulation, and gradient penalty as in AMP [45] to train our discriminator. Task reward r_t^τ is a trajectory-following reward that measures how far away the humanoid’s center \mathbf{c}_t on the xy plane is from the 2D trajectory: $\exp(-2 \times \|\mathbf{c}_t - \boldsymbol{\tau}_t\|^2)$. The energy penalty is expressed as $-0.0005 \cdot \sum_{j \in \text{joints}} |\boldsymbol{\mu}_j \dot{\mathbf{q}}_j|^2$ where $\boldsymbol{\mu}_j$ and $\dot{\mathbf{q}}_j$ correspond to the joint torque and the joint angular velocity, respectively.

Motion Symmetry Loss. During our experiments, we noticed that asymmetric gaits emerge as training progresses. It manifests itself as “limping” where the humanoid produces asymmetric motion, especially at a lower speed. This could be due to the small temporal window used in AMP (10 frames), which is not sufficient to generate symmetrical motion. Compared to AMP, we use a humanoid with more than double the degrees of freedom (69 vs 28), and the complexity of the control problem grows exponentially. This could also contribute to limping behavior as it becomes harder for the discriminator to discern asymmetric gaits. Thus, we utilize the motion symmetry loss proposed in [67] to ensure symmetric gaits. Specifically, we first design two functions Φ_s and Φ_a that can mirror the humanoid state and action along the character’s sagittal plane. Symmetry is then enforced by ensuring that the mirrored states lead to mirrored actions:

$$L_{\text{sym}}(\theta) = \|\pi_{\text{PACER}}(\mathbf{h}_t, \mathbf{o}_t, \boldsymbol{\beta}, \boldsymbol{\tau}_s) - \Phi_a(\pi_{\text{PACER}}(\Phi_s(\mathbf{h}_t, \mathbf{o}_t, \boldsymbol{\beta}, \boldsymbol{\tau}_s)))\|^2, \quad (22)$$

Notice that the motion symmetry loss is not a reward and is directly defined on the policy output. As the loss can be computed in an end-to-end differentiable fashion, we directly optimize this loss through SGD.

B.3. Training

Our training procedures closely follow AMP [45], with notable distinctions in the motion dataset, initialization, termination condition, terrain, and humanoids used. Training takes ~ 3 days to converge on one NVIDIA RTX 3090.

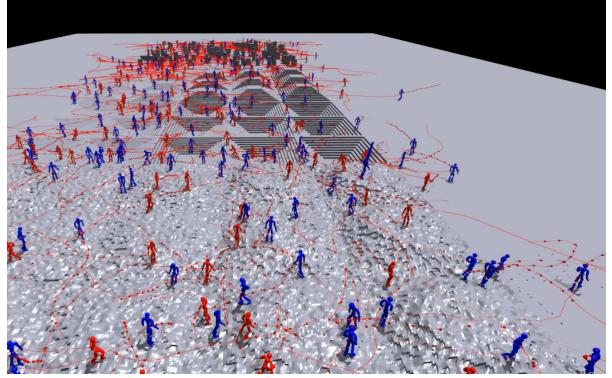


Figure 10. During training, 2048 humanoids are simulated in parallel on our synthetic terrain.

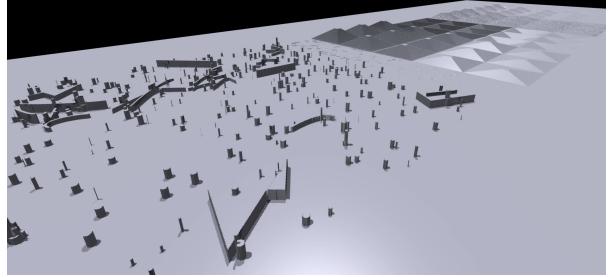


Figure 11. Synthetic terrains used for training PACER. From left to right: obstacles, discrete terrains, stairs (up), stairs (down), uneven terrains, and slopes.

Dataset. We use a small subset of motion sequences from the AMASS dataset [35] to train our humanoid controller. Specifically, we hand-picked ~ 200 locomotion sequences consisting of walking and turning at various speeds, as well as walking up and down stairs. These motions form the reference motion database and provide our AMP Discriminator $D(\mathbf{h}_t, \mathbf{a}_t)$ with “real” samples.

