

Supplementary Material for SLEDGE: Synthesizing Driving Environments with Generative Models and Rule-Based Traffic

Kashyap Chitta* Daniel Dauner* Andreas Geiger

University of Tübingen Tübingen AI Center
<https://github.com/autonomousvision/sledge>

Abstract. In this **supplementary document**, we first provide additional methodological and experimental details. Next, we include supplementary experiments to complement the main paper. Finally, we present additional qualitative examples. The **supplementary video** contains qualitative visualizations of the scenarios generated via SLEDGE and of the increased failure rate of PDM-Closed when tested on these scenarios.

1 Additional Methodological Details

In this section, we provide detailed information regarding the architectures and methodologies introduced in our work.

1.1 Raster-to-Vector Autoencoder

Encoder. The Raster-to-Vector Autoencoder (RVAE) applies a ResNet-50 as its encoder π [7]. This is pretrained on ImageNet [5], with a modified first convolution layer to account for the 12-channel Rasterized Scene Image (RSI) input (denoted as \mathbf{x}_{RSI}). We apply group normalization (with 32 channels) and a 3×3 convolution layer on the ResNet output to extract an 8×8 feature grid with 128 channels that is split into two chunks for the mean $\pi_{\mu}(\mathbf{x}_{\text{RSI}})$ and variance $\pi_{\sigma}(\mathbf{x}_{\text{RSI}})$ of the VAE. After reparameterization [10], the encoder outputs a $8 \times 8 \times 64$ tensor, which we refer to as the Rasterized Latent Map (RLM).

Decoder. The transformer decoder ϕ applies 6 cross-attention based transformer decoder layers with a dimensionality of $d_{\text{model}} = 512$ and $d_{\text{ffn}} = 1024$ for the Feed-Forward-Networks (FFNs), and no dropout. Each grid cell of the latent vector is linearly projected to d_{model} and functions as keys and values in the transformer. The queries are learnable parameters corresponding to each entity in the vector output.

Prediction Heads. For each query output, we apply an FFN with a single hidden layer of size $d_{\text{ffn}} = 1024$ to output the attributes of the entity. The FFN outputs a 20×2 set of point coordinates for polylines (lanes and traffic lights) or a descriptor with 6 attributes (2D position, orientation, 2D extent, and speed)

for vehicles and pedestrians. Similarly, static objects have 5 attributes, omitting the speed scalar. We use tanh activation for 2D positions and orientation and re-scale the attributes to $[-32, 32]^2$ and $[-\pi, \pi]$, respectively. The ego velocity $\mathbf{v} \in \mathbb{R}^2$ has two attributes for the longitudinal and lateral axis. The existence attribute $p \in [0, 1]$ is predicted with a linear layer based on the entity queries with variable count.

Training. During training, we employ a bipartite matching loss that uniquely pairs each ground truth entity (indexed by i) with a corresponding entity $\sigma(i)$. For polylines, we define the loss

$$\mathcal{L}_{\text{line}} = \sum_{i \in \{\mathcal{L}, \mathcal{R}, \mathcal{G}\}} \left(\lambda_{\text{line}}^p \cdot \mathcal{L}_{\text{CE}}(p_i, \hat{p}_{\sigma(i)}) + \lambda_{\text{line}}^{\text{attr}} \cdot \mathbb{1}_{[p_i > 0]} \cdot \|\mathbf{L}_i - \hat{\mathbf{L}}_{\sigma(i)}\|_1 \right), \quad (1)$$

where $\mathcal{L}_{\text{CE}}(\cdot, \cdot)$ denotes the cross-entropy loss, p , and \hat{p} represent the actual and predicted existence probabilities of a polyline. \mathbf{L} and $\hat{\mathbf{L}}$ are the actual and predicted polyline coordinates, both as a $\mathbb{R}^{20 \times 2}$ matrix. The terms λ_{line}^p and $\lambda_{\text{line}}^{\text{attr}}$ are hyperparameters, tuned to balance the components of the loss, set to 10 and 4, respectively. We similarly define the loss for all predicted bounding boxes as

$$\mathcal{L}_{\text{box}} = \sum_{i \in \{\mathcal{V}, \mathcal{P}, \mathcal{O}\}} \left(\lambda_{\text{box}}^p \cdot \mathcal{L}_{\text{CE}}(p_i, \hat{p}_{\sigma(i)}) + \lambda_{\text{box}}^{\text{attr}} \cdot \mathbb{1}_{[p_i > 0]} \cdot \|\mathbf{b}_i - \hat{\mathbf{b}}_{\sigma(i)}\|_1 \right), \quad (2)$$

