

Controllable Motion Generation via Diffusion Modal Coupling

Luobin Wang^{*1}, Hongzhan Yu^{*1}, Chenning Yu¹, Sicun Gao¹, Henrik Christensen¹

Abstract—Diffusion models have recently gained significant attention in robotics due to their ability to generate multi-modal distributions of system states and behaviors. However, a key challenge remains: ensuring precise control over the generated outcomes without compromising realism. This is crucial for applications such as motion planning or trajectory forecasting, where adherence to physical constraints and task-specific objectives is essential. We propose a novel framework that enhances controllability in diffusion models by leveraging multi-modal prior distributions and enforcing strong modal coupling. This allows us to initiate the denoising process directly from distinct prior modes that correspond to different possible system behaviors, ensuring sampling to align with the training distribution. We evaluate our approach on motion prediction using the Waymo dataset and multi-task control in Maze2D environments. Experimental results show that our framework outperforms both guidance-based techniques and conditioned models with unimodal priors, achieving superior fidelity, diversity, and controllability, even in the absence of explicit conditioning. Overall, our approach provides a more reliable and scalable solution for controllable motion generation in robotics.

I. INTRODUCTION

Diffusion models [1] have recently emerged as a powerful class of deep generative models, combining stability, multi-modal expressiveness, and high-fidelity sample generation. These attributes make them well-suited to address the complexity and uncertainty prevalent in robotics applications. Indeed, their impact spans diverse domains, ranging from sensor simulation [2], e.g. synthesizing realistic camera data for domain randomization and sim-to-real transfer [3]), to trajectory generation [4] where they produce diverse motion plans under challenging environmental conditions [5]), and policy learning [6], wherein they model multi-step action distributions for high-dimensional control tasks [7]). These applications highlight how diffusion-based approaches can reshape core robotics paradigms, facilitating robust and adaptive robotic systems through high-fidelity probabilistic modeling.

However, several challenges persist, and the chief one among them is controllability. Generating samples that adhere to specific objectives and domain constraints is non-trivial. Naive sampling may yield physically implausible trajectories or policies lacking necessary structure. Even worse, when controllability is insufficient, practitioners often generate large numbers of samples to ensure they capture desirable outcomes. This poses scalability concerns and adds challenges to mine high-fidelity samples for large-scale synthetic generations.

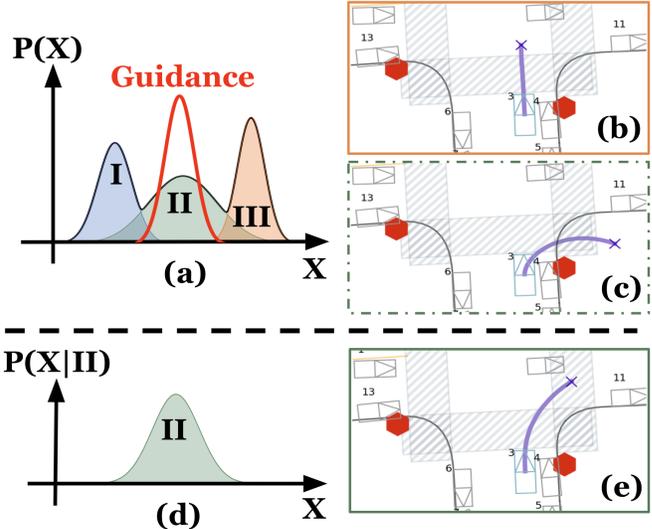


Fig. 1: High-level comparison between guidance-based and our proposed methods. (a) A standard diffusion model captures a multi-modal data distribution. For illustrative purposes, only three modes are visualized. A guidance function (in red) is employed to direct the generation towards mode II, i.e., the less likely yet operational behaviors. (b) When the guidance factor is low, it fails to effectively steer the model, which continues to produce the more probable, dominant modes. (c) Conversely, when the guidance factor is high, sample fidelity and feasibility degrades. Because guidance is applied to intermediate predictions during denoising, it can push these predictions out of the data manifold’s high-fidelity region, causing distribution mismatch. (d)(e) In our method, the less likely data mode is tightly coupled to a corresponding prior mode. During sampling, we simply initiate the denoising process from that specific prior mode, achieving direct controllability without sacrificing sample fidelity.

Existing solutions commonly resort to constraint-based sampling or post-hoc guidance [8], [9]. The core idea is to incorporate domain knowledge and task-specific constraints during the sampling process. While such methods indeed promote better alignment with target objectives, they risk degrading sample fidelity by shifting outputs away from the learned high-fidelity data manifold. This is because manually imposed constraints or guidance signals risk overriding the model’s intrinsic generative dynamics. Consequently, suboptimal or unnatural samples that no longer reflect the true data distribution would be led. Balancing fine-grained control with the preservation of sample fidelity thus remains an open research problem (Figure 1).

In this work, we propose a novel framework that achieves strong controllability in diffusion models, without incurring the distribution drift commonly encountered in guidance-based approaches. The key idea is to replace the unimodal