Initialization. To initialize our humanoids during training, we use reference state initialization [42] to randomly sample a body state \mathbf{h}_0 . The initial root positions are randomly sampled from a “walkable map” that corresponds to all locations that can be used as a valid starting point (e.g. not on top of obstacles). As we use NVIDIA’s Isaac Gym, we create 2048 humanoids that are simulated simultaneously in parallel during training: see Fig. 10.

Random terrain, trajectory, and body shape sampling. To learn a model that can traverse diverse types of terrain that pedestrians may encounter in real life, we train our trajectory-following controller on a variety of different environments. Specifically, we follow ANYmal [53] to create terrain curricula with varying difficulties to train our agents. Six types of terrain are created: slopes, uneven terrain, stairs (down), stairs (up), discrete, and obstacles. The terrains follow a gradual increase in difficulty, where we vary the slope angle, terrain unevenness, slope angle of stairs, and obstacle density, as shown in Fig. 11.

Trajectory samples for training are generated procedu-

rally: τ_s is randomly sampled by generating velocities and turn angles. We limit the velocity to be between [0, 3] m/s and the acceleration to be between [0, 2] m/s².

To train with different body shapes, we extract all unique human body shapes from the AMASS dataset, which amount to 476 shapes (273 male and 200 female). We randomly sample (with replacement) 2048 body shapes to create humanoids at the beginning of the training process. To create reference humanoid states \hat{h}_t for the discriminator, we perform forward kinematics based on the sampled pose and the humanoids' kinematic tree. At the beginning of every 250 episodes, we randomly sample a new batch of pose sequences from the motion dataset and create new reference humanoid states. In this way, we obtain reference states of diverse body types and motions.

Termination condition. To speed up training, we employ early termination [42] and terminate the episode if there is a collision force greater than 50 N on the humanoid body, with either the scene or other humanoids. The ankles and foot joints are exceptions to this rule, as they are in contact with the ground. This condition also serves as a fall detection mechanism, as falling will involve a collision force from the ground. Notice that this termination condition encourages the humanoid to avoid obstacles and other humanoids since a collision will trigger an early termination.

C. Experimental Details

In this section, we include details of the experiments presented in Sec 4 of the main paper.

C.1. Dataset Details

The ORCA dataset contains two distinct subsets, *ORCA-Maps* and *ORCA-Interact*. ORCA-Maps is generated with up to 10 pedestrians and 20 obstacles in each scene. This contains many obstacle interactions, but fewer agent-agent interactions. ORCA-Interact has up to 20 pedestrians, but no obstacles, and therefore has no map annotations. Each data subset contains 1000 scenes that are 10s long, and we split them 0.8/0.1/0.1 into train/val/test splits. The map annotations in the ORCA dataset contain two channels, one representing the walkable area and one representing obstacles. The bounding box diameter for every agent is fixed to 0.8m.

In nuScenes [4], there are seven map layers representing the lane, road segment, drivable area, road divider, lane divider, crosswalk, and sidewalk. The bounding box diameters are given by the dataset. We follow the official trajectory forecasting benchmark for scenes in the train/val/test splits. For ETH/UCY [29, 41], we use the official training splits of each contained dataset for training.

All trajectory data in all datasets is re-sampled to 10 Hz for training and evaluation of TRACE.

C.2. Guidance Metrics

Here, we define the *Guidance Error* for each of the objectives evaluated in Sec 4.1 and 4.2 of the main paper.

Obstacle Avoid. This is the obstacle collision rate as defined below.

Agent Avoid. This is the agent collision rate as defined below.

Waypoint and Perturbed Waypoint. If the objective is to reach a waypoint at a specific timestep, this is simply the distance of the agent from the target waypoint at that specified timestep (in meters). Otherwise, if the objective is to reach at *any* timestep, the error is the minimum distance between the agent and the goal waypoint across the entire trajectory.

Social Groups. For each pedestrian in the group, we measure the mean absolute difference between the specified social distance d_{soc} (see Appendix A.2.2) and the distance to the closest neighbor in the same social group.