where \mathbf{b} and $\hat{\mathbf{b}}$ denote the actual and predicted bounding box attributes. We set the hyperparameters λ_{box}^p and $\lambda_{\text{box}}^{\text{attr}}$ to 5 and 1, respectively. We found this simpler loss to be sufficient, unlike the complex IoU-based losses used in prior work on object detection that mostly deals with axis-aligned bounding boxes [3]. For the velocity vector of the ego vehicle, we define a L_1 -Loss

$$\mathcal{L}_{\text{ego}} = \|\mathbf{v} - \hat{\mathbf{v}}\|_1, \quad (3)$$

where \mathbf{v} and $\hat{\mathbf{v}}$ represent the actual and the predicted velocity, both in \mathbb{R}^2 . Finally, we regularize the RLM with a weighted Kullback-Leibler-term following a standard VAE [10]

$$\mathcal{L}_{\text{KL}} = \left(-\frac{1}{2} \sum 1 + \log(\pi_{\sigma}(\mathbf{x}_{\text{RSI}})^2) - \pi_{\mu}(\mathbf{x}_{\text{RSI}})^2 - \pi_{\sigma}(\mathbf{x}_{\text{RSI}})^2 \right) \cdot \lambda_{\text{KL}}, \quad (4)$$

where λ_{KL} is set to 0.1. The total loss is given by

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{line}} + \mathcal{L}_{\text{box}} + \mathcal{L}_{\text{ego}} + \mathcal{L}_{\text{KL}}. \quad (5)$$

We compute the matching with the Hungarian algorithm [11] for each class type in $\{\mathcal{L}, \mathcal{R}, \mathcal{G}, \mathcal{V}, \mathcal{P}, \mathcal{O}\}$ independently. For calculating the pairwise cost of entities during matching, we use $\mathcal{L}_{\text{line}}$ from Eq. (1) for polylines and a modified version of \mathcal{L}_{box} from Eq. (2) for bounding boxes. Specifically, we only consider the L_1 cost between the bounding box positions while ignoring the remaining attributes

Table 1: DiT Configurations. We experiment with 3 different model sizes.

Model	#Layers	Hidden size	#Heads	#Params	batch/GPU	GPUh
DiT-B	12	768	12	138M	512	96
DiT-L	24	1024	16	487M	192	96
DiT-XL	28	1152	16	718M	128	960

to prevent undesired matching of distant agents based on similar velocities or 2D bounding box extents.

Implementation. We train the RVAE using tuPlan Garage¹. Training takes roughly 24 hours distributed over 4 A100s. We train the RVAE with a learning rate of $5e^{-5}$ for 40 epochs and divide the learning rate by 10 after 35 epochs. Moreover, we use a batch size of 32 per GPU and AdamW as the optimizer with a weight decay of $5e^{-3}$. We employ 30 lane queries, 10 queries each for red and green traffic lights, 30 vehicle queries, 10 pedestrian queries, 20 queries for static objects, and 1 query for the ego velocity. If more entities are in the scene than available queries, we only consider the nearest entities to the ego center as targets. We discard ground-truth vehicles and static objects which are not on the driveable area. We use a threshold of 0.3 for filtering all existence attributes $p \in [0, 1]$ during inference.

1.2 Diffusion Transformer

Architectures. We summarize the DiT-B, DiT-L, and DiT-XL configurations of our work in Table 1. We follow the architectures of [16], where we leave the patch size p constant to 1, resulting in 8×8 input tokens from the RLM.

Implementation. We train the DiT with Diffusers². We trained DiT-B and DiT-L for 1 day on 4 A100s and our final DiT-XL model for 5 days on 8 A100s. All models were trained with an AdamW optimizer, a weight decay of $1e^{-6}$, an exponential moving average (EMA) decay of 0.9999, and a constant learning rate of $1e^{-4}$. We use the DDPM scheduler [9], with a $t_{\max} = 1000$ linear variance schedule from 0.0015 to 0.015. During inference, we exclusively apply 100 denoising steps with a classifier-free guidance scale of 4.0.

1.3 Route Extrapolation

Initial Pose. In the initial step for route extrapolation, we generate a $64m \times 64m$ map with agents without masking and only conditioned on the map ID label. Then, we (a) determine the nearest lane from the current position (initially at $(0, 0)$) and retrieve all available paths using DFS. We prioritize paths that reach

¹ https://github.com/autonomousvision/tuplan_garage

² <https://github.com/huggingface/diffusers>

the tile ends of the 64m×64m map and select a route according to the difficulty (minimizing or maximizing turns).