^{*}Equal Contribution, ¹University of California, San Diego

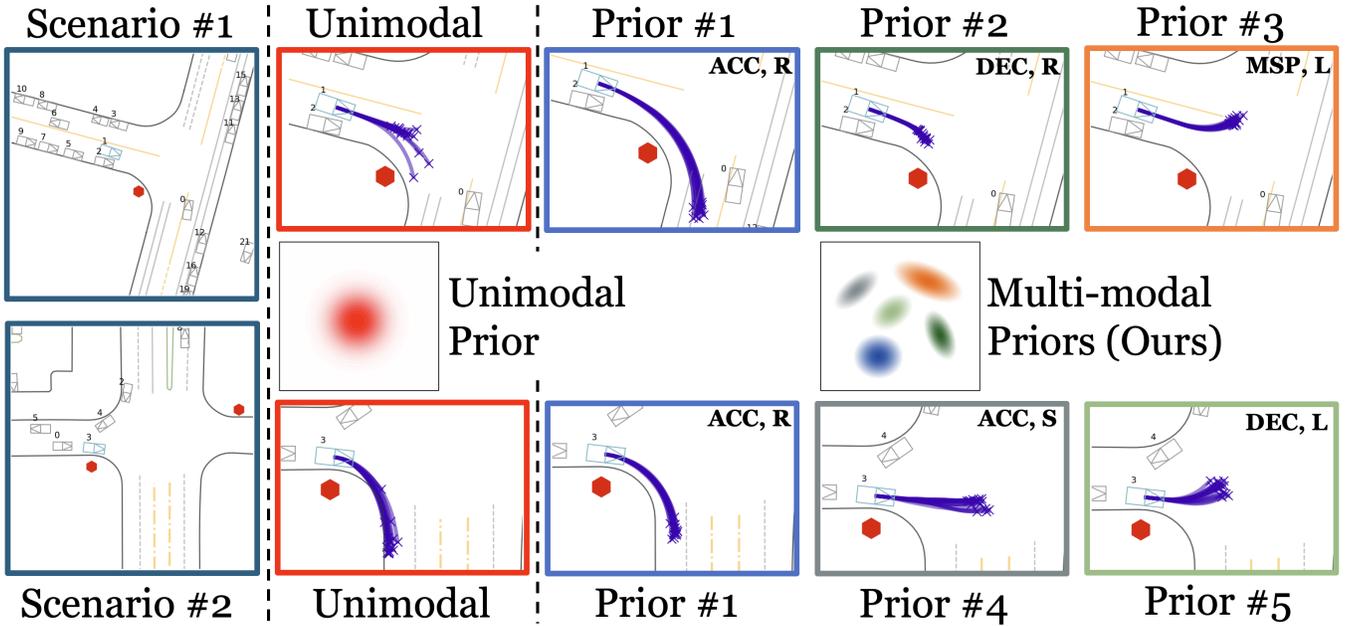


Fig. 2: High-level overview. Conventional diffusion models rely on a unimodal prior distribution, providing no inherent mechanism to control which trajectories get emphasized. This limitation can neglect lower-probability yet operationally crucial motion plans - an important limitation for robotics applications that demand flexible behaviors. In our work, we propose a multi-modal prior distribution and enforcing strong modal coupling between prior and data distributions through a novel diffusion process. By denoising directly from the relevant prior mode, we preserve trajectory fidelity and realism without risking distribution mismatch. Crucially, even with an unconditioned diffusion model, our framework offers straightforward control over which data mode to attend, facilitating precise and adaptive motion generation. In the figure, each prior corresponds to one system behavior. “ACC”, “DEC” and “MSP” refer to speed modes (acceleration, deceleration, and maintaining speed), while “R”, “L” and “S” represent steering modes (right, left, and straight).

prior in a standard diffusion model with a multi-modal distribution that is tightly coupled to the principle modes of the target data. Therefore at sampling, rather than applying post-hoc guidance to steer the generative process toward a specific mode, we initiate the denoising process directly from the corresponding prior mode. This ensures that sampling remains inherently aligned with the training distribution, eliminating the need for external guidance and mitigating the risk of distribution mismatch (Figure 2).

Our current formulation assumes that the data modes subject to control are explicitly known, yet the diffusion model itself does not require conditioning on these mode labels. While this assumption may appear restrictive, it serves as a foundation step towards enabling strong controllability over data with unknown key modes, wherein the central challenge shifts to accurately identifying the appropriate prior mode based on target constraints. Nevertheless, we demonstrate that by adopting a *multi-modal* prior distribution, strong *modal coupling*, and a careful *prior parametrization*, our method significantly outperforms guidance-based techniques and even conditioned modeling with a unimodal prior, in terms of both fidelity and controllability. We validate these claims on the Waymo dataset for motion prediction, and in Maze2D for multi-task control. The paper is organized as follows. Section II and III cover the related work and the necessary background. We detail the proposed method in Section IV. Section V presents the experimental results, and Section VI concludes the paper.

II. RELATED WORK

The multi-modal nature of human behaviors poses a great challenge for predicting realistic trajectories and control sequences. Diffusion models have proven effective in capturing this multi-modality within driving scenarios while closely adhering to real-world behavior distributions. SceneDM [10] utilizes a diffusion-based framework to model joint-distributions of all agents in a scene. SceneDiffuser [11] employs a latent diffusion architecture derived from Bird’s Eye View representations, whereas MotionDiffuser [4] demonstrates its capabilities of predicting realistic future trajectories that align with true data distribution via PCA-compressed trajectory representations. Additionally, VBD [12] jointly optimizes a motion predictor and a denoiser that share the same scene encoder, which further improves realism and versatility. However, achieving such realism and broad distribution coverage often requires drawing many random samples from a standard Gaussian prior that is unimodal. Our approach enhances the realism of generated trajectories by incorporating a multi-modal prior that more effectively captures distinct data modes.

The typical strategy to control diffusion-based generation is incorporating domain-specific objectives into the generation process. Classifier guidance [13] guides the diffusion model with a separate cost function that encodes the objectives during sampling. Recent works [4] propose different analytical guidance functions to achieve either realism or safety-critical objectives. On the other hand, Classifier-free

guidance [8] additionally optimizes a time-dependent conditional model to obtain guidance. For instance, VBD [12] achieves this using the motion predictor jointly optimized for goal-guided trajectory generation. However, the constraints imposed by guidance often degrade the realism of model generation due to distribution mismatch. Our approach completely avoids this issue by coupling prior and data modes.

III. BACKGROUND

Diffusion models are probabilistic generative models that synthesize new data by iteratively denoising an initial noise sample. These models begin by defining a forward stochastic process that progressively adds noise to real data x_0 , eventually converting it into pure noise. Formally, this forward noising process is defined as:

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}), \quad (1)$$

$$q(x_t|x_{t-1}) = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon_t, \quad (2)$$

where $x_{1:T}$ denotes the sequence of noised samples from step 1 to T , ϵ_t is standard Gaussian noise, and β_t is the pre-defined forward variance. Let $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{i=1}^t \alpha_i$. This forward process allows one to directly sample x_t at any intermediate step t in close form:

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I). \quad (3)$$

In a typical diffusion setup, the variance schedule $\beta_{1:T}$ is chosen such that $\bar{\alpha}_T$ approaches 0 under a sufficiently large T . This ensures the final noised state x_T to converge to pure noise, establishing a standard unimodal Gaussian *prior*.

