

DTU

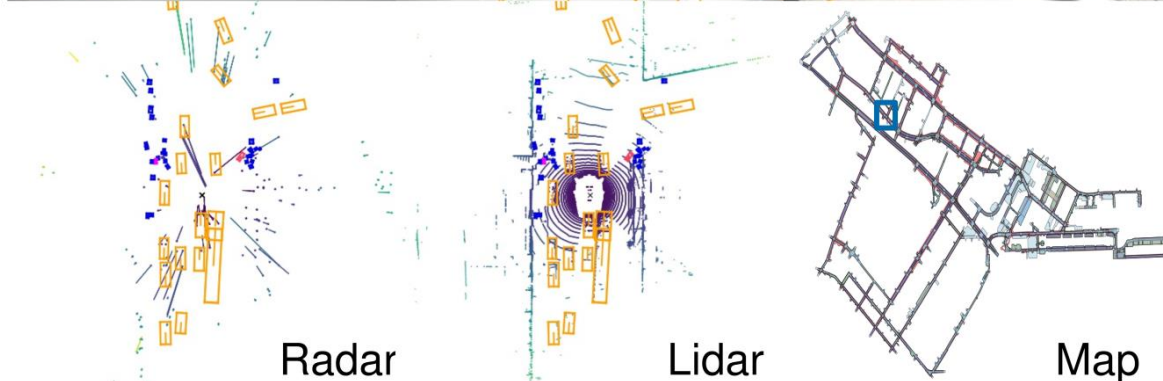
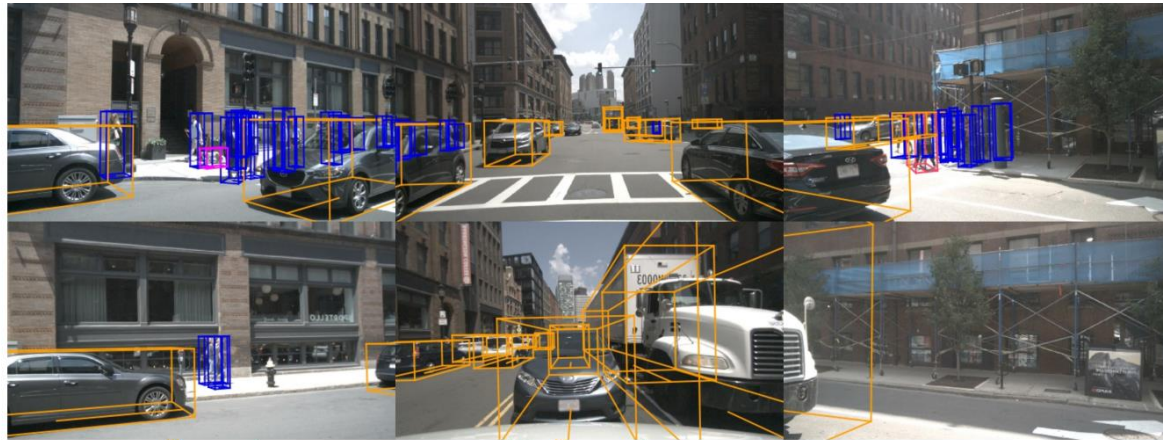


Enhancing Sparse Radar Point Clouds using Diffusion Models  
For Radar-Only 3D Detection

**Hassan Hotait**  
**Msc Thesis**

# Outline

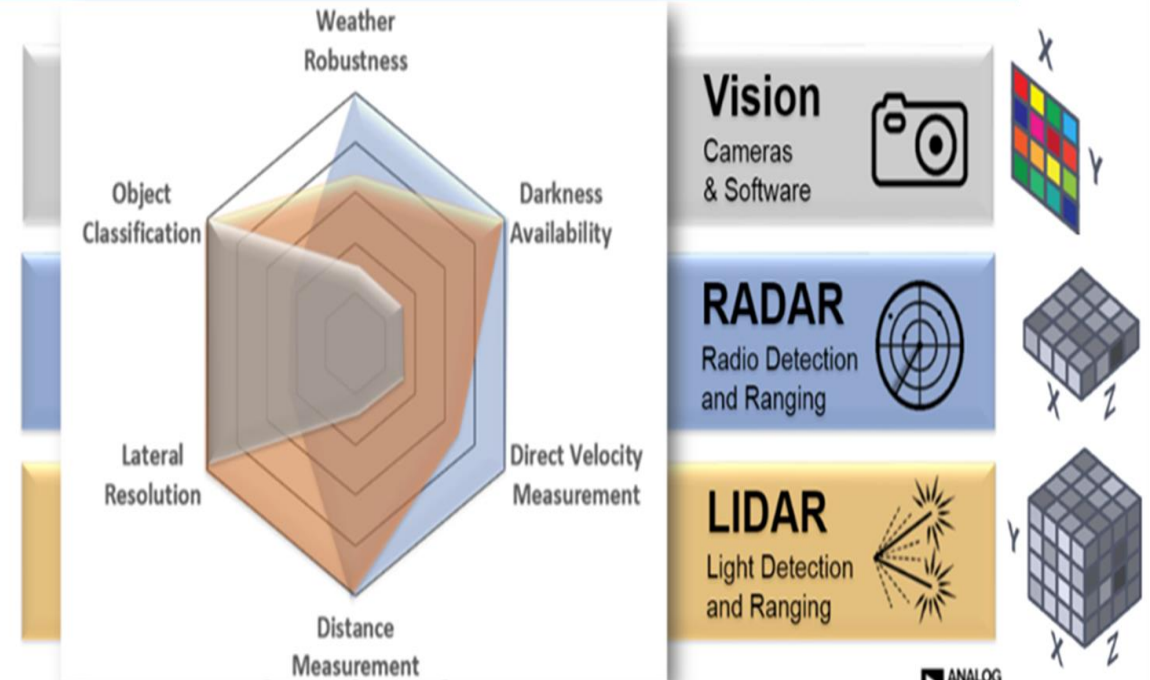
- 3D Detection in Autonomous Driving
- Thesis Scope
- Recap on Image Based Diffusion & Design Choices
- Our Method
- Qualitative & Quantitative Results
- Future Work



"Ped with pet, bicycle, car makes a u-turn, lane change, peds crossing crosswalk"

## Perception Sensors for Autonomous Driving

A single sensor cannot do the entire job!

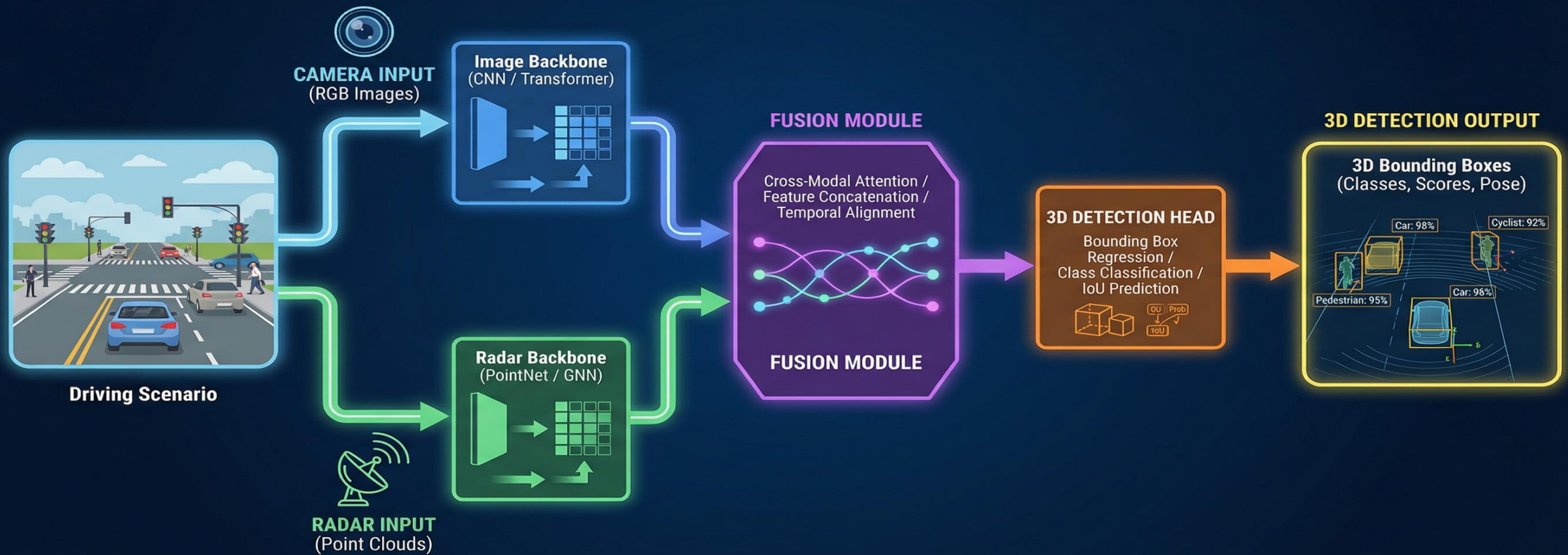


7

\* Source: dSPACE Technology Conference 2017

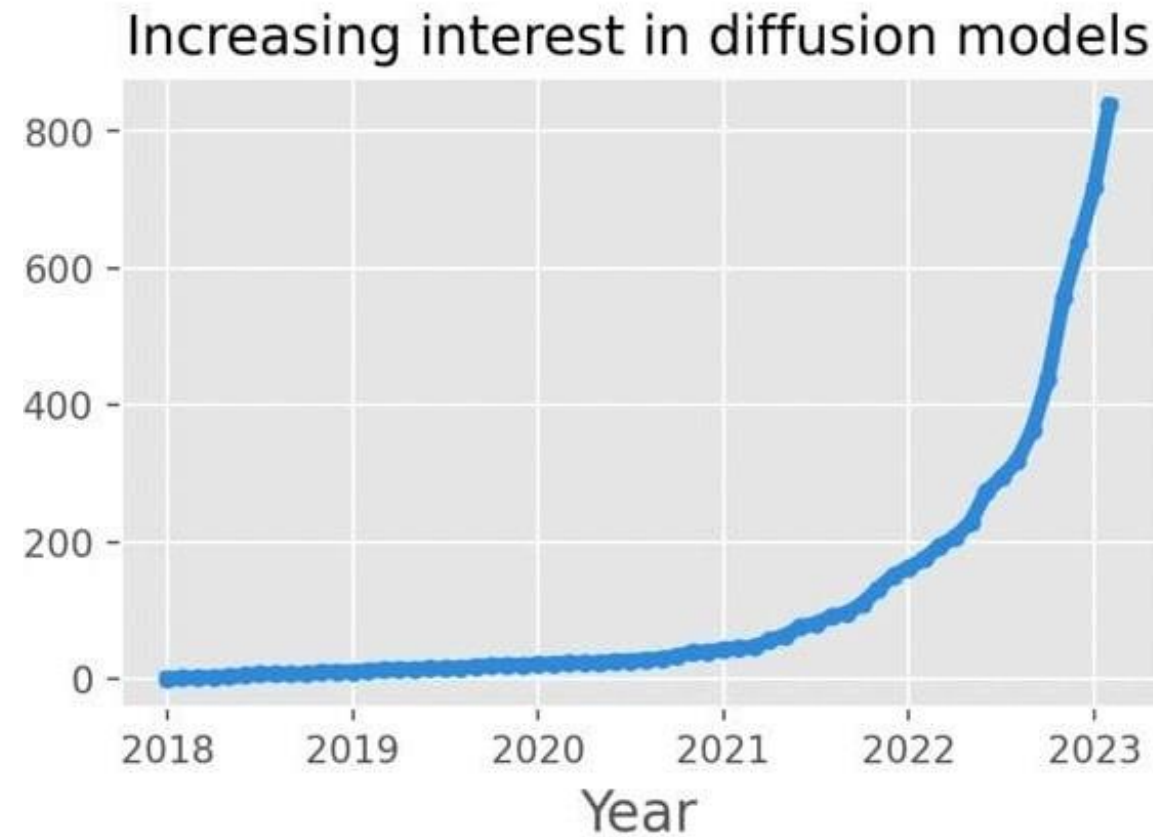


## CAMERA-RADAR FUSION MODEL FOR 3D DETECTION (DRIVING SCENARIOS)



Method						
Date	Name	Modalities	Map data	External data	mAP	
<input type="text"/>		Radar ▾	All ▾	All ▾		
> 2024-03-08	RadarDistill	Radar	no	no	0.205	
> 2022-02-03	KPConvPillars	Radar	no	no	0.049	
> 2020-09-22	Radar-PointGNN	Radar	no	no	0.005	