RLM Proposals. Next, we **(b)** transform the previously predicted map poly-lines into the local coordinate frame of route end pose and rasterize the map as RSI. We forward the RSI in our encoder π , obtaining the $8 \times 8 \times 32$ latent grid of the previously generated map. In **(c)**, we sample $N = 8$ random $8 \times 8 \times 32$ RLMs from $\mathcal{N}(0, \mathbf{I})$, which are to be denoised.

Inpainting. We concatenate the $4 \times 8 \times 32$ region of the encoded latent, corresponding to half a frame, where $x \in [-32, 0)$ (behind the ego vehicle), with the $4 \times 8 \times 32$ region of the random RLMs, corresponding to half a frame, where $x \in (0, 32]$ (ahead of the ego vehicle). We then **(d)** apply 100 denoising steps (conditioned on the map label with a guidance scale of 4.0) only on the randomly initialized sections of the RLM. We found this strategy slightly better in enforcing map coherency than simultaneously denoising masked and unmasked sections of the RLMs.

Proposal Scoring. After decoding the RLM’s via ϕ , we **(e)** match the lane poly-lines for all N new samples with the previous map on the overlapping frame section (where $x \in [-32, 0)$) and assign a cost for each new proposal. The cost is based on average point-wise distances of lines (clipped to 3m), where we add a fixed cost of 3m if lines of the current or proposed graph remain unmatched. We **(f)** select the proposal with minimum cost, extract the lanes and agents in the newly generated section (i.e. where $x \in [0, 32]$), and convert the new tile in global coordinates of the initial frame. We repeat step (a)-(f), starting from the current route end pose, until the desired route length is reached.

Limitations. The route extrapolation has several modes of error. If no sufficient path is found during DFS expansion (i.e., due to errors in connectivity), the newly generated tile overlaps with the existing map. Specifically, tile edges that cross intersections are difficult for DiT-XL to inpaint. Despite these challenges, the DiT-XL model can generate large maps autoregressively, enabling long simulations, which we view as a foundation for future work.

1.4 Lane Graph Representations

Vector Lane Representation. The native lanes in nuPlan are arranged into groups of parallel blocks. This results in many individual poly-lines with only a single adjacent line in the forward or backward direction. Our lane representation deviates slightly from nuPlan, where consecutive lane nodes are summarized into one individual line, as shown in Fig. 1. We first collect all nodes with incoming degree unequal 1 (called start nodes) or outgoing degree unequal 1 (called end nodes). We further extend the start node set with outgoing lanes of end nodes. Next, we apply Depth-First-Search (DFS) on all start nodes until an end node or start node is reached. If a new path was found, we add the polyline to our line set. We found this strategy effective, as it requires about 40% fewer lane queries and simplifies inferring the lane adjacency.

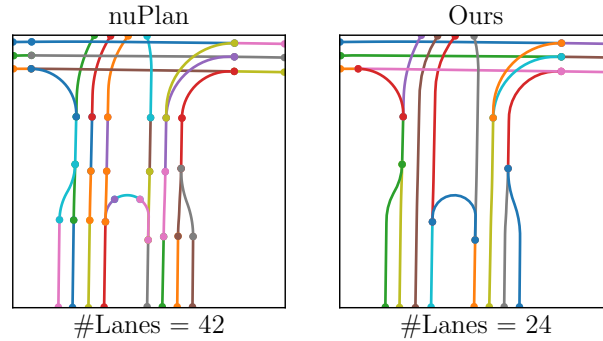


Fig. 1: Lane Graphs. We use a more compact lane graph representation for generation and simulation by combining adjacent lines that only permit a single traversable path.

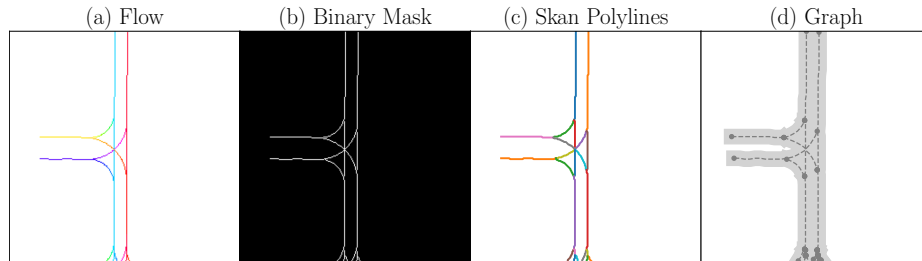


Fig. 2: Skan. Our pipeline uses functions from established image processing libraries to skeletonize and extract a lane graph from the RSI.