Multi-Objective. For the multi-objective guidance presented in Tab 1 of the main paper (Waypoint + Avoidance), the reported guidance error is the waypoint error since obstacle and agent collision rates are already reported in other columns.

C.3. Other Metrics

Next, we describe in more detail the metrics used to evaluate the standalone TRACE model in Sec 4.1 and 4.2 of the main paper.

Obstacle Collision Rate. Measures the average fraction of time that an agent (as represented by a single disk) is overlapping with an obstacle on the map within a rollout. Note that a disk is used to represent each agent because this is how they are represented in the ORCA simulator, so the ground truth data contains no collisions using this representation.

Agent Collision Rate. Measures the average fraction of agents involved in an agent-agent collision within each scene rollout. This again uses the disk representation of each agent.

Realism (EMD). Compares the histogram of statistics over the entire test set between generated and ground truth trajectories. This is done for velocity, longitudinal acceleration, and lateral acceleration. In particular, the statistics at each timestep of the test set are aggregated together into a histogram. The histogram is then normalized such that it sums to 1. The earth mover's distance (EMD) between the ground truth and generated histograms is then computed² and reported. Note that this metric is computed wrt *the dataset being evaluated on*. For example, in Sec 4.1 of the

²using pyemd

main paper, even though TRACE is trained on both ORCA-Maps and ORCA-Interact, the metric is only computed for ORCA-Maps since this is the test data.

Realism (Mean). Measures the average longitudinal and lateral acceleration within a generated trajectory in m/s^2 . Similar metrics are commonly used in the vehicle planning literature [64] as a proxy for how *comfortable* a ride is. In our case of pedestrian motion, this is still relevant since people tend to move in smooth motions without sudden changes in speed or direction.

C.4. VAE Baseline Details

We adapt a conditional VAE model similar to the idea of STRIVE [49] for controlling trajectories through latent space optimization. We adapt the VAE design to our setting.

Architecture. The architecture operates in an agent-centric manner as in TRACE. It is a fairly standard conditional VAE (CVAE) where the conditioning (map and past trajectories) is processed into a single conditioning vector \mathbf{c} that is given to the decoder. At training time, the decoder also takes in a latent vector \mathbf{z} from the encoder (posterior), while at test time the latent vector is sampled from the prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$. To make the methods comparable, the map conditioning is encoded with the same ResNet-18 backbone that TRACE uses; ego and neighbor past trajectories also use the same architecture as TRACE. Since the model is agent-centric rather than scene-centric, the decoder D is simply an MLP that maps the conditioning and sampled latent to an output action trajectory (instead of a graph network as in STRIVE) as $\tau_a = D(\mathbf{z}, \mathbf{c})$. In all experiments, the latent dimension is 64, while the conditioning feature vector is 256-dim.

Training. Training is done using a standard VAE loss consisting of a reconstruction and a KL divergence term. The KL term is weighted by 1e-4. The model is trained with the same batch size as TRACE (400) and for the same number of iterations (40k) with a learning rate of 2e-4.

Test-Time Optimization. The idea of test-time optimization is to search for a latent vector that is likely under the prior (*i.e.* represents a plausible future trajectory) but also meets the desired guidance objective. Concretely, the optimization objective is

$$\min_{\mathbf{z}} \alpha \mathcal{J}(D(\mathbf{z}, \mathbf{c})) - \log p(\mathbf{z}) \quad (23)$$

where \mathcal{J} is a guidance loss as described in the main paper and α balances the prior term with the guidance loss. Optimization is performed with Adam [25] using a learning rate of 0.02. For experiments in Sec 4.1 of the main paper, optimization uses 100 iterations (same as the number of diffusion steps K). For Sec 4.2, the iteration budget is increased to 200 to accommodate more difficult out-of-distribution objectives.

Discussion on VAE Comparison. The VAE with test-time optimization is generally a very strong baseline. Given a large enough compute budget, the optimization can usually faithfully meet the desired objective. However, the number of optimization iterations needed to meet an objective can be large; *e.g.* in Sec 4.2 it requires twice the number of diffusion steps, making it slower than TRACE. Moreover, when optimizing for a long time to closely meet objectives, the diversity of optimized samples becomes low as they converge to similar minima (Fig. 4 in the main paper). This is due to the prior term in Eq. (23), which always drives the trajectory towards the mean.