Method	Input	KD	Car(AP)↑	mAP↑
BEVFusion* [20]	C,R		65.9	38.3
RadarDistill-CR	C,R	✓	<b>67.7</b>	<b>39.6</b>



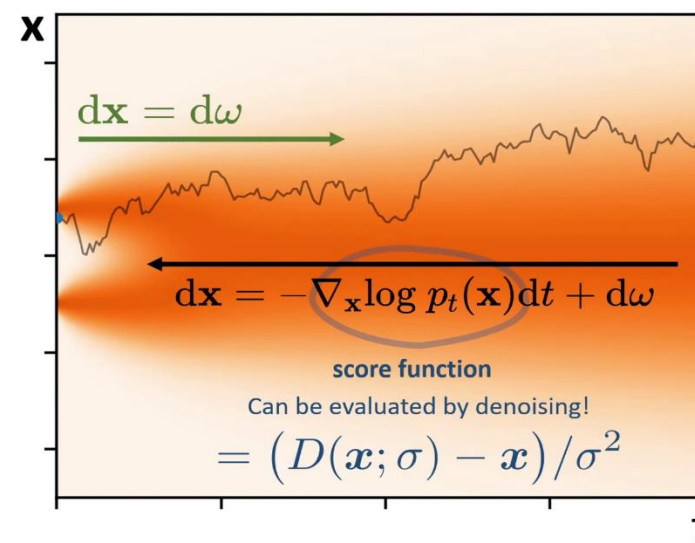
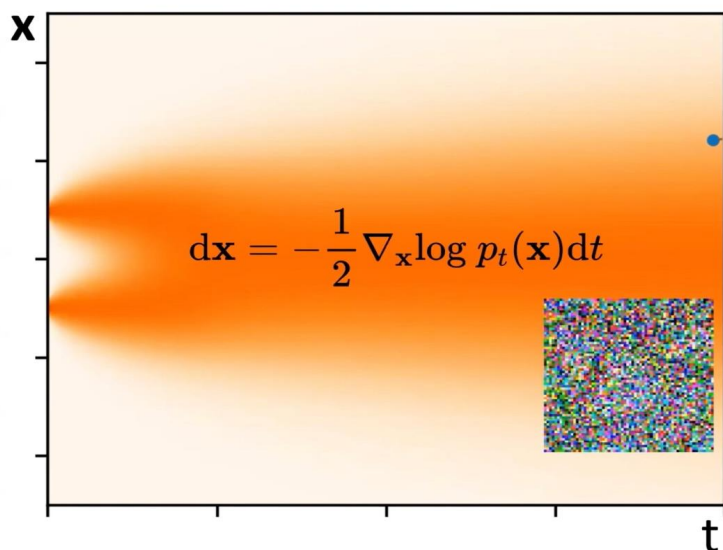
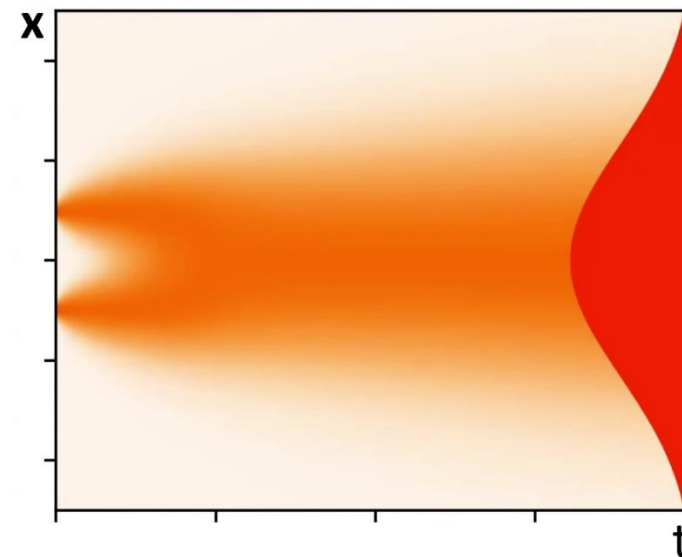
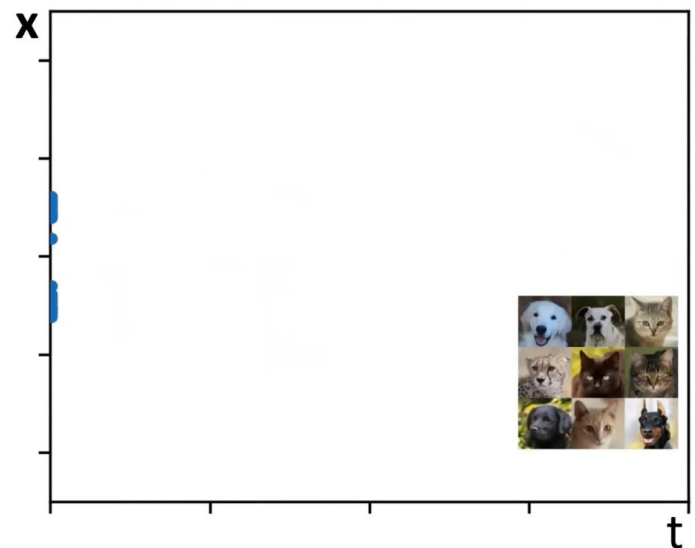
**Thesis Scope:** *Using diffusion models for enhancing sparsity of radar pointclouds for radar-only 3D detection*

# Related Work – Image Based Diffusion Generative Methods

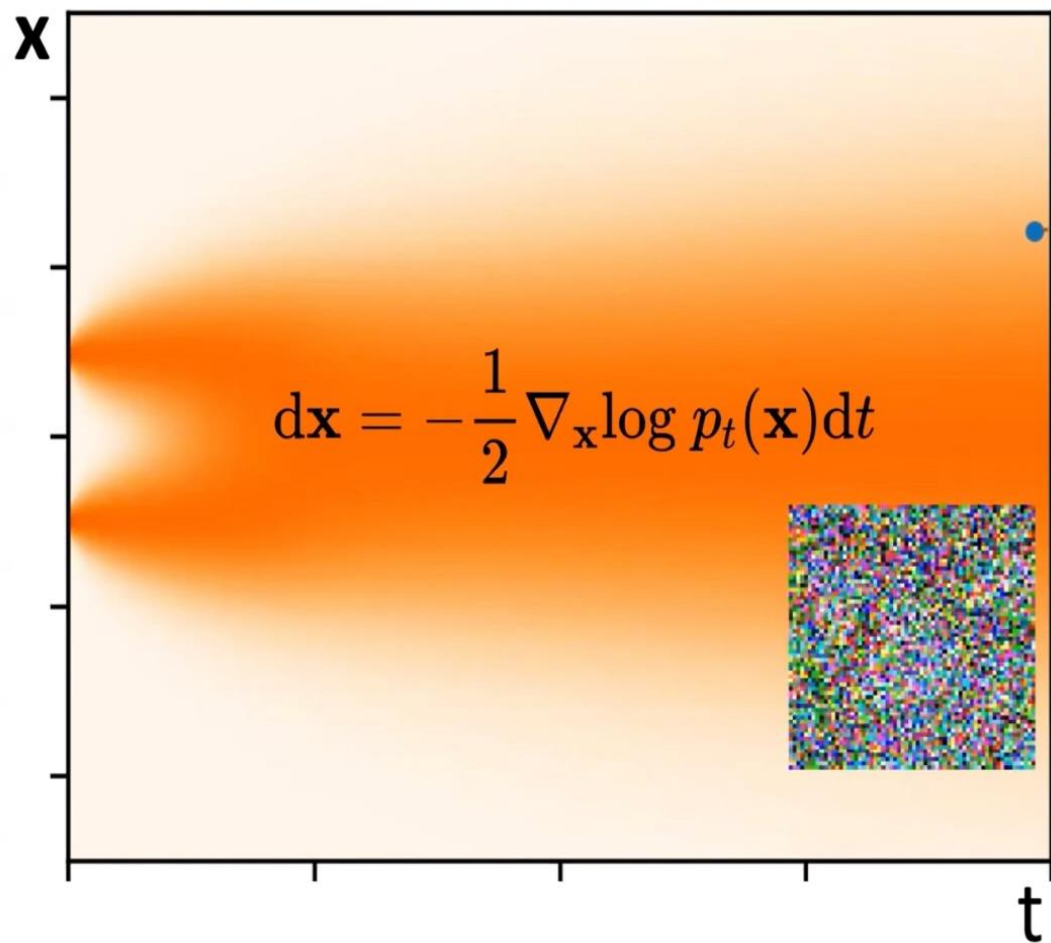
	VP	VE	iDDPM + DDIM
<b>Sampling</b>			
ODE solver	Euler	Euler	Euler
Time steps	$t_{i < N} \quad 1 + \frac{i}{N-1}(\epsilon_s - 1)$	$\sigma_{\max}^2 (\sigma_{\min}^2 / \sigma_{\max}^2)^{\frac{i}{N-1}}$	$u_{j_0 + \lfloor \frac{M-1-j_0}{N-1} i + \frac{1}{2} \rfloor}$ , where $u_M = 0$ $u_{j-1} = \sqrt{\frac{u_j^2 + 1}{\max(\alpha_{j-1}/\alpha_j, C_1)}} - 1$
Schedule	$\sigma(t) \quad \sqrt{e^{\frac{1}{2}\beta_d t^2 + \beta_{\min} t} - 1}$	$\sqrt{t}$	$t$
Scaling	$s(t) \quad 1/\sqrt{e^{\frac{1}{2}\beta_d t^2 + \beta_{\min} t}}$	1	1
<b>Network and preconditioning</b>			
Architecture of $F_\theta$	DDPM++	NCSN++	DDPM
Skip scaling $c_{\text{skip}}(\sigma)$	1	1	1
Output scaling $c_{\text{out}}(\sigma)$	–	$\sigma$	$-\sigma$
Input scaling $c_{\text{in}}(\sigma)$	1	–	$1/\sqrt{\sigma^2 + 1}$
Noise cond. $c_{\text{noise}}(\sigma)$	–	–	$\min_j  u_j - \sigma $
<b>Training</b>			
Noise distribution	–	–	$\mathcal{U}\{0, M-1\}$
Loss weighting $\lambda(\sigma)$	$1/\sigma^2$	–	–
<b>Parameters</b>			
	$\beta_d = 19.9, \beta_{\min} = 0.1$	$\sigma_{\min} = 0.02$	$\frac{j}{2M(C_2+1)}$
	$\epsilon_s = 10^{-3}, \epsilon_t = 10^{-5}$	$\sigma_{\max} = 100$	$C_1 = 0.001, C_2 = 0.008$
	$M = 1000$	–	$M = 1000, j_0 = 8^\dagger$