Skan for Lane Graph Extraction from RSI. Extracting lane graphs from raster images is challenging, as it requires inferring precise polyline locations and connectivity from thin structures in a dense pixel grid. Our pipeline, involving multiple processing steps, is shown in Fig. 2. First, we create a binary lane mask (using a threshold of 0.1) and apply a skeletonization algorithm [19] to gather a one-pixel wide line structure. Using the skan library [15], we extract a sequence of pixel coordinates for individual lines in the mask, which we align to the flow direction and convert into the ego coordinate frame in $[-32, 32]^2$. Finally, following the heuristic applied in our proposed approach, we connect polylines if the start- and end-poses have less than 1.5m displacement and differ less than 60 degrees in orientation.

Reconstruction Metrics. Given the varying formulations of lane graphs used in nuPlan, our models, and the baselines, we define a uniform representation as a foundation for all graph-based metrics. We interpolate along all polylines while inferring the orientation to extract poses at fixed 1.5m intervals. The poses are connected along polylines and between lane adjacencies at start- and endposes. We assign predicted poses to ground-truth poses using a Hungarian matching algorithm [11], of poses which differ less than 60 degrees in orientation and 1.5m in displacement. Given the assignment, we calculate (1) **F1** as the harmonic mean of the precision and recall and (2) **Lateral L2 (Lat.)** as the lateral offset

between matched pairs in the Euclidean distance. We introduce the latter metric to disentangle longitudinal interpolation mismatches from structural errors in the predicted lanes. Next, we measure the **(3) Chamfer** distance using all poses by averaging the squared distance between nearest-neighbor pairs between both point sets and summing the averages. Following literature for aerial lane extraction [8], we calculate the metrics in two settings. For the **(1) GEO** scores, we utilize the complete graph in the $64\text{m} \times 64\text{m}$ frames. The **(2) TOPO** scores compute the reachability graph of every tenth node and up to a radius of 50m. We compute F1, Lateral L2, and Chamfer for each sub-graph and average the metrics. Consequently, connectivity errors can lead to substantial deductions in TOPO results.

1.5 Lane Graph Generation

VAE (RSI). We follow the encoder architecture as described in Section 1.1 with two channels for the lane flow. We input the $256 \times 256 \times 2$ RSI and use the same latent map shape as our RVAE. The decoder applies 5 up-convolutions, each followed by batch normalization and a ReLU activation. With a last convolution layer, the model reconstructs the two-channel flow image. We train the VAE with a binary cross entropy loss (given the lane channels are normalized in $[0, 1]$), while adding a KL term weighted by $\lambda_{\text{KL}} = 1.0$. We use the same training setting from our RVAE (see Section 1.1), including optimizer, learning rate schedules, and number of epochs. During inference, we extract the vectorized lane graph according to Section 1.4.

HMapGen. Our re-implementation of HMapGen [14] is trained using the GRAN codebase³. The graph representation is similar to our lane formulation in Section 1.4, where the locations of global key points are placed on the start, end, or transitioning positions of polylines. We further formulate local nodes as 20 equispaced points interpolated along the polylines (i.e. connections of global key points). The autoregressive model predicts the adjacency matrix and node coordinate at each step. The adjacency matrix is supervised with the loss from [13]. The coordinate predictions are Gaussian mixtures with 20 components, supervised with the Mixture Density Network [1] formulation. We train this reimplementation of HMapGen with a batch size of 512, an Adam optimizer, and a learning rate of $5e^{-4}$ for 25 epochs on a single A100.

DiT-L (RSI). We consider a DiT-L diffusion model applied on the $256 \times 256 \times 2$ lane flow map of the RSI. In order to accommodate the large pixel grid, we use a patch size of 16 for the DiT-L tokenization. We clip the predicted samples to $[0, 1]$ for the DDPM scheduler and use a linear variance schedule from 0.02 to 0.0001 following [6]. We train DiT-L on 4 A100s for two days (approximately the same compute of training the RVAE and DiT-L RLM variant) with a batch size of 48. The remaining parameters for training and inference of DiT-L are according to Section 1.2.

³ <https://github.com/lrjconan/GRAN>

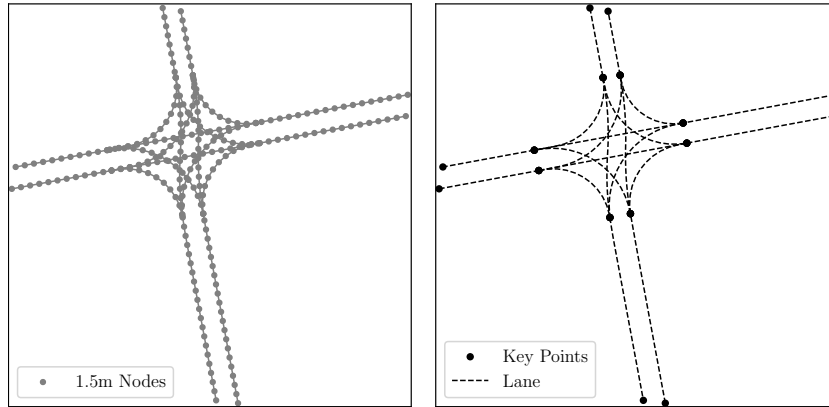


Fig. 3: Key Points and Lanes used for Frechet Metrics. We rasterize the graph generated with the key points and lanes into an RSI for the CNN metrics, and use the extracted key points for the urban planning metrics.