C.5. Additional Experiment Details

Finally, we include miscellaneous details of the setup for each experiment in Sec. 4 of the main paper.

Augmenting Crowd Simulation (Sec 4.1). For the no guidance rows in Tab 1, we actually run the evaluation three times and report the averaged metrics. This is because when there is no guidance, no filtering is performed, so a random sample is chosen. Running with several random samples gives a more faithful evaluation of performance. In this experiment, TRACE uses $w = 0.0$ for classifier-free sampling. 20 samples are drawn and guided from the model for each pedestrian before filtering. The weighting α for each guide (in Eq. 6 of the main paper) is tuned manually to meet objectives while maintaining realistic trajectories. The *Waypoint* guidance used in Tab 1 is the *Local Waypoint at Any Time* introduced in Appendix A.2.2. The agent avoidance guidance uses an additional social distance buffer of 0.2m.

Real-world Data Evaluation (Sec 4.2). In this experiment, 10 samples are drawn from the model before filtering and the waypoint guidance is *Global Waypoint at Any Time* since the model operates in a closed loop for longer than the planning horizon. This waypoint guidance uses an urgency of $u = 0.7$ and a preferred speed of $v_{\text{pref}} = 1.25m/s$. The perturbed waypoint objective randomly perturbs the target ground truth waypoint with Gaussian noise with a standard deviation of 2m. The social group guidance uses $d_{\text{soc}} = 1.5$ and cohesion $c = 0.3$. In each nuScenes scene, social groups are determined heuristically by forming a scene graph where edges are present if two pedestrians are within 3m of each other and moving in a similar direction (velocities have a positive dot product): the connected components of this graph with more than one agent form the social groups.

Controllable Pedestrian Animation (Sec 4.3). In this experiment, 10 samples are drawn from the model before filtering. Waypoint guidance uses *Global Waypoint at Specific Time* with the waypoint randomly placed at a reasonable distance ([7, 12] meters in front of the user and up to 5 meters to either side) 9 sec in the future. Agents are initial-

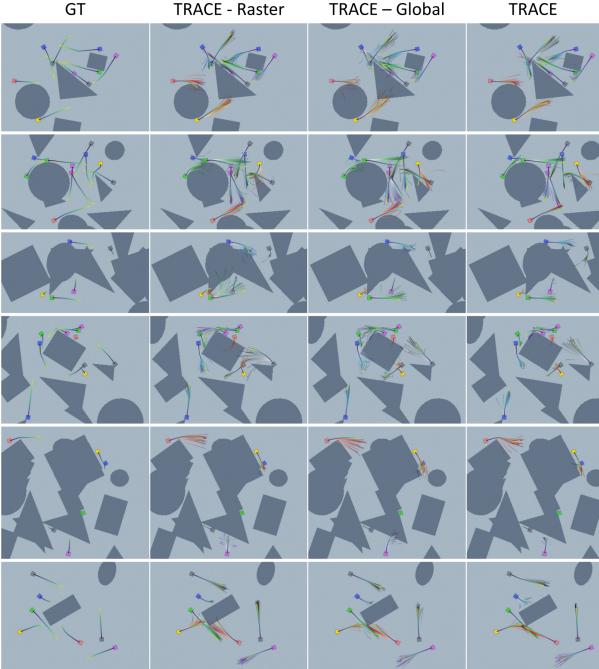


Figure 12. Random sampling with no guidance from different TRACE architecture ablations. Using the learned feature map has apparent benefits in subtle interaction with obstacles. 20 random samples are visualized for each pedestrians with the (arbitrarily) chosen plan in bold.

ized in a standing pose with a uniform random initial root velocity in $[1.2, 2.0] \text{m/s}$.

D. Supplementary Results

In this section, we include additional experimental results omitted from the main paper due to page limits.

D.1. Qualitative Results

Extensive qualitative video results for both TRACE and PACER are provided [on the supplementary webpage](#).