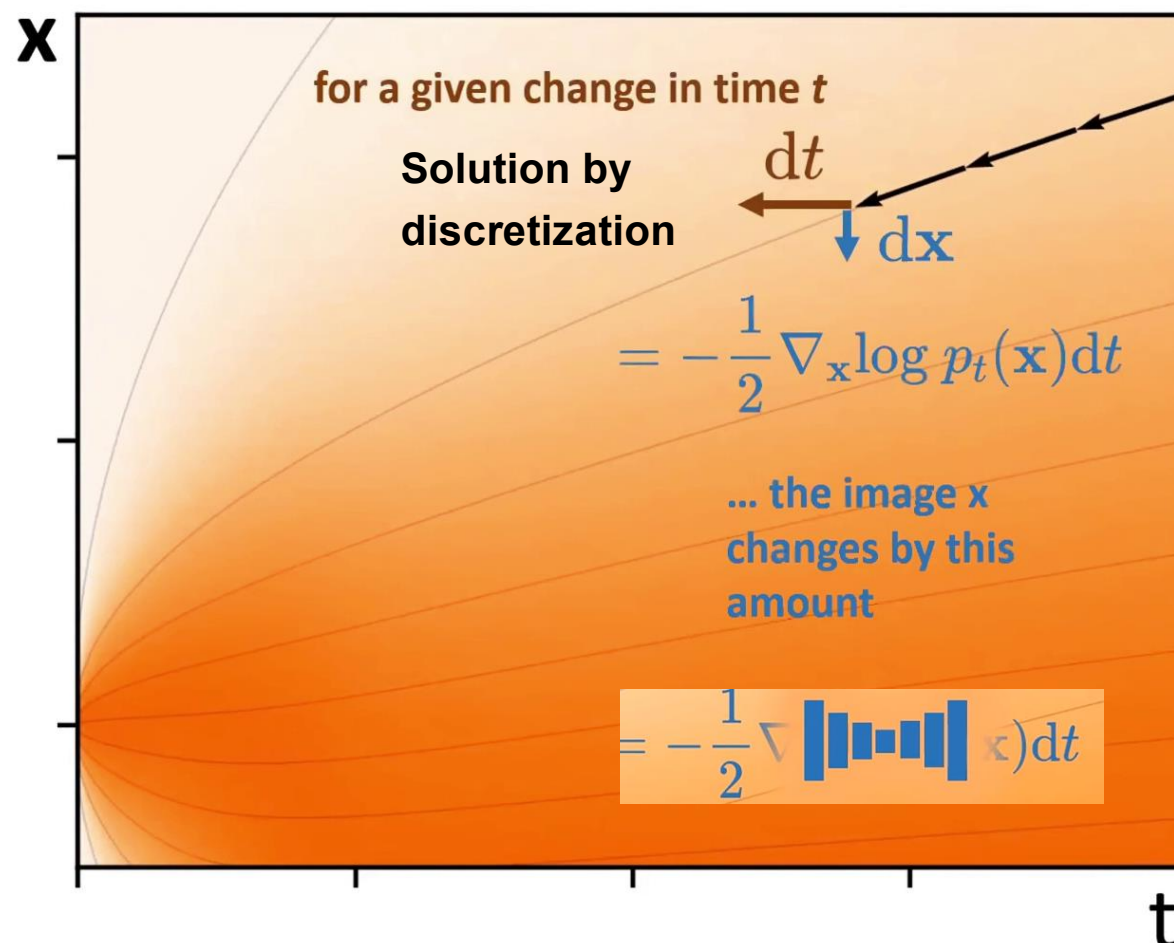
# Related Work – Diffusion Based Generative Methods



# Related Work – Diffusion Based Generative Methods



Forward & Reverse Process Loop



Score Function Predicted by NN

# Related Work – Diffusion Based Generative Methods



## Training

Agnostic of usage as part of ODE/SDE stepping

1. Network fails to approximate the **true score of data**

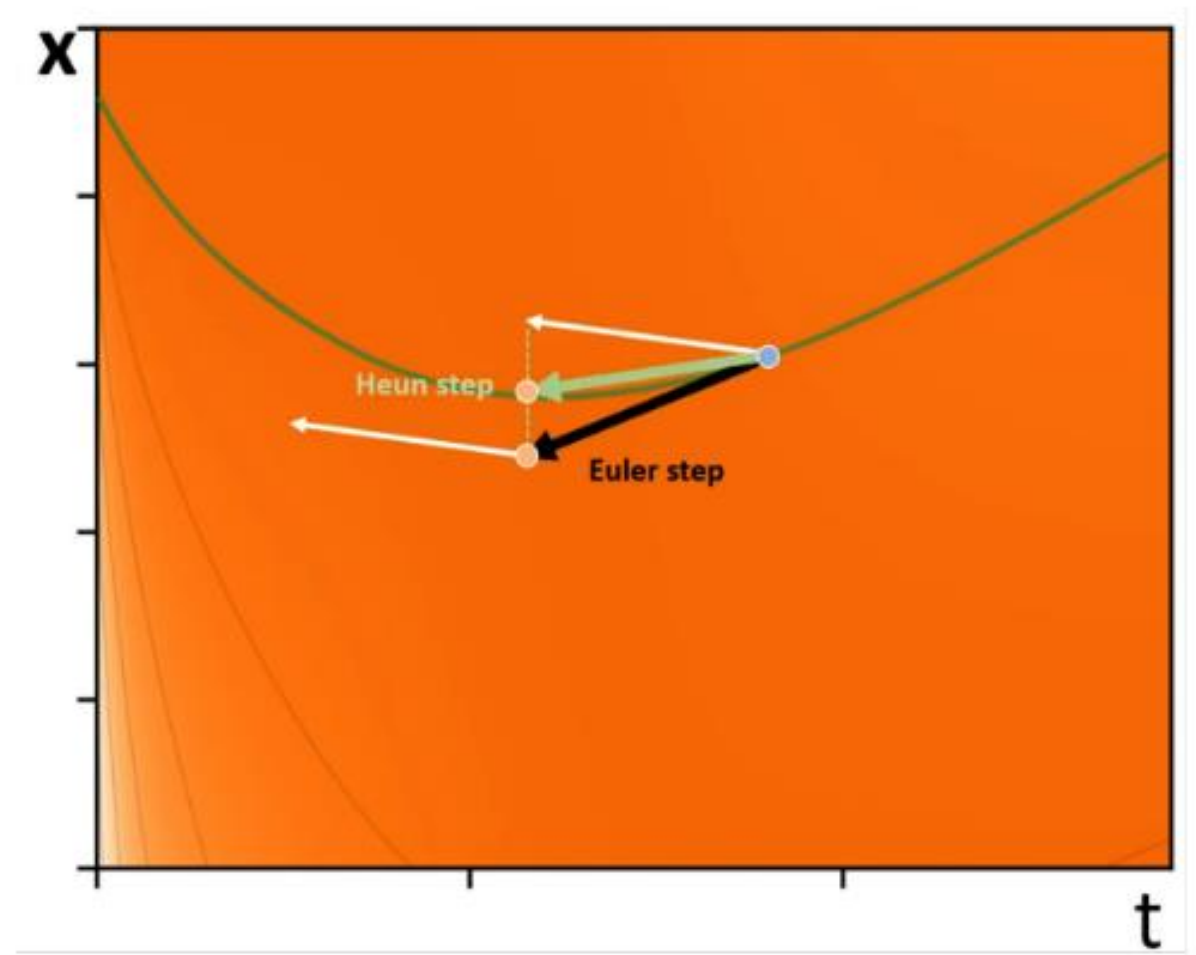
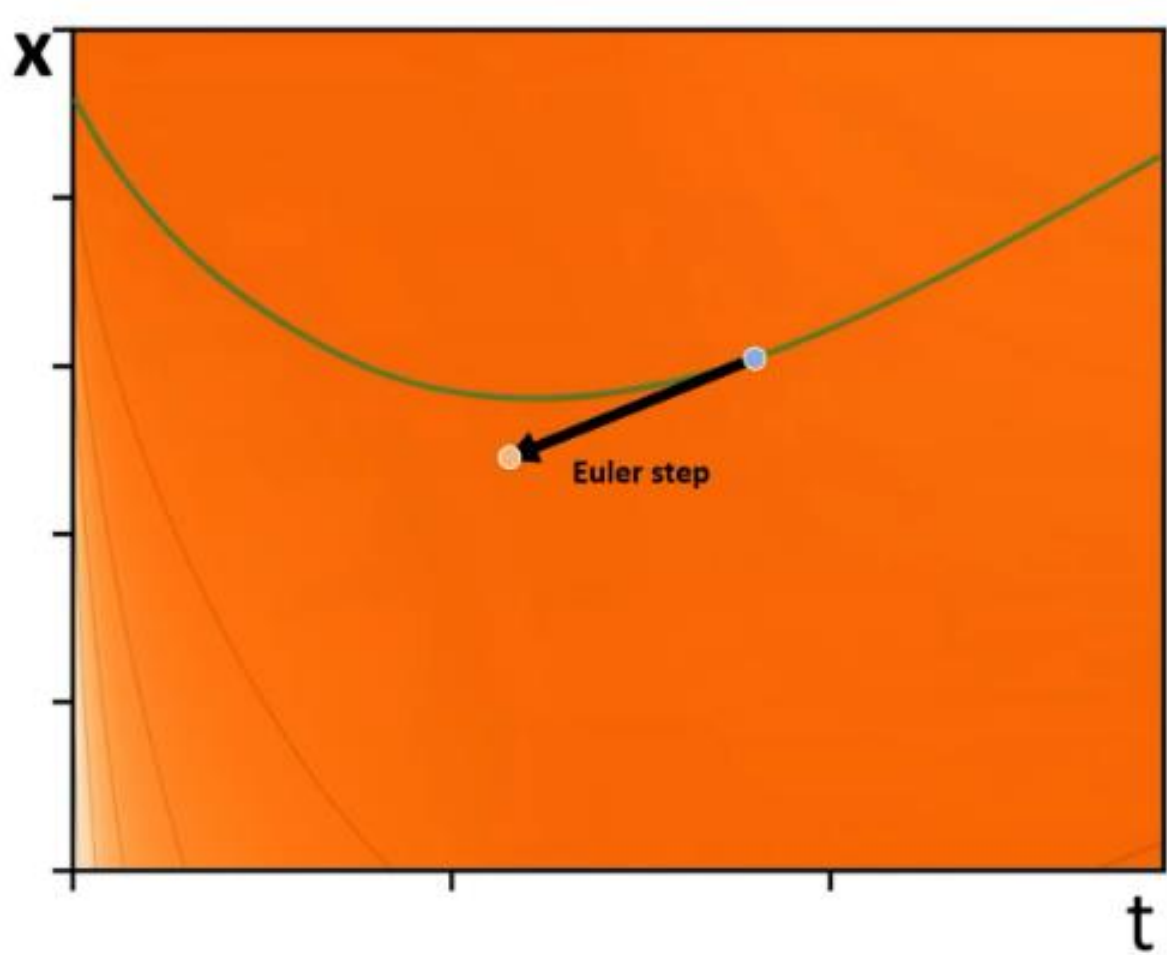


## Sampling

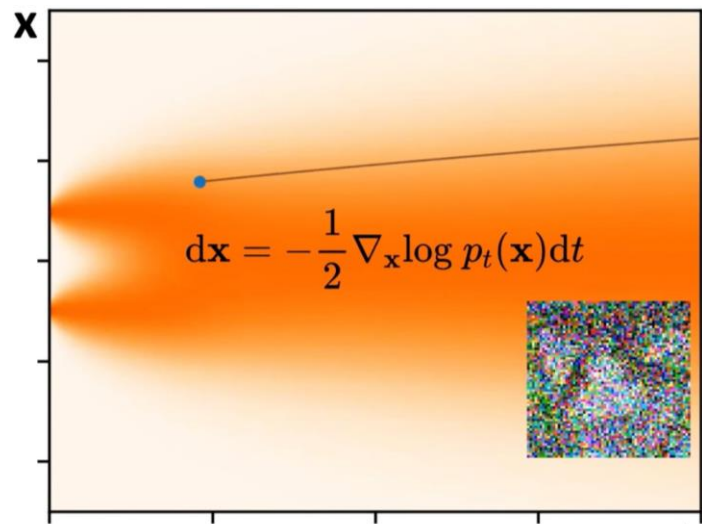
Similarly, agnostic to how the network was trained