Generative Metrics. For the assessment of generation performance, we use the same uniform graph representation of the reconstruction metrics in Section 1.4. Consequently, we extract a lane representation where nodes with degree unequal 1 are marked as start or endpoints of polylines. As shown in Fig. 3, we extract individual lines with an analogous logic based on DFS expansions of start points (see Section 1.4), where we obtain the connectivity between polylines from the adjacency of the 1.5m nodes. Thereby, we can formulate generative metrics on lane polylines and key points based on the uniform graph. For the learned Frechet metrics, we rasterize the lane polylines as an RSI and fill the remaining flow channels with 0. For **ImNet** and **RVEnc**, we encode the RSI (either 3 or 12 color channels) with a ResNet-50 encoder (trained on ImageNet) or the encoder π of our RVAE. We then calculate the ImNet and RVEnc Frechet distances (scaled by 10^3 and 10^2 , respectively). Moreover, we follow the Frechet distances in [14], namely **Connectivity** as degree of key points (scaled by 10), **Density** as number of key points (not scaled), **Reach** number available paths from each key point (not scaled), and **Convenience** defined as length of available paths from each key point (scaled by 10). We additionally search the longest path with DFS in the uniform graph, starting from the node pose closest to (0,0), and report average length and standard deviation.

1.6 SLEDGE Simulation of PDM-Closed

General. We use the simulation pipeline of nuPlan [2], where the output trajectory of a planner is simulated with a Linear Quadratic Regulator controller [12] and a kinematic bicycle model [17]. All simulations have a step frequency of 10 Hz. As the planner under test, we use PDM-Closed, with the exact configuration proposed in [4]. During the ‘Replay’ and ‘Lane→Agent’ simulation, the planner has access to nuPlan’s map interface. We created a custom map interface for the

generated lane graph in ‘Lane & Agent’ simulations. We set the speed limit of generated lanes to 15m/s.

Reactive Agents. All reactive vehicles are projected onto the nearest lane polyline and simulated with the Intelligent Driver Model [18]. During initialization, we discard vehicles if the bounding box overlaps with other agents or is off the driveable area. However, overlaps may occur over the course of a simulation. Each vehicle iteratively plans a path, following subsequent lanes with minimum curvature. We use the same IDM parameters of nuPlan [2] and extend the minimum future path length to 30m to facilitate long simulations. If a vehicle is within $\alpha = 64m$ of distance to the ego center, the future path of the vehicle is updated (up to 30m), and the vehicle state is propagated according to IDM. We simulate pedestrians with a constant velocity model, which can lead to uncommon movement during long simulation periods. Thus, we set the simulation radius of pedestrians to 10m.

Traffic Lights. For the ‘Replay’ simulation, we load the traffic light state of each lane at the first occurrence. During the ‘Lane→Agent’ and ‘Lane & Agent’ simulation, we match the generated traffic light polyline to the nearest lane based on an average point-wise Euclidean distance. We flip the traffic lights at 15s intervals to allow the planner to progress naturally. If the nuPlan’s map interface is available, we further set all traffic lights on an intersection to red when a lane does not follow the target route.

Simulation Metrics. The Planner Failure Rate (PFR) is based on nuPlan’s Closed Loop Score (CLS) metric. Specifically, we consider failures as infractions in multiplier metrics of nuPlan’s CLS based on the same implementation. These include at-fault collisions, driving off-road, driving in opposite traffic direction for more than 6m, or staying below 20% of the target progress. The original nuPlan closed loop score metric assigns a large multiplicative penalty for these violations. Since both our planner and agent simulation logic are deterministic, the metrics computed do not vary with different evaluation runs. The metrics are computed after the simulation is complete. For more details, please see [4].

Termination. We set the simulation duration to 30s and 150s for route lengths of 100m and 500m, respectively. Restricting the simulation length allows us to compute the nuPlan metrics without adding additional heuristics (i.e. to determine when a planner stops progressing due to mistakes). In practice, we find that a planner has sufficient time to reach the threshold of 20% progress to prevent a failure with these selected durations.