A reverse process recovers clean samples by removing noise step by step, modeled as a Gaussian distribution whose mean is given by a neural network and whose variance is fixed based on the forward variance schedule. Concretely,

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t), \sigma_t^2 I), \quad (4)$$

$$\sigma_t^2 = \frac{1 - \bar{\alpha}_{t-1}\beta_t}{1 - \bar{\alpha}_t}, \quad (5)$$

where θ denotes the parameters of the neural denoiser trained to maximize the variational lower bound [1] on the log-likelihood of the observed data x_0 : $\max_\theta -\log p_\theta(x_0|x_1) + \sum_t D_{KL}[q(x_{t-1}|x_t, x_0)||p_\theta(x_{t-1}|x_t)]$. To simplify learning, one can reparameterize μ_θ in terms of noise prediction: $\mu_\theta(x_t) = \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(x_t, t))$, where ϵ_θ is optimized to predict the true noise, approximating the *score*, i.e., the gradient of log density $\nabla_{x_t} \log q(x_t)$, across different noise scales. Alternatively, one can predict the clean sample x_θ^0 , using it directly to derive close-form posterior mean:

$$\mu_\theta(x_t) = \frac{\sqrt{\bar{\alpha}_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}x_\theta^0(x_t, t). \quad (6)$$

Controllability in diffusion models refers to the ability to impose particular constraints on the generated samples. A common strategy for achieving this is through guidance [8], where an additional cost function f is integrated into each

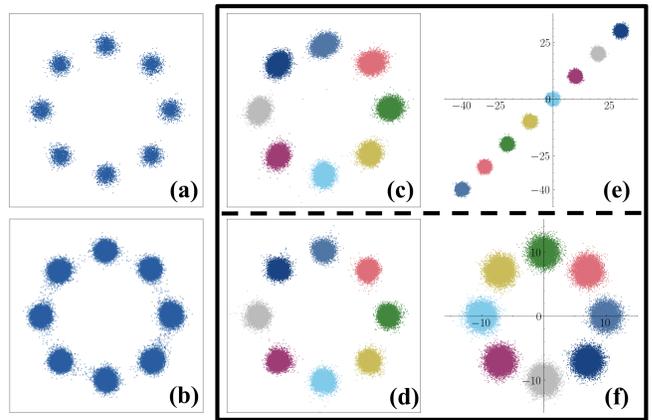


Fig. 3: 2D toy example. (a) The data distribution with eight distinct modes. (b) Results from DDPM [1], highlighting spurious samples in the gaps between modes. (c)-(d) Results from our method, which trains an unconditioned diffusion model but explicitly pairs each prior mode with a data mode. This yields significantly fewer spurious samples and provides direct control over which mode is generated. (e)-(f) The pre-defined prior parametrizations. In particular, (e) shows that when the prior means lie far from the origin (i.e., have large magnitudes), the diffusion model struggles to learn those corresponding modes effectively.

step of the reverse denoising process. Formally, the learned scores are modified as follows:

$$\hat{\epsilon}_\theta(x_t, t) = \epsilon_\theta(x_t, t) + \omega \nabla_{x_t} f(x_t), \quad (7)$$

where ω controls the influence of f . This cost function may encode a variety of task-specific constraints, such as enforcing speed limits or steering preferences in autonomous driving, restricting actions in control tasks, or nudging the model towards a particular label or style in image generation.

IV. ENHANCED DIFFUSION CONTROLLABILITY VIA MODAL COUPLING

A common approach to controlling diffusion-based generation is through guidance. However, as the guidance function is not integrated into the core training objective, it imposes post-hoc constraints that may incur a distribution mismatch, causing samples to drift away from high-fidelity regions of the data manifold.

To address these challenges, we propose adopting a multi-modal prior, and pose the following question: *if each prior mode is coupled with a corresponding mode in the target distribution, can we run the reverse process from different prior modes for strong controllability without guidance?* Not only is the answer “yes”, but it also promises to eliminate reliance on post-hoc constraints. As a result, each mode of the target distribution is naturally accommodated during sampling, with no mismatch risk posed by guidance.

In this work, we assume that the modes of data distribution are explicitly known, and model the multi-modal prior as a Gaussian mixture model:

$$x_T \sim \sum_{i=1}^k r_i \cdot \mathcal{N}(\mu_i, \sigma_i^2 I), \quad (8)$$

where k is the number of modes, r_i denotes the proportion of data with mode label i , and μ_i and σ_i^2 specify the mean and variance of the i -th Gaussian component, respectively. Notably, when conditioning on a specific data x_0 with label L , the prior reduces to a unimodal Gaussian:

$$x_T|x_0 \sim \mathcal{N}(\mu_L, \sigma_L^2 I). \quad (9)$$

This structure allows us to explicitly account for different modes in the data while retaining a unimodal form given a specific label. In the sequel, we will demonstrate that, by *coupling* these prior modes with the corresponding data modes and by carefully *parameterizing* each prior mode, we can achieve robust mode control even using unconditioned diffusion models.

A. Modal Coupling

We begin by defining forward and reverse diffusion processes that accommodate a general Gaussian prior $x_T|x_0 \sim \mathcal{N}(\mu, \sigma^2 I)$.

Lemma 1. Let $\eta_t := 1 + \sum_{m=1}^{t-1} \left(\sqrt{\prod_{n=m+1}^t \alpha_n} \right)$, and consider the forward noising process

$$q(x_t|x_{t-1}) = \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\sigma\epsilon_t + \frac{\mu}{\eta_T}. \quad (10)$$

where $\epsilon_t \sim \mathcal{N}(0, I)$. Then, for any step t ,

$$q(x_t|x_0) = \mathcal{N}\left(x_t|\sqrt{\bar{\alpha}_t}x_0 + \frac{\eta_t\mu}{\eta_T}, (1 - \bar{\alpha}_t)\sigma^2 I\right). \quad (11)$$

Under the standard assumption that $\bar{\alpha}_T \rightarrow 0$ as T grows large, it follows that $q(x_T|x_0) = \mathcal{N}(x_T; \mu, \sigma^2 I)$.