2. Approximating ideal trajectory by finite steps

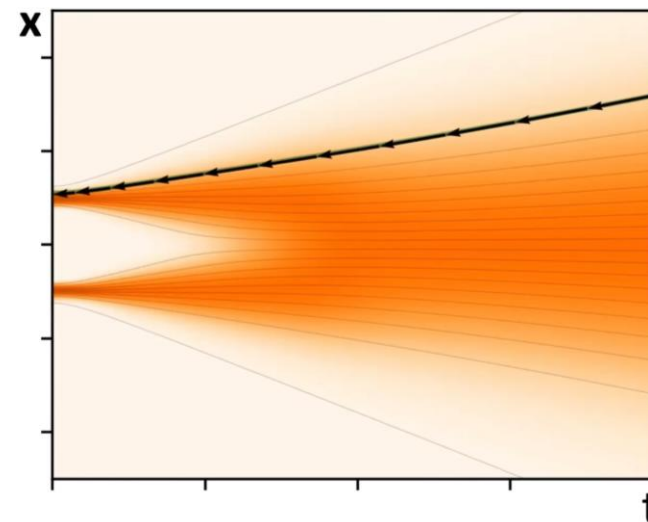
# Sampling Design Choices – ODE Solver



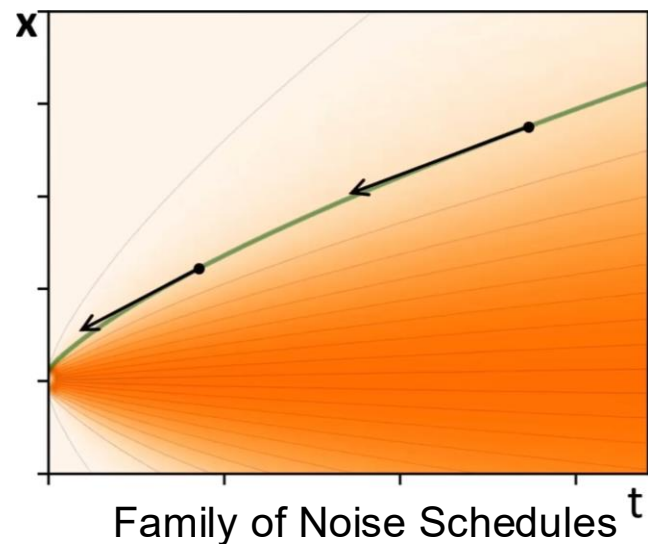
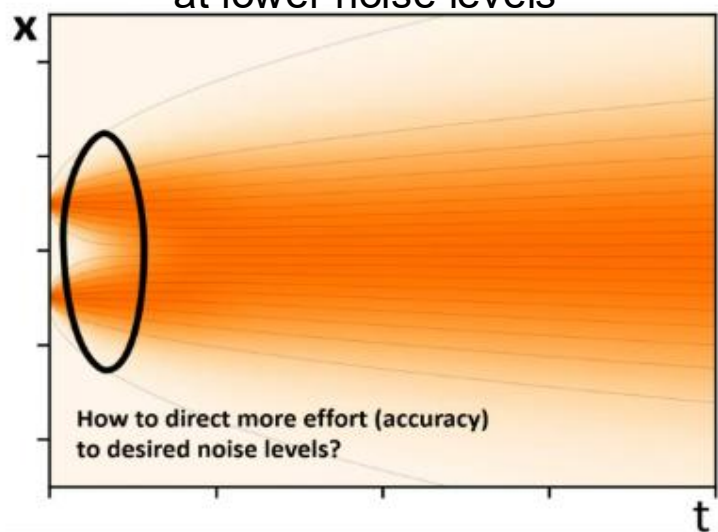
# Related Work – Diffusion Based Generative Methods



Most Details are constructed at lower noise levels

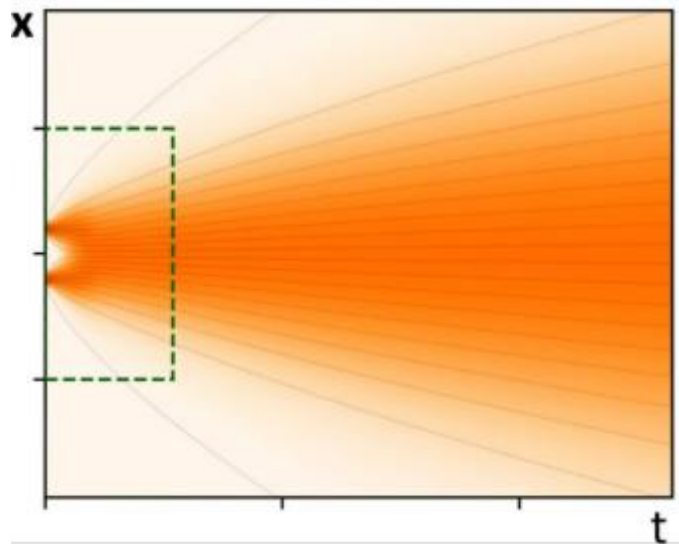


Varying Length Timesteps

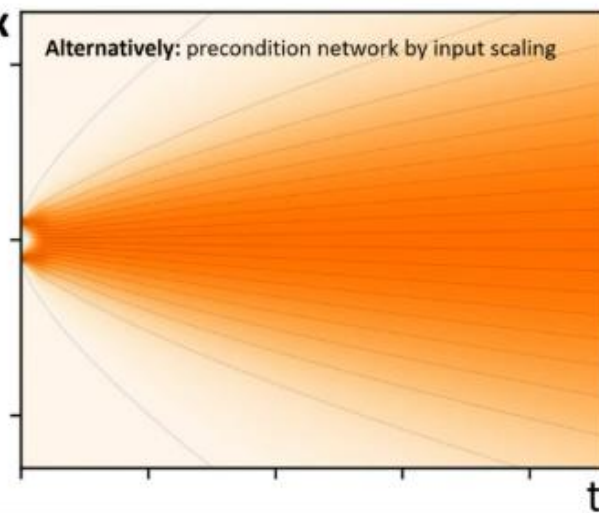
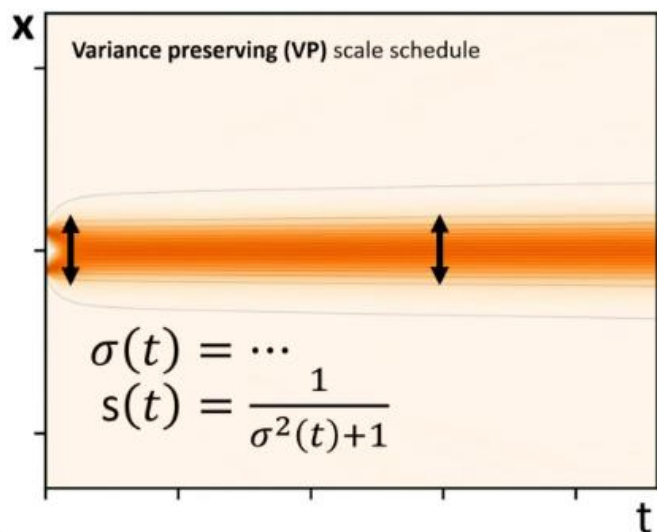


Family of Noise Schedules

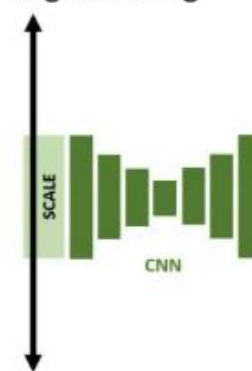
# Sampling Design Choices – Scale Schedule & Signal Scaling



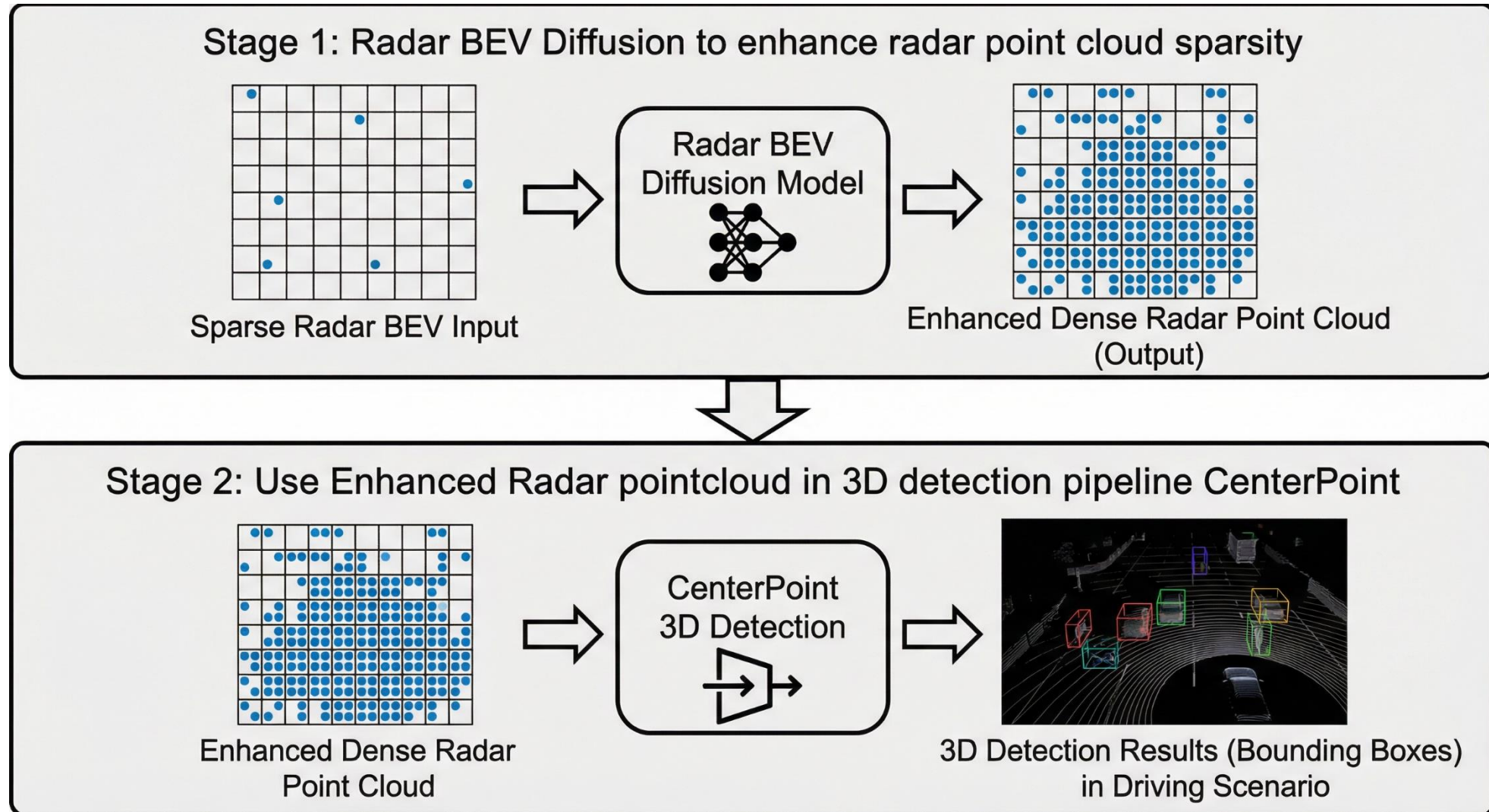
Noise magnitude grows unbounded at higher noise  $x$  levels