D.2. TRACE– No Guidance Ablation Study

The focus of our work is on controllability using guidance. However, it is still desired that the model performs well even when guidance is not used. In particular, random samples from the model should be robust (avoid collisions), realistic (similar to the training data distribution), and accurate (capture the ground truth future trajectory). To this end, we evaluate our architecture and training approach compared to baselines and ablations while using no guidance. We evaluate in the open-loop 5s rollout setting on ORCA-Maps as used in Sec. 4.1 of the main paper. Two additional metrics common in trajectory forecasting are evaluated: the average displacement error (ADE) and the final displacement error (FDE) [69]. For a trajectory sample defined from timesteps 1 to T , these are defined as

$$\text{ADE} = \frac{1}{T} \sum_{t=1}^T \|\hat{\mathbf{y}}_t - \mathbf{y}_t\|_2 \text{ and } \text{FDE} = \|\hat{\mathbf{y}}_T - \mathbf{y}_T\|_2 \text{ with } \hat{\mathbf{y}} \text{ the sample from the model and } \mathbf{y} \text{ the ground truth.}$$

First, we evaluate the TRACE architecture, which uses a feature grid to condition denoising on the map input. An alternative way to condition trajectory generation is to encode the map into a single (“global”) feature vector using a convolutional backbone. This global feature can then be given in the same way as the past trajectory features. The VAE baseline and *TRACE-Global* ablation do this using a ResNet-18 backbone. The *TRACE-Raster* ablation is similar to CTG [72], which rasterizes both the map *and* agent histories and encodes them into a single global feature instead of encoding the trajectory states separately. Tab. 5 shows the results comparing these methods. In this experiment, we take 20 samples from each model and evaluate the one that is closest to the ground truth wrt the ADE. We see that using the feature grid map provides the lowest obstacle collision rate while maintaining competitive accuracy and realism. As qualitatively shown in Fig. 12, the use of the feature grid gives local cues to the model to inform subtle obstacle interactions and avoid collisions. Though agent collisions are slightly worse with the grid map, no model does particularly well, and exploring improved agent-agent interactions is an important direction for future work.

In the main paper, we discuss how mixed training data and classifier-free sampling (*i.e.* training with random dropping on conditioning) are important to enable flexibility for guidance. To ensure this training approach does not negatively affect base model performance without guidance, we compare to (1) an ablation that uses the ORCA-Maps data only to train (rather than a mix of ORCA-Maps+ORCA-Interact) and (2) ablations that use varying levels of dropping. Tab. 6 shows results, where again the sample closest to ground truth is evaluated. Interestingly, training with mixed data allows for increased accuracy compared to training only on the ORCA-Maps dataset. Increasing the drop probability past 5% has very little effect on performance and comes with the added benefit of using classifier-free sampling to get flexible guidance at test time.

D.3. TRACE– Effect of Classifier-Free Sampling

Next, we examine how weight w affects model performance when using classifier-free sampling both with and without guidance. First, we analyze the effect when evaluating on the ORCA-Maps dataset with no guidance in the open-loop setting (like Sec 4.1 of the main paper). In this case, we evaluate $w \geq 0$ which increases emphasis on the input conditioning to the model. Quantitative results are shown in Tab. 7: for each value of w , the evaluation is run 3 times with different random samples and the metrics are averaged. For $w \in [0, 1]$, the collision rates and acceleration realism remain similar or slightly improve, which we expect since input conditioning such as the obstacle map is empha-

Method	History	Map	Accuracy		Collision Rate		Realism (EMD)		
	Input	Feature	ADE	FDE	Obstacle	Agent	Vel	Lon Acc	Lat Acc
VAE [49]	States	Global	0.340	0.774	0.062	0.115	0.041	0.038	0.039
TRACE-Raster	Raster	Global	0.337	0.808	0.052	0.100	0.027	0.013	0.014
TRACE-Global	States	Global	0.280	0.686	0.056	0.094	0.022	0.013	0.016
TRACE	States	Grid	0.318	0.757	0.046	0.110	0.028	0.020	0.020

Table 5. No guidance evaluation on *ORCA-Maps* dataset. Ablation on architecture design choices.