The proof is in Appendix VI-A. Lemma 1 establishes that introducing the constant shift term μ/η_T at each forward step and scaling the Gaussian noise by σ ensure the final marginal $q(x_T|x_0)$ to align with the desired Gaussian prior, which sets the foundation for the corresponding reverse process.

Lemma 2. For the diffusion model with the forward process defined in (10), the reverse process is:

$$p(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}|\mu(x_t), \beta(x_t)), \quad (12)$$

where $\beta(x_t) = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \sigma^2$, and

$$\begin{aligned} \mu(x_t) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\sqrt{\alpha_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 \\ &\quad + \frac{\eta_{t-1}(1 - \alpha_t) - \sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{(1 - \bar{\alpha}_t)\eta_T} \cdot \mu, \end{aligned} \quad (13)$$

$$= \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\mu}{\eta_T} - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \sigma \epsilon_t \right). \quad (14)$$

We include the proof in Appendix VI-B. This concludes the derivation of the proposed forward and reverse processes that support an arbitrary Gaussian, thereby enabling any prior mode to be directly coupled with a specific mode in the target distribution. Given a labeled dataset \mathcal{X} and a set of pre-defined prior mode parameterizations, we train the diffusion

model x_θ^0 by:

$$\min_{\theta} E_{\substack{t \in [1, T] \\ \epsilon \sim \mathcal{N}(0, I)}} \left[\sum_{(x_0, L) \in \mathcal{X}} \|x_\theta^0(\hat{x}_t(x_0, L, \epsilon), t) - x_0\|^2 \right], \quad (15)$$

$$\hat{x}_t(x_0, L, \epsilon) = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\sigma_L\epsilon + \frac{\eta_t\mu_L}{\eta_T}. \quad (16)$$

Here, we adopt the reparameterization of μ_θ , i.e., the mean of posterior $p(x_{t-1}|x_t)$, in term of clean sample prediction. Switching to noise prediction is straightforward but omitted due to performance considerations.

At sampling, we can either draw from the multi-modal prior distribution (8) to generate across all data modes, or select a unimodal prior (9) to concentrate on a specific mode, allowing fine-grained control over the diffusion process.

B. Prior Parametrization

Next, we present our parametrization of the proposed multi-modal prior distribution. As shown in Figure 3, it is crucial to keep the prior modes well-separated (i.e., not overlapping) so that the diffusion model can clearly distinguish the target mode from the others. Additionally, each mode must be constrained from generating excessively large samples that could complicate the diffusion process.

To specify the Gaussian mean of each prior mode, we draw on the concept of placing *evenly spaced* points on a high-dimensional sphere. Let d be the data dimension, and k the number of modes. If $k \leq d + 1$, which is often the case in complex data distributions suited to diffusion models, this arrangement can be achieved by placing the vertices of a $(k - 1)$ -simplex in \mathbb{R}^d [14]. Let $e_{1:k} \subset \mathbb{R}^k$ be the standard basis vectors. We define the means $\mu_{1:k}$ as follows:

$$w_i = \delta \cdot \sqrt{\frac{k}{k+1}} \cdot (e_i - \mathbb{1}_k \cdot \frac{1}{k}) \quad (17)$$

$$\mu_i = [w_i^1, \dots, w_i^k, 0, \dots, 0]^T \in \mathbb{R}^d, \quad (18)$$

where δ is a hyper-parameter for adjusting the sphere radius and w_i^j denotes the j -th component of w_i . One can verify that $\|\mu_i\|^2 = \delta^2$ for all $i \in [1, k]$, and that these points lie at a pairwise distance of $\delta \cdot \sqrt{2 + 2/(k - 1)}$. Finally, we choose the variance of each Gaussian so that their confidence intervals are at a certain level do not overlap with respect to the above pairwise distance. That is, for all $i \in [1, k]$,

$$\sigma_i \sqrt{\mathcal{X}_{d,c}^2} \leq \delta \cdot \sqrt{2 + \frac{2}{(k - 1)}}, \quad (19)$$

where c denotes the desired confidence level and $\mathcal{X}_{d,c}^2$ is the c quantile of the chi-square distribution with d degree of freedom [15].

V. EXPERIMENTS

A. Waymo - Controllable Motion Prediction

We evaluate the proposed method on the motion prediction task using the Waymo Open Motion Dataset [16] that features diverse and realistic autonomous driving scenarios. Because motion prediction is inherently probabilistic and

Method		minADE (↓)	meanADE (↓)	minFDE (↓)	meanFDE (↓)	minMR (↓)	meanMR (↓)	minOR (↓)	meanOR (↓)
VBD, $S = 1$		1.707	1.707	4.593	4.593	0.551	0.551	0.170	0.170
VBD, $S = 3$		1.265	1.717	3.294	4.597	0.394	0.549	0.126	0.172
VBD, $S = 9$		0.968	1.708	2.422	4.557	0.279	0.544	0.102	0.178
cVBD, $k = 9, S = 1$		1.425	1.425	3.667	3.667	0.437	0.437	0.173	0.173
cVBD, $k = 9, S = 3$		1.025	1.418	2.530	3.650	0.290	0.430	0.126	0.178
Ours, $S = 1$	$k = 3, \delta = 8$	1.267	1.267	3.523	3.523	0.378	0.378	0.107	0.107
	$k = 3, \delta = 15$	1.216	1.216	3.221	3.221	0.367	0.367	0.117	0.117
	$k = 9, \delta = 8$	1.019	1.019	2.678	2.678	0.260	0.260	0.070	0.070
	$k = 9, \delta = 15$	0.940	0.940	2.423	2.423	0.212	0.212	0.057	0.057
Ours, $S = 3$	$k = 3, \delta = 8$	1.093	1.281	2.986	3.565	0.299	0.382	0.090	0.109
	$k = 3, \delta = 15$	1.117	1.219	2.937	3.223	0.322	0.364	0.103	0.139
	$k = 9, \delta = 8$	0.892	1.017	2.323	2.662	0.210	0.262	0.050	0.065
	$k = 9, \delta = 15$	0.857	0.938	2.190	2.419	0.182	0.209	0.046	0.060