2 Additional Results

In this section, we provide additional quantitative and qualitative results.

RVAE Ablation Study. As shown in Table 2, We observe a small drop in performance with a larger encoder. Modifying the RLM size, regularization, and decoder architecture had no significant impact.

Table 2: RVAE Ablation Study. The default configuration uses a ResNet-50 encoder, $8 \times 8 \times 64$ latent with $\lambda_{\text{KL}} = 0.1$, channel group masking and 6 transformer decoder layers.

Experiment	Config	GEO			TOPO		
		F1 \uparrow	Lat. \downarrow	Ch. \downarrow	F1 \uparrow	Lat. \downarrow	Ch. \downarrow
Encoder	ResNet-101	0.977	0.172	0.477	0.936	0.302	23.865
RLM	$8 \times 8 \times 128$	0.983	0.158	0.369	0.948	0.276	19.571
	$32 \times 32 \times 4$	0.978	0.178	0.489	0.941	0.301	21.460
	$\lambda_{\text{KL}} = 1.0$	0.979	0.188	0.446	0.941	0.318	21.499
Decoder	No Masking	0.981	0.161	0.399	0.945	0.282	20.096
	3 Layers	0.979	0.173	0.425	0.938	0.303	23.789
-	Default	0.980	0.164	0.411	0.944	0.288	20.624

Table 3: $64\text{m} \times 64\text{m}$ Lane Graph Generation. Frechet distances between samples from the validation set and model outputs, based on penultimate feature vectors of ResNet-50 encoders. We consider variants pre-trained on ImageNet (ImNet) and the encoder of the RVAE (RVEnc). *Trained with $\sim 6\times$ more compute than others, which already converge at the lower compute budget.

Architecture	Representation	Frechet (CNN) \downarrow	
		ImNet	RVEnc
VAE	RSI	1.77	60.61
RVAE	Vector	1.38	28.94
HDMaGen	Vector	0.97	23.90
DiT-L	RSI	1.00	25.22
	RLM	0.66	10.97
DiT-XL*	RLM	0.64	9.37

Generative Metrics. In Table 3, we further present the Frechet distances based on the encoded ground truth and generated lane graphs, with two ResNet-50 encoder. We consider a CNN pre-trained on ImageNet and from the RVAE. DiT-L as latent diffusion models on the RLM substantially outperforms all methods, which directly operate on raster or vector presentations. Additional compute resources and model size in DiT-XL further improves generative quality.

Uncurated Random Samples. We visualize more generated lane graphs for all approaches in our study (Fig. 4) and additional qualitative results of initial scene states generated by our final DiT-XL model (Fig. 5).