Train Data	Drop Rate	Accuracy		Collision Rate		Realism (EMD)		
		ADE	FDE	Obstacle	Agent	Vel	Lon Acc	Lat Acc
ORCA-Maps	10%	0.351	0.819	0.040	0.112	0.030	0.023	0.024
Mixed	0%	0.303	0.719	0.042	0.123	0.028	0.019	0.020
Mixed	5%	0.307	0.712	0.040	0.108	0.024	0.020	0.023
Mixed	10%	0.318	0.757	0.046	0.110	0.028	0.020	0.020

Table 6. No guidance evaluation on *ORCA-Maps* dataset. Ablation on training routine.

w	Collision Rate		Realism (EMD)		
	Obstacle	Agent	Vel	Lon Acc	Lat Acc
0.0	0.051	0.131	0.019	0.012	0.014
0.3	0.051	0.130	0.027	0.008	0.010
0.5	0.050	0.132	0.029	0.008	0.009
0.7	0.050	0.132	0.033	0.009	0.008
1.0	0.049	0.130	0.040	0.010	0.009
2.0	0.051	0.132	0.063	0.017	0.015
3.0	0.052	0.138	0.087	0.025	0.022
4.0	0.051	0.145	0.102	0.033	0.028

Table 7. Classifier-free sampling analysis on *ORCA-Maps* dataset with no guidance.

w	Waypoint Error	Realism (Mean)	
		Lon Acc	Lat Acc
0.0	1.129	0.233	0.218
-0.3	0.972	0.213	0.199
-0.5	0.802	0.212	0.204
-0.7	0.670	0.240	0.233
-1.0	0.546	0.345	0.348

Table 8. Classifier-free sampling analysis on nuScenes dataset using perturbed waypoint guidance.

sized. For $w > 1$, the guidance tends to be too strong and the trajectory samples are almost deterministic. Though the quantitative difference is not large as w increases, in Fig. 13 we see that increasing does have a considerable qualitative effect.

Second, we look at results on the nuScenes dataset using perturbed waypoint guidance in the closed-loop setting (the same as in Sec. 4.2 of the main paper). In Sec. 4.2 of the main paper, we saw that using $w < 0$ improves susceptibility to guidance and allows the model to achieve out-of-distribution objectives. This is further confirmed in Tab. 8. We see that the smaller the w , the better the waypoint reaching error. However, for $w < -0.5$ the mean accelerations of pedestrians start to deviate more from those observed in the ground truth nuScenes data, as the model is capable of

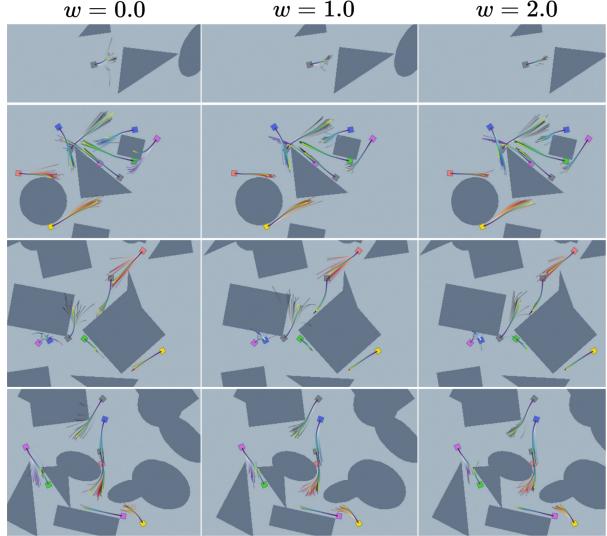


Figure 13. Sampling using increasing classifier-free weights w . 20 samples are visualized for each pedestrian. Larger w tends to emphasize collision avoidance and reduces the variance of the sampled trajectory distribution, especially for pedestrians near obstacles where conditioning has a large effect on motion.

producing more extreme trajectories to reach waypoints.

D.4. PACER– Ablation Study

In this experiment, we demonstrate the importance of multiple design decisions in the PACER model. First, we choose to make the animation controller agent-aware by including neighboring pedestrians in the heightmap given to the model. For comparison, we train a model that is agent *unaware*, *i.e.*, the input height map only contains obstacles. As seen in the top half of Tab. 9, even though TRACE is already agent-aware, having PACER endowed with awareness is highly beneficial. Both with and without agent avoidance guidance on TRACE, the agent-aware model greatly improves the collision rate.