TABLE I: Quantitative validation for motion prediction accuracy on Waymo. We evaluate the proposed method using two configurations for the number of modes: $k = 3$, which considers only steering modes, and $k = 9$, which incorporates both steering and speed modes. S denotes the number of sampled trajectories per scenario for each setting. Due to limited computing resources, the training size was reduced to one-tenth of its original size. We anticipate that the reported metrics would improve further if trained on the full dataset.

multi-modal, it is crucial to capture an unbiased distribution of possible future trajectories. In our experiments, we focus on the *single-agent* setting where distinct modes in the target distribution are more clearly defined, thus better highlighting the advantages of our approach.

We adopt Versatile Behavior Diffusion (VBD) [12] as the backbone of our model. A query-centric, Transformer-based scene encoder first processes the scene context, including scenario information such as traffic lights, map polylines, and agent trajectories, to capture interrelationships among all scene components. Building on this contextual understanding, a denoiser model then predicts a clean control sequence from a noised input. Under the single-agent setting, the denoiser considers only the ego-agent’s control trajectory, while other agents remain visible to the scene encoder, given their potential influence on the ego-agent’s future behavior.

Next, we describe the parameter settings of the proposed method. We adopt a two-layered mode design: the first layer corresponds to steering modes (left turn, right turn, and go straight), while the second layer addresses speed control (acceleration, deceleration, and maintaining speed). If both layers are considered, we take the Cartesian product of these two layers, which yields $k = 9$ total modes. During training, we set $\sigma_i = 1$ for all $i \in [1, k]$. The data dimension is 80, corresponding to 2D control predictions over 40 future steps (i.e., 4 seconds).

1) *Motion Prediction Accuracy*: Table I summarizes the primary metrics used to evaluate prediction accuracy on the test set. We report four standard metrics. *Average Displacement Error (ADE)* measures the mean distance between the ground-truth future ego-agent trajectories and the model predictions, averaged over the future time horizon. *Final Displacement Error (FDE)* measures the displacement error at the final time step. *Miss Rate (MR)* measures the recall of the trajectory predictions. And *OffRoad (OR)* assesses prediction consistency in how often the predicted trajectories

drive off the road. For each metric, we report the minimum and mean values computed from the sampled trajectories. In addition, we implement a variant of VBD, referred to as *conditioned VBD (cVBD)*, which is given explicit mode labels as part of its input.

Although enhancing prediction performance is not the primary goal of this project, Table I demonstrates notable improvements in the reported metrics, especially in the mean statistics. We attribute these gains primarily to our use of a multi-modal prior distribution, mirroring the positive results observed on the 2D toy set (Figure 3). This benefit is further underscored by the performance increase when the number of modes k rises from 3 to 9. Meanwhile, the configurations with $\delta = 15$ outperform those with $\delta = 8$, emphasizing the importance of ensuring that prior mode distributions do not overlap. Lastly, as cVBD still relies on a unimodal prior distribution, its ability to effectively learn each mode in the target distribution is limited. By contrast, our method leverages strong modal coupling to better capture multi-modal behavior, ultimately leading to more robust predictions.

2) *Controllable Trajectory Synthesis*: We assess controllability in Table II by measuring how well sampled trajectories remain realistic and feasible despite deviating from the reference trajectories (Figure 4). To generate Table II, we manually label potential ego-agent futures with an emphasis on steering feasibility. For instance, if the ego-agent is in a lane that allows turning, yet the reference trajectory continues straight, we count the turning maneuver as feasible. However, identifying valid speed modes that differ from the reference is more challenging, and thus, occasional mislabeling may occur. The final evaluation set comprises 140 scenarios, which we provide in our codebase ¹.

¹<https://github.com/RobinWangSD/Diffusion-with-Multi-Modal-Priors.git>

Method		minADE (\uparrow)	meanADE (\uparrow)	minOR (\downarrow)	meanOR (\downarrow)	ACC[ST] (\uparrow)	ACC[SP] (\uparrow)	ACC (\uparrow)
VBD+G	$\omega = 1$	1.049	1.814	0.070	0.108	0.877	0.450	0.387
	$\omega = 10$	1.790	2.717	0.070	0.113	0.952	0.624	0.601
	$\omega = 100$	2.294	3.443	0.121	0.184	0.944	0.711	0.701
cVBD,	$k = 9$	1.996	2.901	0.111	0.205	0.883	0.806	0.707
Ours	$k = 9, \delta = 8$	1.932	2.233	0.035	0.050	0.906	0.702	0.647
	$k = 9, \delta = 15$	2.162	2.369	0.023	0.026	0.917	0.701	0.662

TABLE II: Quantitative validation for trajectory synthesis controllability, using 9 sampled trajectories per scenario. The test set consists of 140 scenarios, each manually labeled with potentially feasible ego-agent futures that deviate from the dataset’s original trajectories. Note that higher *ADE* values imply a greater deviation from the reference trajectories, not indicating the degraded prediction performance.