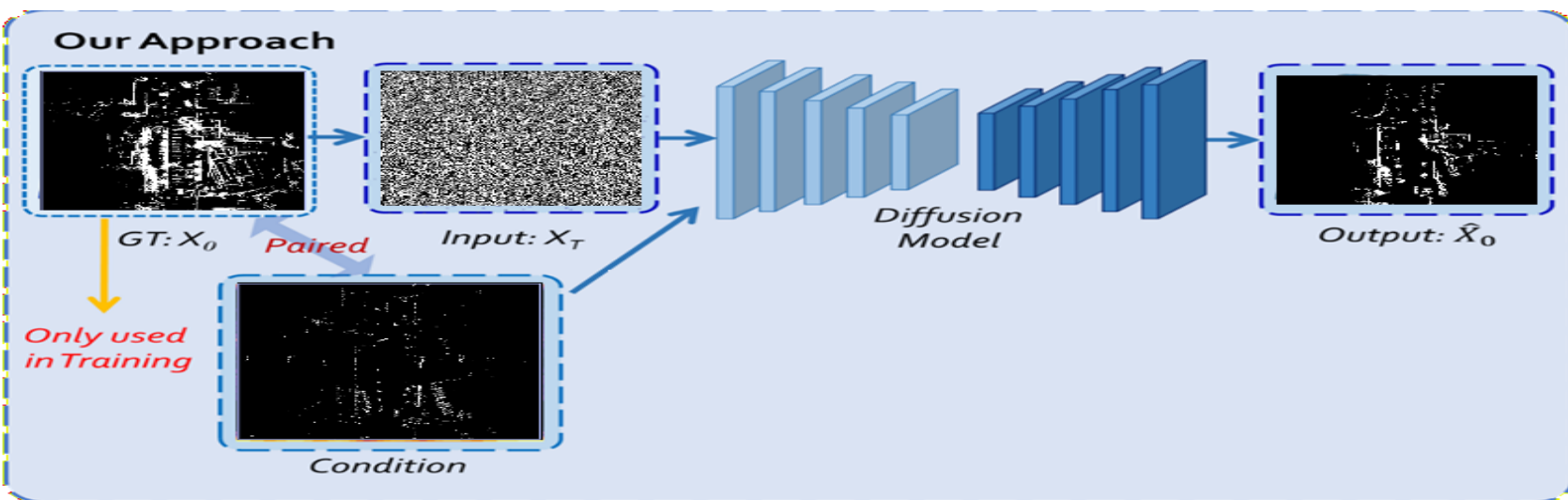
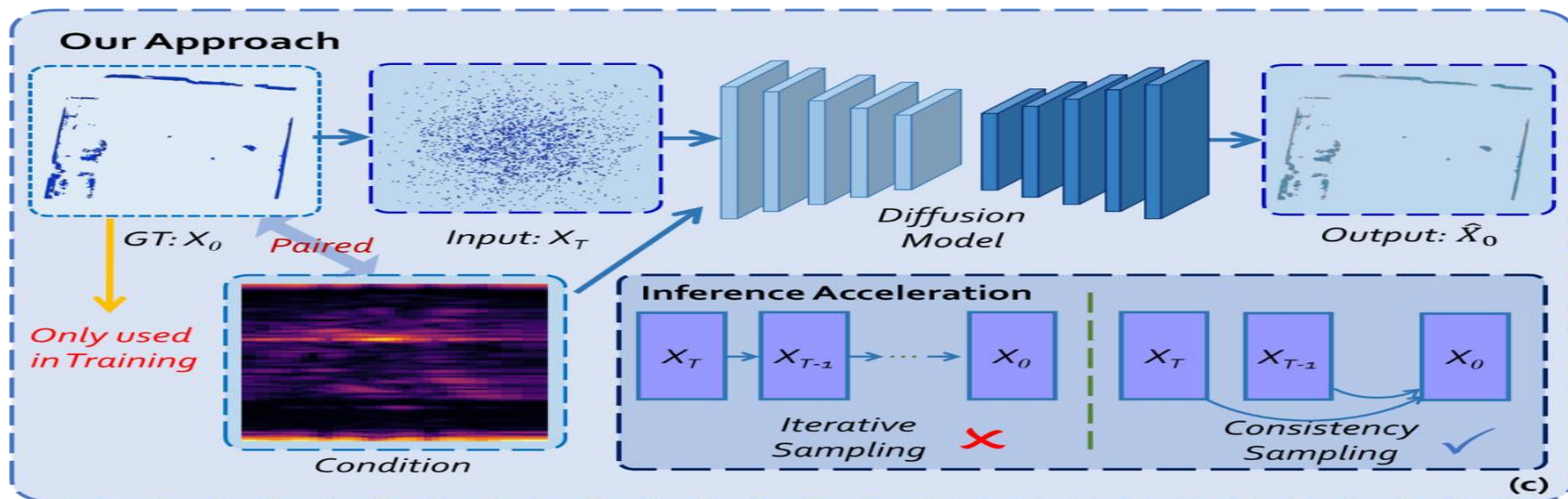
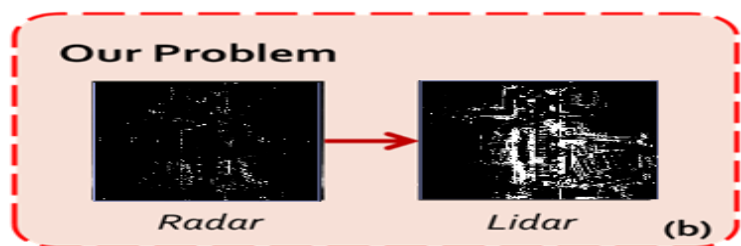
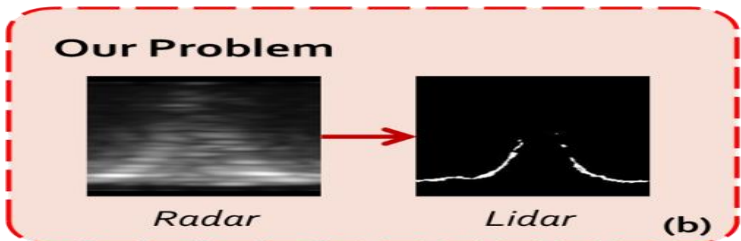
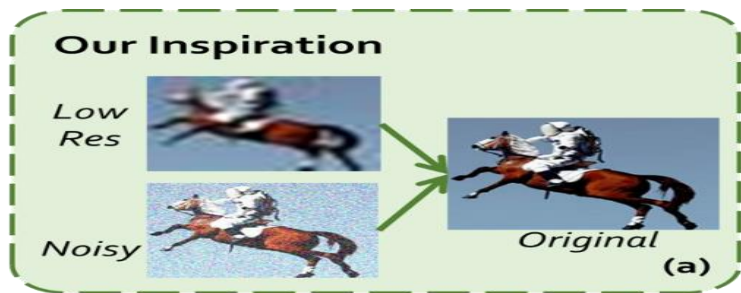
Signal scaling



# 2 Stage Approach For Radar-Only 3D Detection

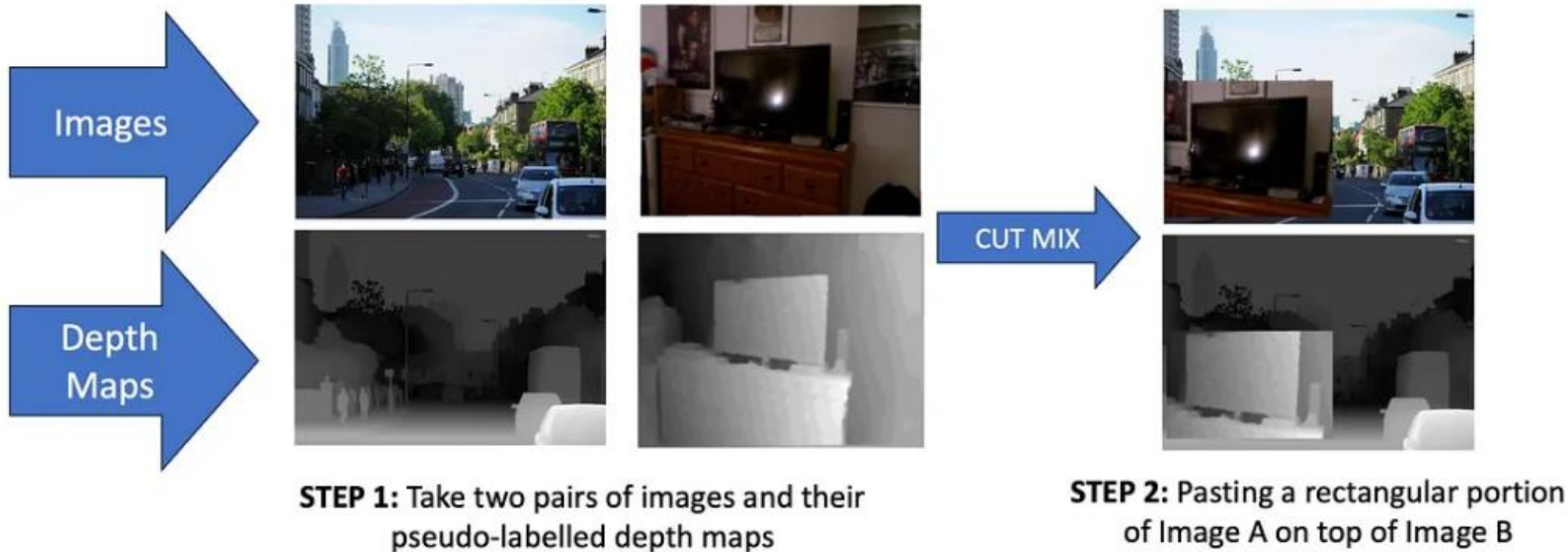


# Methodology – Starting Point (Related Work)



## CUT MIX AUGMENTATION

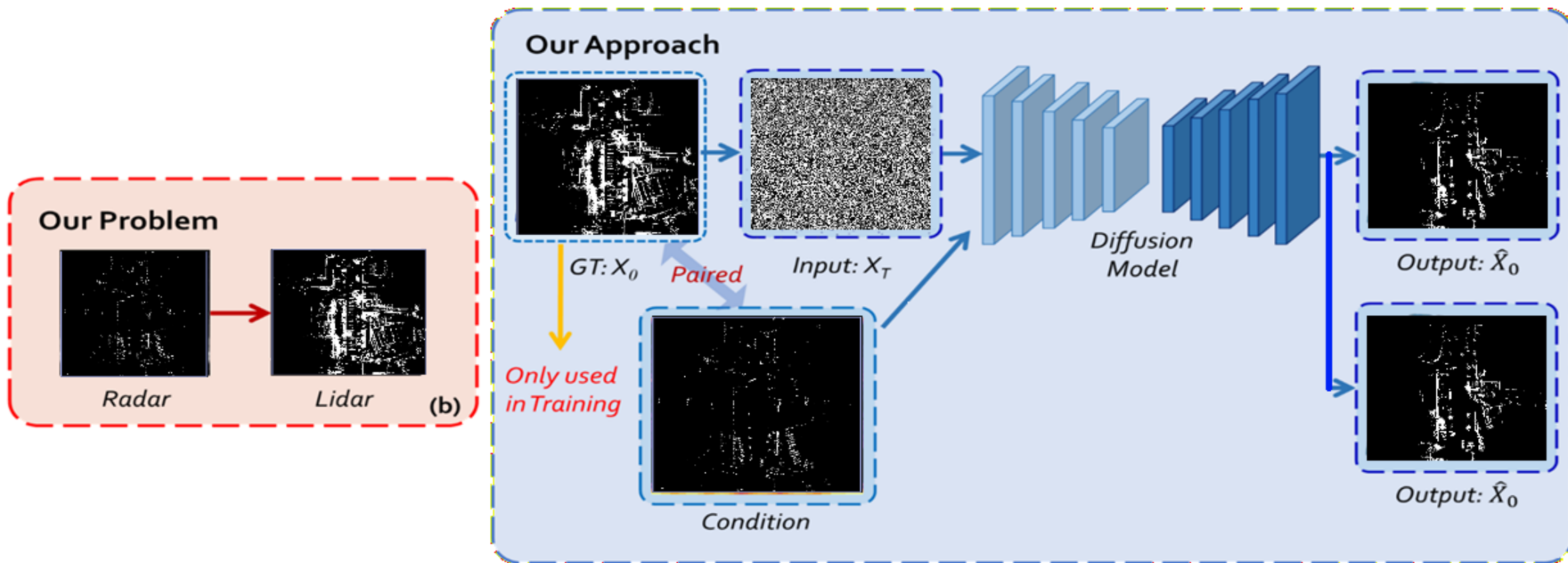
Cut Mix Augmentation increases the **SPATIAL DISTORTION** in the training images, making the Student Network learn more general and robust representations!