References

1. Bishop, C.M.: Mixture density networks (1994)

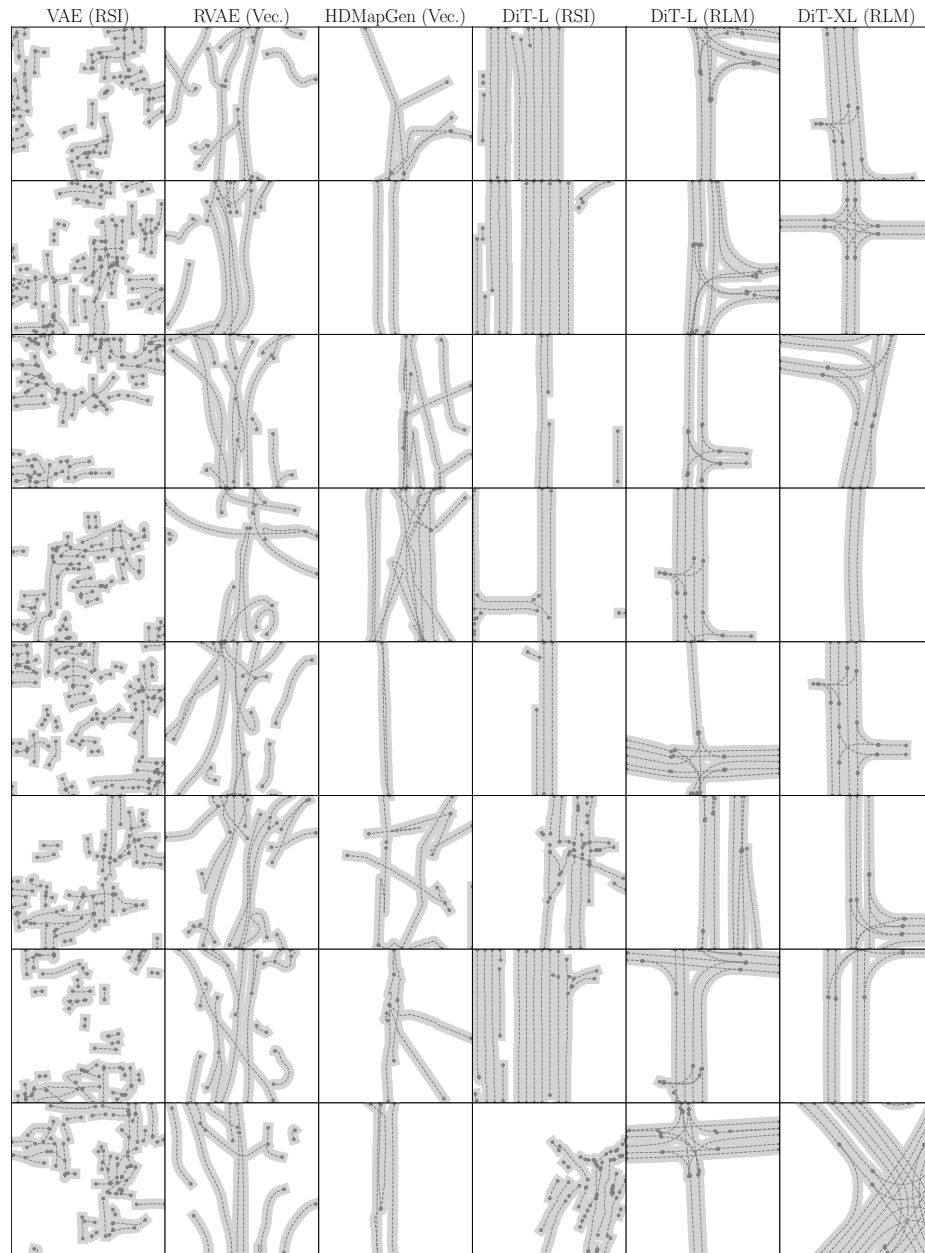


Fig. 4: Uncurated Random Samples. We visualize lane graphs generated by the models from our study, in addition to the samples shown in the main paper.

- Caesar, H., Kabzan, J., Tan, K.S., Fong, W.K., Wolff, E.M., Lang, A.H., Fletcher, L., Beijbom, O., Omari, S.: nuplan: A closed-loop ml-based planning benchmark for autonomous vehicles. In: Proc. IEEE Conf. on Computer Vision and Pattern