Terrain	Model	Guide	Fail Rate	Traj Follow Error
Flat (Crowd)	Agent Unaware	None	0.252	0.102
	Agent Aware	None	0.087	0.082
	Agent Unaware	Agt Avoid	0.060	0.067
	Agent Aware	Agt Avoid	0.013	0.071
Random	Body Unaware	None	0.125	0.105
	Body Aware	None	0.093	0.104
	Body Unaware	Waypoint	0.103	0.102
	Body Aware	Waypoint	0.107	0.111

Table 9. PACER ablation study while using TRACE as the trajectory planner.

● Flat, No Guide ◆ Flat, Agt Avoid ■ Obstacles, Obs Avoid ▲ Random, Waypoint

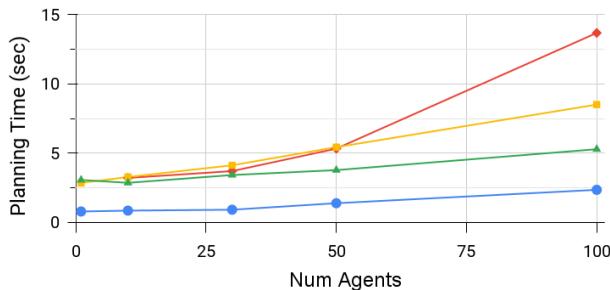


Figure 14. TRACE planning time within the end-to-end pedestrian animation system for varying terrains, guidance, and number of simulated agents.

Second, we evaluate whether making PACER body-aware is necessary, *i.e.* if it needs to take the parameters of the SMPL body β shape as input, as it is already training in simulation with a variety of body shapes. The bottom half of Tab. 9 shows that while traversing random terrains, body awareness helps to improve the failure rate when no guidance is used. When waypoint guidance is added, performance is essentially unchanged.

As motion is best seen in videos, we also include videos of how the symmetry loss and body shape conditioning affect motion quality; please see the supplementary webpage.

D.5. Runtime Analysis

Fig. 14 shows an analysis of the average runtime for one planning step of TRACE within the end-to-end animation system (on an NVIDIA TITAN RTX). Varying numbers of simulated humanoids are tested using the terrains and guidance introduced in Sec 4.3 of the main paper. With ≤ 50 agents, TRACE planning takes ≤ 5 sec, but becomes more costly with 100 agents, especially using agent avoidance guidance. Since collision avoidance requires pairwise comparisons between many agents, it can be costly.

The standalone PACER model is real time, running at ~ 30 fps for 1 humanoid and ~ 25 fps for 100.

E. Discussions and Limitations

TRACE Efficiency. The main limitation of using our system in a real-time setting is the speed of the denoising process. This is a well-known issue with diffusion models, and the community is actively working to address it. For example, recent work on distilling diffusion models [38] could be applied here to greatly speed up sampling.

Multi-Objective Guidance. One challenge with using several objectives simultaneously to guide TRACE is balancing the weight α for each. Though it is not difficult to tune each weight individually, we found that when combined, the guidance strength can be too much depending on the scene. Intuitively, if two guidance objectives are pushing a trajectory in the same direction (*e.g.* avoiding obstacle collision and going to a waypoint), the combined guidance will have compounded strength that may push the trajectory to diverge off-manifold. Work in image generation has noticed similar effects when using strong guidance, which manifests itself as saturated images. To avoid this, various forms of dynamic clipping during sampling have been introduced [54]. While this makes sense for images that have been normalized in a fixed range, it is not trivial for trajectories and we think this is an interesting problem for future work.

PACER Motions. Though PACER is robust and traverses diverse terrains while driving humanoids with different body shapes, it struggles with large obstacles when there is no way around them. The motion generated at low speed can also be unnatural as our motion database contains few samples where the humanoid is traveling at extremely low speed. Our humanoids also lack motion diversity, since most body types will have similar walking gaits and will not manifest common pedestrian behaviors such as talking on the phone or with each other. More research is needed to improve the quality and diversity of the motion.