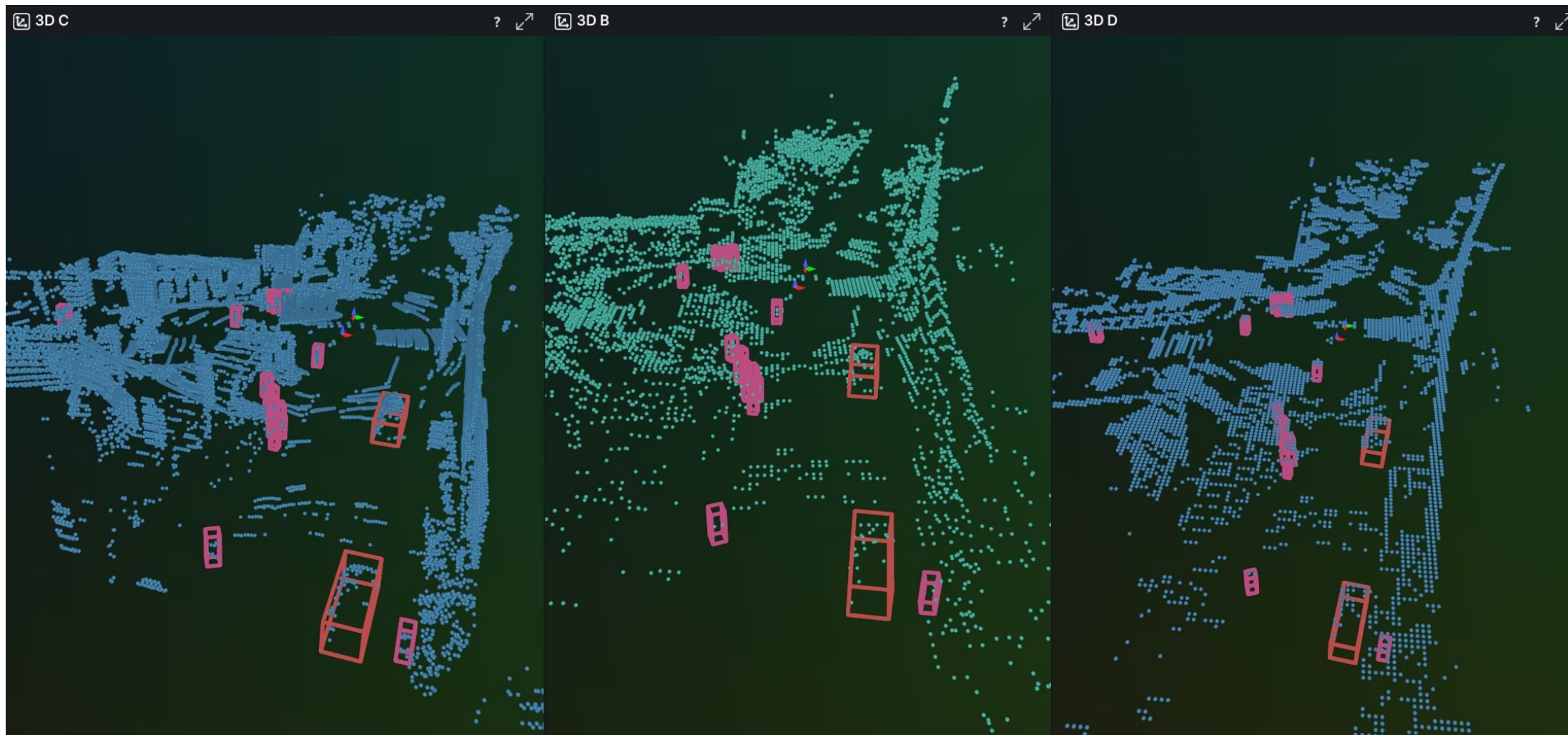


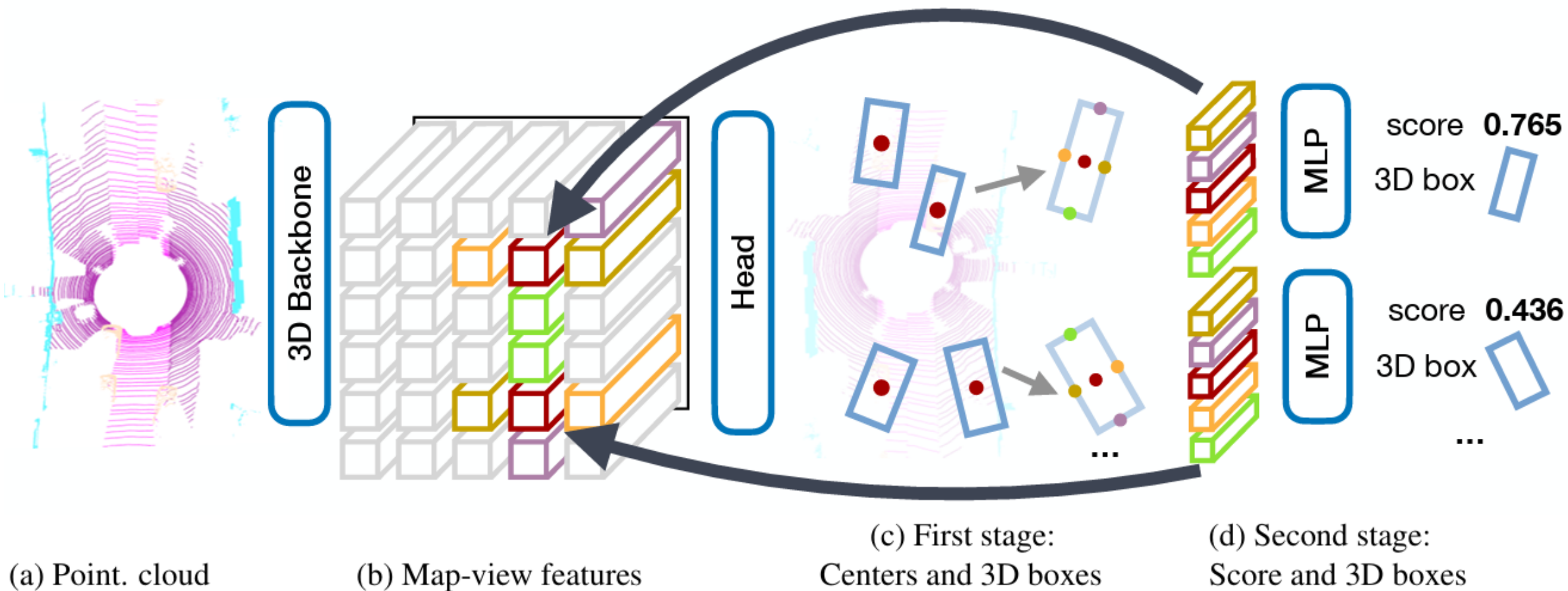
Model	Preprocessing Pipeline	mAP
CenterPoint (Benchmark)	<ul style="list-style-type: none"> <li>• Points From LiDAR MultiSweeps</li> <li>• No Motion Compensation</li> <li>• Temporal Encoding</li> </ul>	56%
RadarDistill (Target)	-	20%
CenterPoint (UpperBound)	<ul style="list-style-type: none"> <li>• Points From LiDAR MultiSweeps</li> <li>• Motion Compensation</li> <li>• Ground Point Removal</li> <li>• Points to BEV               <ol style="list-style-type: none"> <li>1. Occupancy BEV Map</li> <li>2. Constant Zero Height Map</li> </ol> </li> <li>• BEV To Points</li> </ul>	13%
RadarBEVDiff(Cut&Mix) + CenterPoint (Our 2 Stage Method)	<ul style="list-style-type: none"> <li>• Points from Denoised BEV Occupancy Map</li> <li>• BEV To Points</li> </ul>	6%

# *Question:*

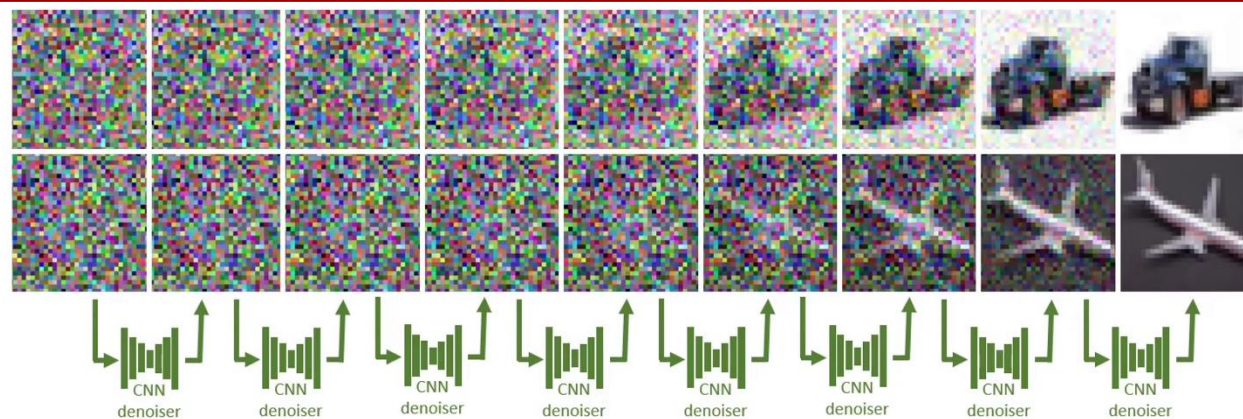
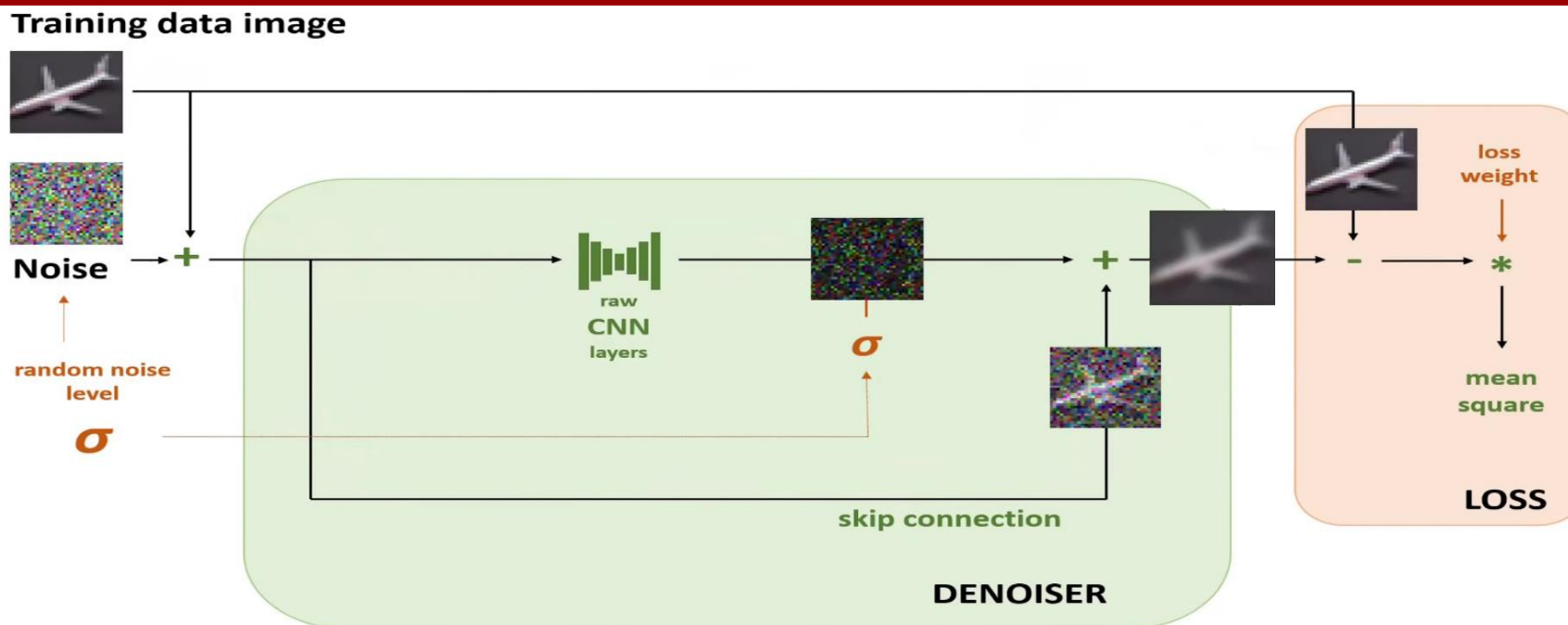
How to preserve height  
information in BEV  
representation?