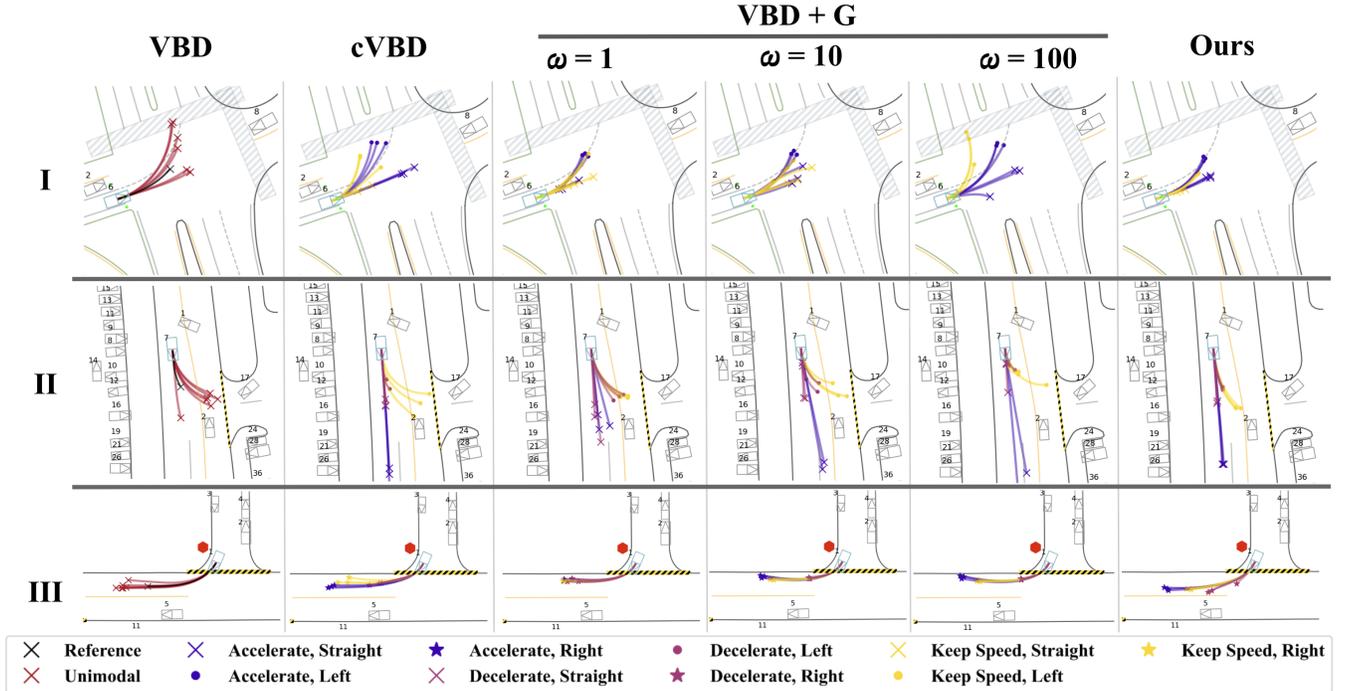


Fig. 4: Qualitative results for motion prediction. In standard diffusion, trajectories are sampled randomly from a unimodal prior, offering no inherent controllability. cVBD takes the intended mode as part of its model input, but still employs a unimodal prior, preventing it from effectively presenting distinct data modes. VBD-G applies guidance to steer generation, but it relies heavily on the guidance influence factor, making it difficult to balance sample fidelity against controllability. Our method integrates modal coupling with a multi-modal prior distribution, yielding notable improvements in both sample fidelity and controllability.

In addition to the conditioned VBD baseline, we examine a post-hoc guidance strategy for the unconditioned VBD model, referred to as *VBD+G*. At each denoising step $t \in [2, T]$, we update the posterior mean:

$$\hat{\mu}_\theta(x_t) = \mu_\theta(x_t) + \omega \cdot \nabla_{\mu_\theta(x_t)} f(x_\theta^0(\mu_\theta(x_t), t - 1)), \quad (20)$$

where $\mu_\theta(x_t)$ is the posterior mean estimate computed from the predicted $x_\theta^0(x_t, t)$ via (6), ω is a parameter controlling guidance strength, and f is the guidance function applied to the predicted x_0 , treating $\mu_\theta(x_t)$ as x_{t-1} . These steps follow the procedures described in [12]. In our experiments, f aligns with the mode list used during training for both the proposed method and the conditioned VBD. Specifically, given a target steer-speed label pair, f penalizes any input trajectory that

exceeds the specified steer-speed range.

To evaluate performance, we continue using *Average Displacement Error (ADE)* with respect to reference trajectories, which measures how far sampled trajectories differ from the reference. We also include *OffRoad (OR)*, an important metric that captures fidelity and feasibility of sampled trajectories. Finally, we introduce an *overall accuracy (ACC)* metric to measure how often the samples match the intended steer-speed mode, and further decompose it into separate measures for steering (*ACC[ST]*) and speed (*ACC[SP]*).

As shown in Table II, the guidance-based method struggles to balance sample fidelity and controllability. Although increasing ω does improve overall accuracy, it also significantly worsens *OffRoad*. In contrast, our proposed method

Method		U-Maze	Medium	Large
Diffuser	U-Maze	113.9	N/A	N/A
	Medium	N/A	121.5	N/A
	Large	N/A	N/A	123.0
Ours, $k = 3, \delta = 30$		119.5	121.4	120.9

TABLE III: Quantitative evaluations on Maze2D. We adopt the baseline performance from [17], using normalized accumulative reward returns as the evaluation metric. Notably, the baseline trains a separate diffusion model for each layout (or mode), unlike in Table I where our comparisons focus on multi-modal data modeling and the baseline there is a single model handling various modes.

consistently maintains high sample fidelity while achieving reasonably high accuracy. Furthermore, as δ increases, thus separating each prior mode more distinctly, performance continues to improve. It is worth noting that VBD-G requires extensive backpropagation through the diffusion model, making its generation noticeably slower than our method which simply performs a standard reverse process from the specified prior mode. Meanwhile, the conditioned VBD achieves comparable performance to VBD-G at $\omega = 100$.

B. Maze2d - Multi-task Control

The proposed method naturally extends to multi-task control. Rather than confining control-level constraints to a single task, we can also view each task itself as a separate mode. This perspective allows our method to compose multiple tasks within a single, unconditioned diffusion model. We illustrate this with evaluations in Maze2D, where a unified model is trained to perform long-horizon path planning across various maze configurations.

We define the diffusion model to predict 384-step state-control trajectories conditioned on the given initial and goal positions. Given that the state and control dimensions are 3 and 2, respectively, the target distribution has dimension 1920. We set $\sigma_i = 1$ for $i \in [1, 3]$, and choose $\delta = 30$. As shown in Table III, the proposed method achieves performance comparable to the baseline across all tested maze configurations. Note that our aim is not to outperform the existing methods. Instead, we demonstrate that by treating each distinct task as an individual data mode, the proposed method establishes effective modal coupling over a multi-modal prior, which makes it possible to handle various tasks using only a single unconditioned diffusion model.