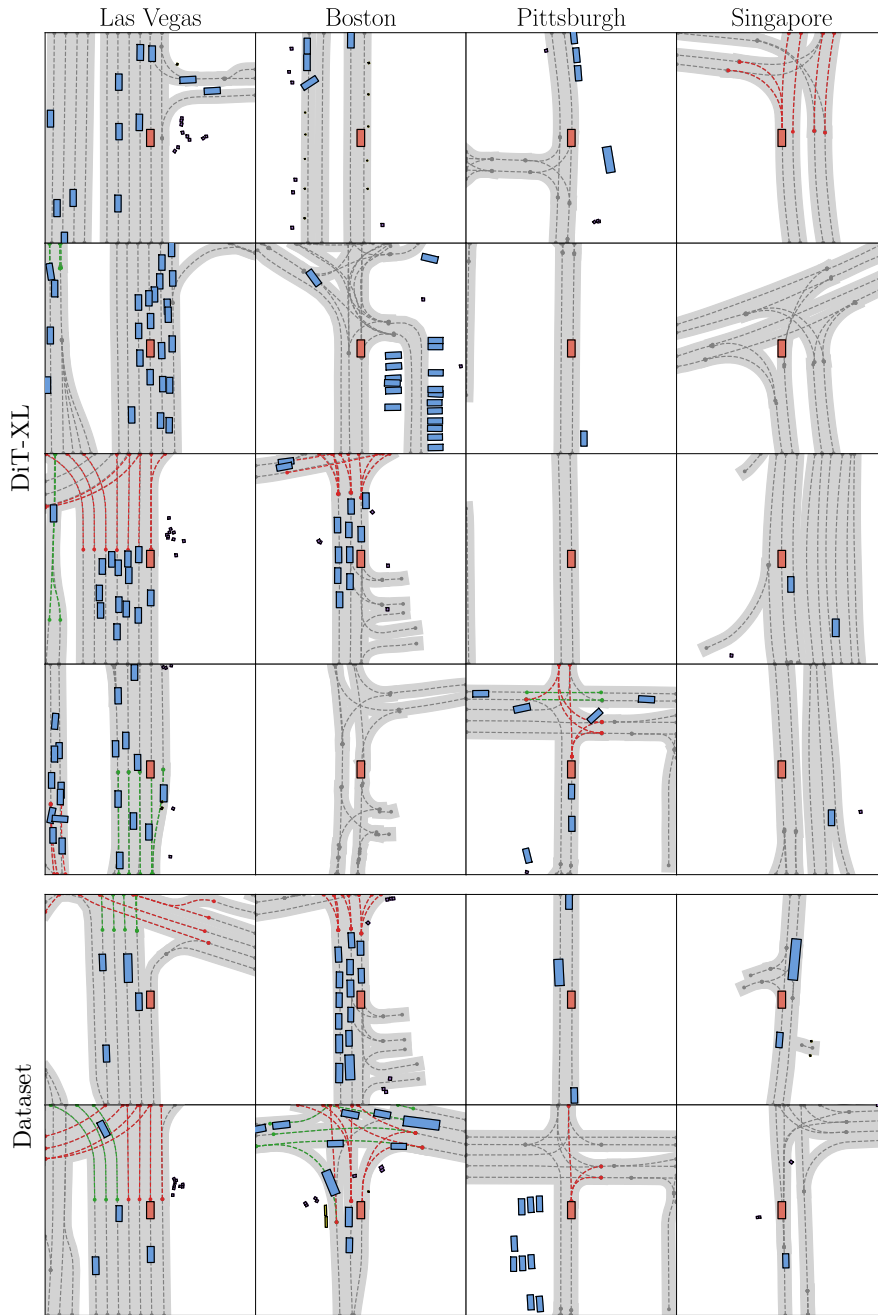


Fig. 5: Uncurated Random Samples. We visualize lanes & agents generated by our final DiT-XL model for 4 cities. Samples from the dataset are included for reference.

- Recognition (CVPR) Workshops (2021)
3. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: Proc. of the European Conf. on Computer Vision (ECCV) (2020)
 4. Dauner, D., Hallgarten, M., Geiger, A., Chitta, K.: Parting with misconceptions about learning-based vehicle motion planning. In: Proc. Conf. on Robot Learning (CoRL) (2023)
 5. Deng, J., Dong, W., Socher, R., jia Li, L., Li, K., Fei-fei, L.: Imagenet: A large-scale hierarchical image database. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2009)
 6. Dhariwal, P., Nichol, A.: Diffusion models beat gans on image synthesis. In: Advances in Neural Information Processing Systems (NeurIPS) (2021)
 7. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2016)
 8. He, S., Balakrishnan, H.: Lane-level street map extraction from aerial imagery. In: Proc. of the IEEE Winter Conference on Applications of Computer Vision (WACV) (2022)
 9. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. arXiv.org **2006.11239** (2020)
 10. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. Proc. of the International Conf. on Learning Representations (ICLR) (2014)
 11. Kuhn, H.W.: The Hungarian method for the assignment problem. *Naval Research Logistics Quarterly* **2**, 83–97 (1955)
 12. Lehtomaki, N., Sandell, N., Athans, M.: Robustness results in linear-quadratic gaussian based multivariable control designs. *IEEE Transactions on Automatic Control* **26** (1981)
 13. Liao, R., Li, Y., Song, Y., Wang, S., Nash, C., Hamilton, W.L., Duvenaud, D., Urtasun, R., Zemel, R.: Efficient graph generation with graph recurrent attention networks. In: Advances in Neural Information Processing Systems (NeurIPS) (2019)
 14. Mi, L., Zhao, H., Nash, C., Jin, X., Gao, J., Sun, C., Schmid, C., Shavit, N., Chai, Y., Anguelov, D.: Hdmapgen: A hierarchical graph generative model of high definition maps. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2021)
 15. Nunez-Iglesias, J., Blanch, A.J., Looker, O., Dixon, M.W.A., Tilley, L.: A new python library to analyse skeleton images confirms malaria parasite remodelling of the red blood cell membrane skeleton. *PeerJ* (2018)
 16. Peebles, W., Xie, S.: Scalable diffusion models with transformers. In: Proc. of the IEEE International Conf. on Computer Vision (ICCV) (2023)
 17. Polack, P., Altché, F., d’Andréa Novel, B., de La Fortelle, A.: The kinematic bicycle model: A consistent model for planning feasible trajectories for autonomous vehicles? In: Proc. IEEE Intelligent Vehicles Symposium (IV) (2017)
 18. Treiber, M., Hennecke, A., Helbing, D.: Congested traffic states in empirical observations and microscopic simulations. *Physical review E* (2000)
 19. Zhang, T.Y., Suen, C.Y.: A fast parallel algorithm for thinning digital patterns. *Communications of the ACM* **27**(3), 236–239 (1984)