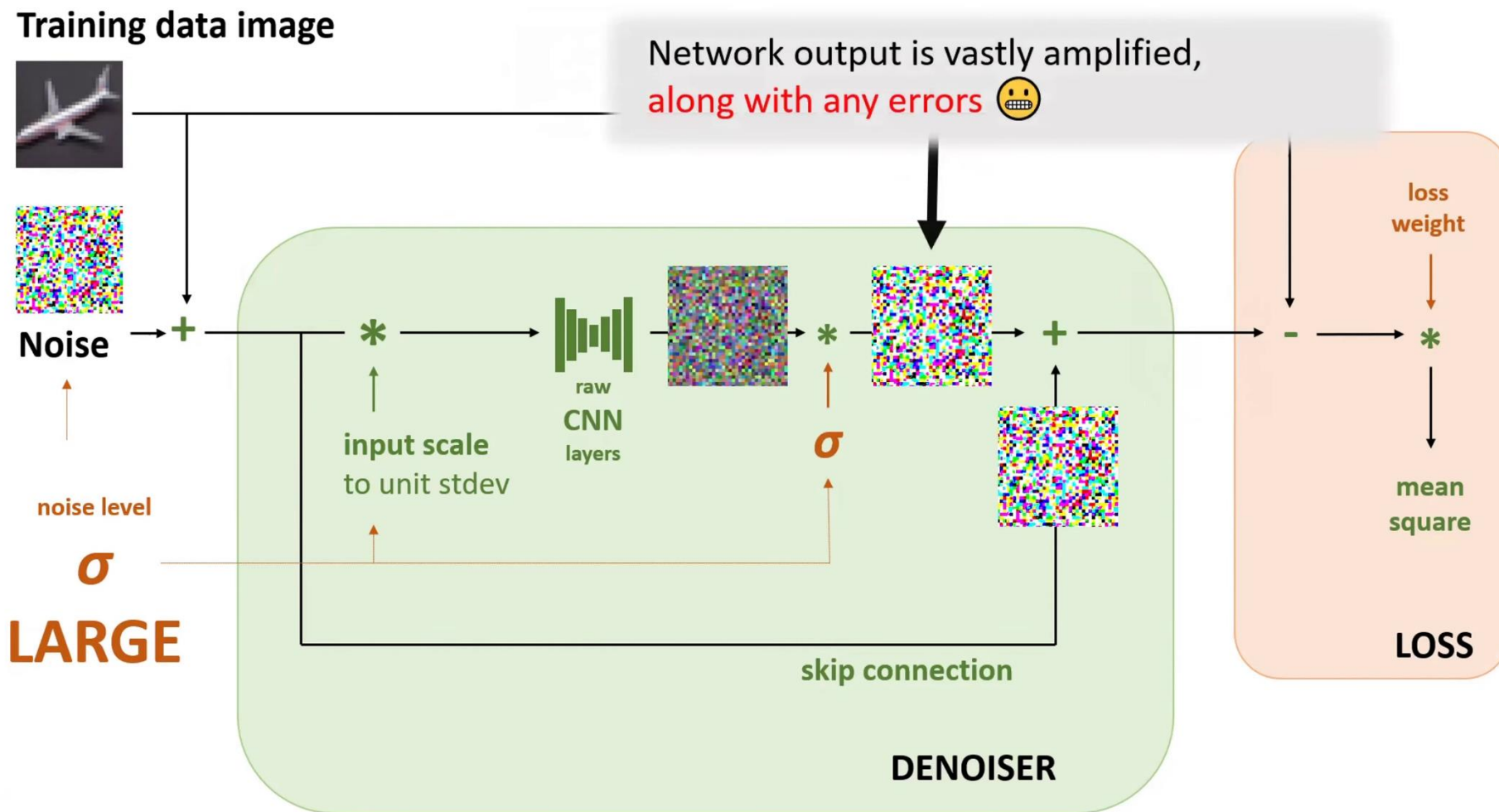




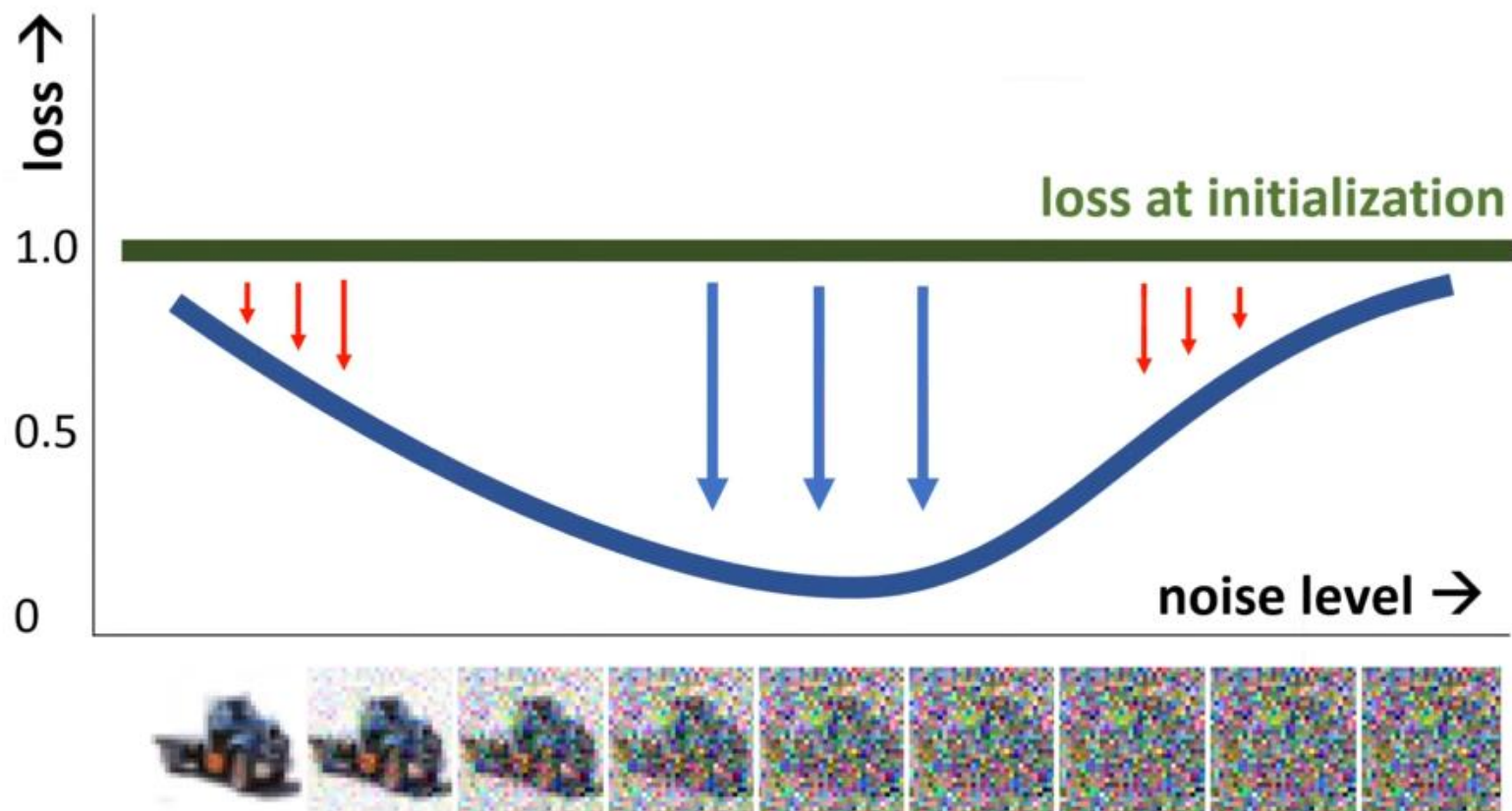
# End-To-End Approach



# CenterPoint – 3D Detection Benchmark



## Loss weighting and noise level distribution



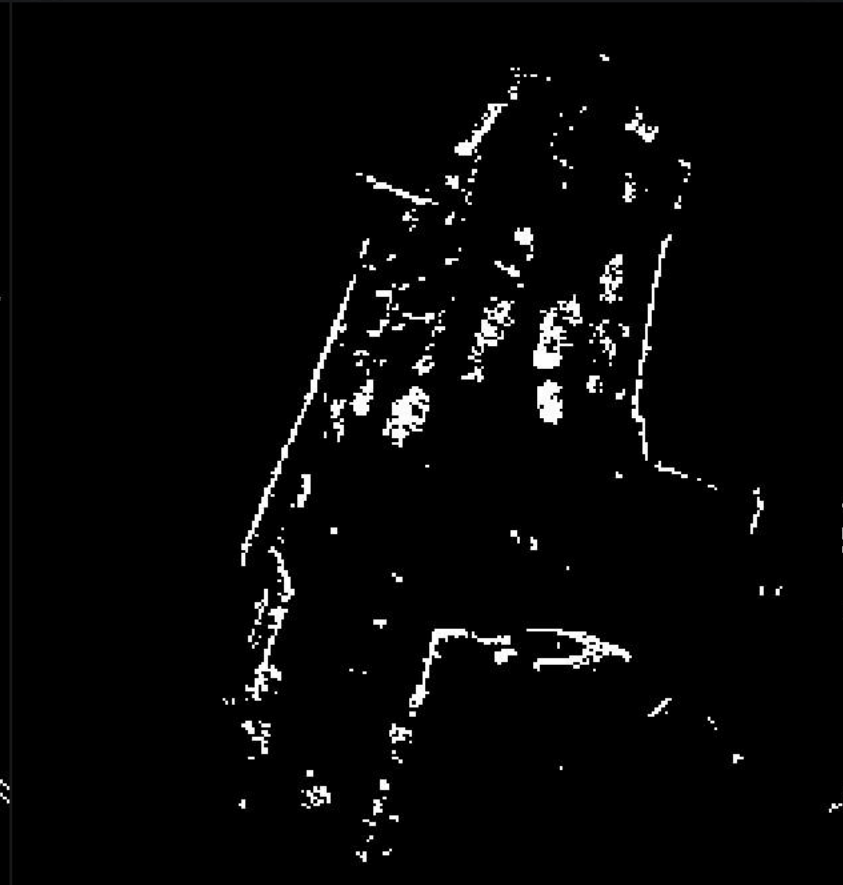