VI. CONCLUSION

This paper presents a novel framework that enables fine-grained control over diffusion models while preserving high-fidelity sample generation. By aligning the sampling process with key data modes from the outset, our method avoids the distribution drift common in post-hoc guidance approaches. Experimental evaluations show that the proposed method consistently outperforms existing techniques on both quantitative and qualitative measures. Moreover, this work lays a strong foundation for future research aimed at relaxing the

assumption of explicit known data modes, thereby advancing towards more controllable diffusion models.

APPENDIX

A. Proof of Lemma 1

The proof for (11) proceeds by induction. We begin with the base case using the proposed forward process (10):

$$x_1 = \sqrt{\alpha_1}x_0 + \sqrt{1 - \alpha_1}\sigma\epsilon_1 + \frac{\mu}{\eta_T} \quad (21)$$

$$\sim \mathcal{N}(\sqrt{\alpha_1}x_0 + \frac{\eta_1\mu}{\eta_T}, (1 - \bar{\alpha}_1)\sigma^2I) \quad (22)$$

The derivation from (21) to (22) is based on the fact that $\alpha_1 = \bar{\alpha}_1$ and $\eta_1 = 1$ by definition. Next, we assume that for an arbitrary $t \in [2, T]$, it holds true that:

$$x_{t-1} = \sqrt{\bar{\alpha}_{t-1}}x_0 + \sqrt{1 - \bar{\alpha}_{t-1}}\sigma\epsilon_{t-1} + \frac{\eta_{t-1}\mu}{\eta_T}. \quad (23)$$

From the forward process (10), we have:

$$x_t = \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\sigma\epsilon_t + \frac{\mu}{\eta_T} \quad (24)$$

$$= \sqrt{\alpha_t} \left(\sqrt{\bar{\alpha}_{t-1}}x_0 + \sqrt{1 - \bar{\alpha}_{t-1}}\sigma\epsilon_{t-1} + \frac{\eta_{t-1}\mu}{\eta_T} \right) \quad (25)$$

$$+ \sqrt{1 - \alpha_t}\sigma\epsilon_t + \frac{\mu}{\eta_T}$$

$$= \sqrt{\bar{\alpha}_t}x_0 + \frac{\sqrt{\alpha_t}\eta_{t-1} + 1}{\eta_T} \cdot \mu \quad (26)$$

$$+ \left(\sqrt{\alpha_t(1 - \bar{\alpha}_{t-1})}\sigma\epsilon_{t-1} + \sqrt{1 - \alpha_t}\sigma\epsilon_t \right).$$

Note that by the definition of η_t :

$$\eta_t = 1 + \sum_{m=1}^{t-1} \left(\sqrt{\prod_{n=m+1}^t \alpha_n} \right) \quad (27)$$

$$= 1 + \sum_{m=1}^{t-1} \left(\sqrt{\prod_{n=m+1}^{t-1} \alpha_n \cdot \sqrt{\alpha_t}} \right) \quad (28)$$

$$= 1 + \sqrt{\alpha_t} \cdot \left[1 + \sum_{m=1}^{t-2} \left(\sqrt{\prod_{n=m+1}^{t-1} \alpha_n} \right) \right] \quad (29)$$

$$= 1 + \sqrt{\alpha_t}\eta_{t-1}. \quad (30)$$

Meanwhile, to handle the last term in (26), we essentially merge two zero-mean Gaussian distributions with distinct variances. That is, merging $\mathcal{N}(0, \sigma_1^2I)$ and $\mathcal{N}(0, \sigma_2^2I)$ leads to the new distribution $\mathcal{N}(0, (\sigma_1^2 + \sigma_2^2)I)$. Here, the merged standard deviation is:

$$\alpha_t(1 - \bar{\alpha}_{t-1})\sigma^2 + (1 - \alpha_t)\sigma^2 = (1 - \bar{\alpha}_t)\sigma^2. \quad (31)$$

Substituting everything back into (26), we have:

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \frac{\eta_t}{\eta_T}\mu + \sqrt{1 - \bar{\alpha}_t}\sigma\epsilon^*, \quad (32)$$

where ϵ^* denotes an arbitrary standard Gaussian sample. This concludes the proof of Lemma 1.

B. Proof of Lemma 2

First, the reverse probability is tractable only when conditioned on x_0 . By Bayes' theorem, we have:

$$p(x_{t-1}|x_t, x_0) = q(x_t|x_{t-1}, x_0) \cdot \frac{q(x_{t-1}|x_0)}{q(x_t|x_0)}. \quad (33)$$

Then by Lemma 1:

$$p(x_{t-1}|x_t, x_0) \quad (34)$$

$$\propto \exp\left(-\frac{1}{2}\left[\frac{(x_t - \sqrt{\alpha_t}x_{t-1} - \mu/\eta_T)^2}{(1 - \alpha_t)\sigma^2} + \frac{(x_{t-1} - \sqrt{\bar{\alpha}_{t-1}}x_0 - (\eta_{t-1}/\eta_T)\mu)^2}{(1 - \bar{\alpha}_{t-1})\sigma^2} - \frac{(x_t - \sqrt{\bar{\alpha}_t}x_0 - (\eta_t/\eta_T)\mu)^2}{(1 - \bar{\alpha}_t)\sigma^2}\right]\right) \quad (35)$$

$$\propto \exp\left(-\frac{1}{2}\left[\left(\frac{\alpha_t}{(1 - \alpha_t)\sigma^2} + \frac{1}{(1 - \bar{\alpha}_{t-1})\sigma^2}\right)x_{t-1}^2 + 2\left(\frac{-\sqrt{\alpha_t}x_t + \sqrt{\alpha_t}\mu/\eta_T}{(1 - \alpha_t)\sigma^2} + \frac{-\sqrt{\bar{\alpha}_{t-1}}x_0 - (\eta_{t-1}/\eta_T)\mu}{(1 - \bar{\alpha}_{t-1})\sigma^2}\right)x_{t-1}\right]\right) \quad (36)$$

From (35) to (36), the constant terms that do not involve x_{t-1} are all omitted. Following the standard Gaussian density function, the mean and variance of $p(x_{t-1}|x_t, x_0)$ can be parameterized as $\mathcal{N}(\tilde{\mu}, \tilde{\beta})$ where:

$$\tilde{\beta} = 1/\left(\frac{\alpha_t}{(1 - \alpha_t)\sigma^2} + \frac{1}{(1 - \bar{\alpha}_{t-1})\sigma^2}\right) \quad (37)$$

$$= \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \cdot (1 - \alpha_t)\sigma^2. \quad (38)$$

Then we derive $\tilde{\mu}$ as follows:

$$\tilde{\mu}(x_t, x_0) = -\left(\frac{-\sqrt{\alpha_t}x_t + \sqrt{\alpha_t}\mu/\eta_T}{(1 - \alpha_t)\sigma^2} + \frac{-\sqrt{\bar{\alpha}_{t-1}}x_0 - (\eta_{t-1}/\eta_T)\mu}{(1 - \bar{\alpha}_{t-1})\sigma^2}\right) \cdot \tilde{\beta}_t \quad (39)$$

$$= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)}{1 - \bar{\alpha}_t}x_0 + \frac{\eta_{t-1}(1 - \alpha_t) - \sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{(1 - \bar{\alpha}_t)\eta_T} \cdot \mu \quad (40)$$

Furthermore, we can parameterize x_0 in terms of x_t and ϵ_t based on (32):

$$\begin{aligned} \tilde{\mu}_t(x_t) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t \\ &+ \frac{\sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)}{1 - \bar{\alpha}_t} \cdot \frac{x_t - \sqrt{1 - \bar{\alpha}_t}\sigma\epsilon_t - (\eta_t/\eta_T)\mu}{\sqrt{\bar{\alpha}_t}} \\ &+ \frac{\eta_{t-1}(1 - \alpha_t) - \sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{(1 - \bar{\alpha}_t)\eta_T} \cdot \mu \\ &= \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\sigma\epsilon_t\right) \\ &+ \frac{\mu}{(1 - \bar{\alpha}_t)\eta_T} \cdot \left(\eta_{t-1}(1 - \alpha_t) - \sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1}) - \frac{(1 - \alpha_t)\eta_t}{\sqrt{\alpha_t}}\right) \end{aligned} \quad (41)$$

$$= \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\sigma\epsilon_t\right) + \frac{\mu}{(1 - \bar{\alpha}_t)\eta_T} \cdot \frac{-(1 - \bar{\alpha}_t)}{\sqrt{\alpha_t}} \quad (43)$$

$$= \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{\mu}{\eta_T} - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\sigma\epsilon_t\right) \quad (44)$$

Last three equations are due to $\eta_t = 1 + \sqrt{\alpha_t}\eta_{t-1}$.

REFERENCES

- [1] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [2] Haoxi Ran, Vitor Guizilini, and Yue Wang. Towards realistic scene generation with lidar diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14738–14748, 2024.
- [3] Anthony Zhou, Zijie Li, Michael Schneier, John R Buchanan Jr, and Amir Barati Farimani. Text2pde: Latent diffusion models for accessible physics simulation. *arXiv preprint arXiv:2410.01153*, 2024.
- [4] Chiyu Jiang, Andre Cornman, Cheolho Park, Benjamin Sapp, Yin Zhou, Dragomir Anguelov, et al. Motiondiffuser: Controllable multi-agent motion prediction using diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9644–9653, 2023.
- [5] Brian Yang, Huangyuan Su, Nikolaos Gkanatsios, Tsung-Wei Ke, Ayush Jain, Jeff Schneider, and Katerina Fragkiadaki. Diffusion-es: Gradient-free planning with diffusion for autonomous driving and zero-shot instruction following. *arXiv preprint arXiv:2402.06559*, 2024.
- [6] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, page 02783649241273668, 2023.
- [7] Guangyao Zhou, Sivaramakrishnan Swaminathan, Rajkumar Vasudeva Raju, J Swaroop Guntupalli, Wolfgang Lehrach, Joseph Ortiz, Antoine Dedieu, Miguel Lázaro-Gredilla, and Kevin Murphy. Diffusion model predictive control. *arXiv preprint arXiv:2410.05364*, 2024.
- [8] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- [9] Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 843–852, 2023.
- [10] Zhiming Guo, Xing Gao, Jianlan Zhou, Xinyu Cai, and Botian Shi. Scenedm: Scene-level multi-agent trajectory generation with consistent diffusion models. *arXiv preprint arXiv:2311.15736*, 2023.
- [11] Ethan Pronovost, Meghana Reddy Ganesina, Noureldin Hendy, Zeyu Wang, Andres Morales, Kai Wang, and Nick Roy. Scenario diffusion: Controllable driving scenario generation with diffusion. *Advances in Neural Information Processing Systems*, 36:68873–68894, 2023.
- [12] Zhiyu Huang, Zixu Zhang, Ameya Vaidya, Yuxiao Chen, Chen Lv, and Jaime Fernández Fisac. Versatile behavior diffusion for generalized traffic agent simulation. *arXiv preprint arXiv:2404.02524*, 2024.
- [13] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- [14] Harold Scott Macdonald Coxeter. *Regular polytopes*. Courier Corporation, 1973.
- [15] Edwin B Wilson and Margaret M Hilferty. The distribution of chi-square. *Proceedings of the National Academy of Sciences*, 17(12):684–688, 1931.
- [16] Scott Etinger, Shuyang Cheng, Benjamin Caine, Chenxi Liu, Hang Zhao, Sabeek Pradhan, Yuning Chai, Ben Sapp, Charles R Qi, Yin Zhou, et al. Large scale interactive motion forecasting for autonomous driving: The waymo open motion dataset. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9710–9719, 2021.
- [17] Michael Janner, Yilun Du, Joshua B Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. *arXiv preprint arXiv:2205.09991*, 2022.