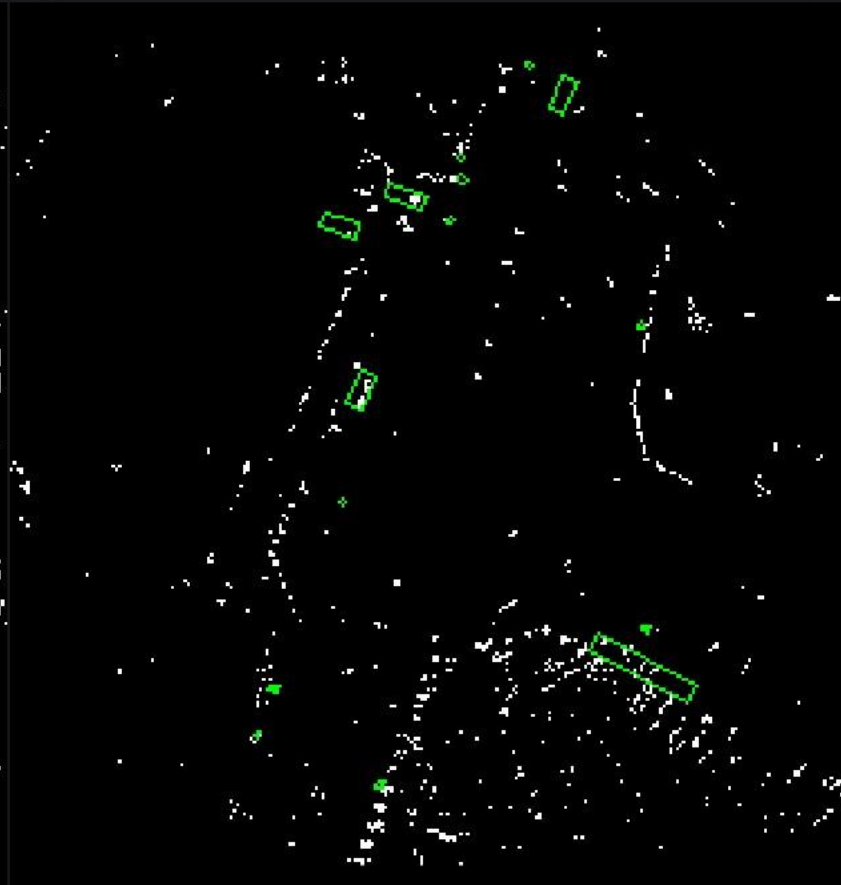
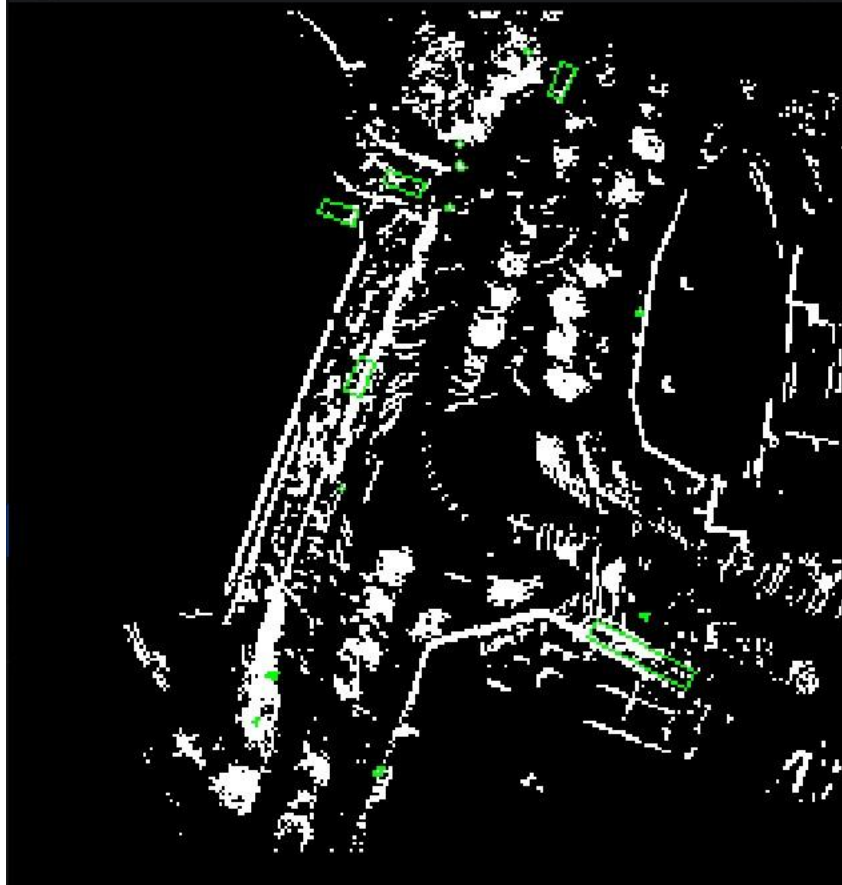
BEV/lidar\_bev\_object\_compensated



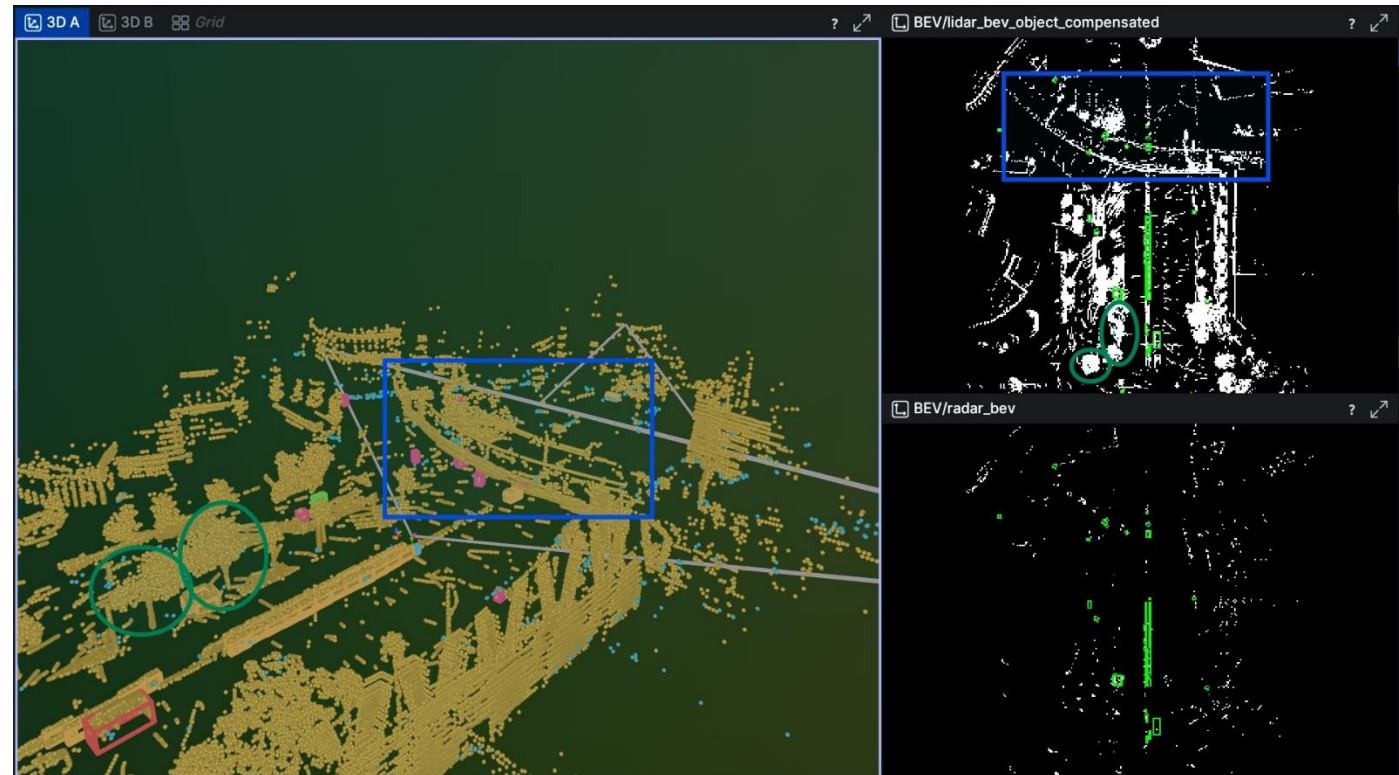
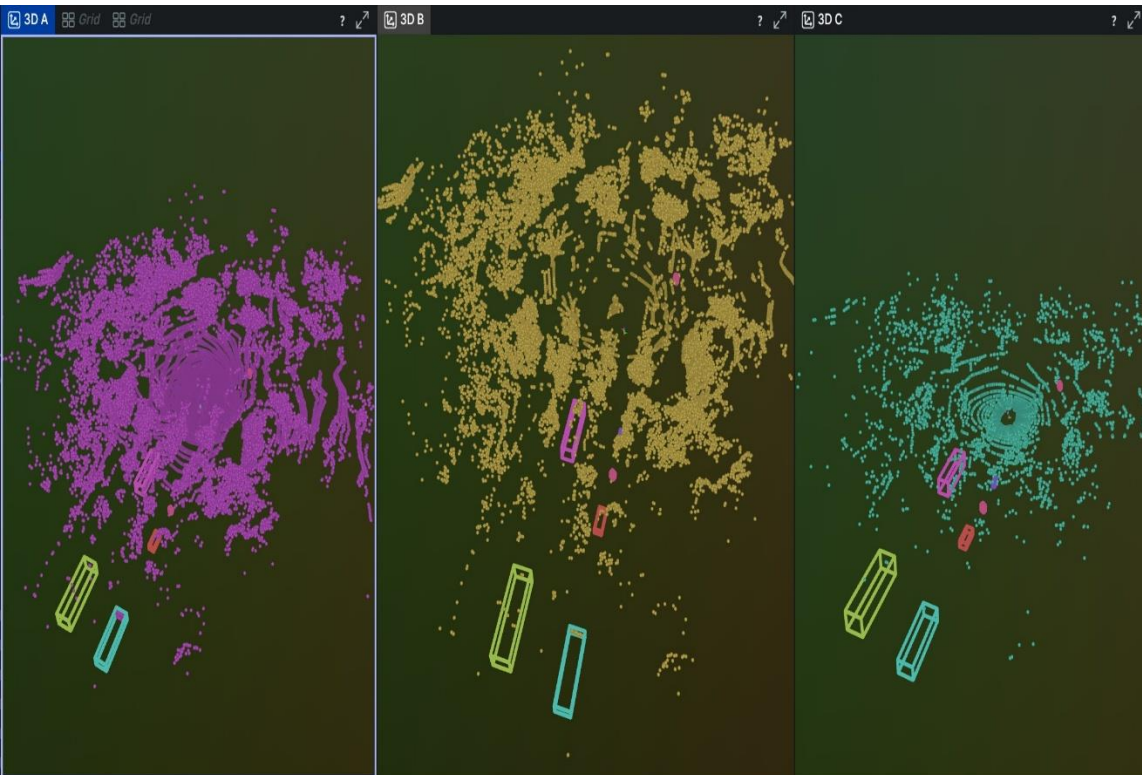
BEV/radar\_bev

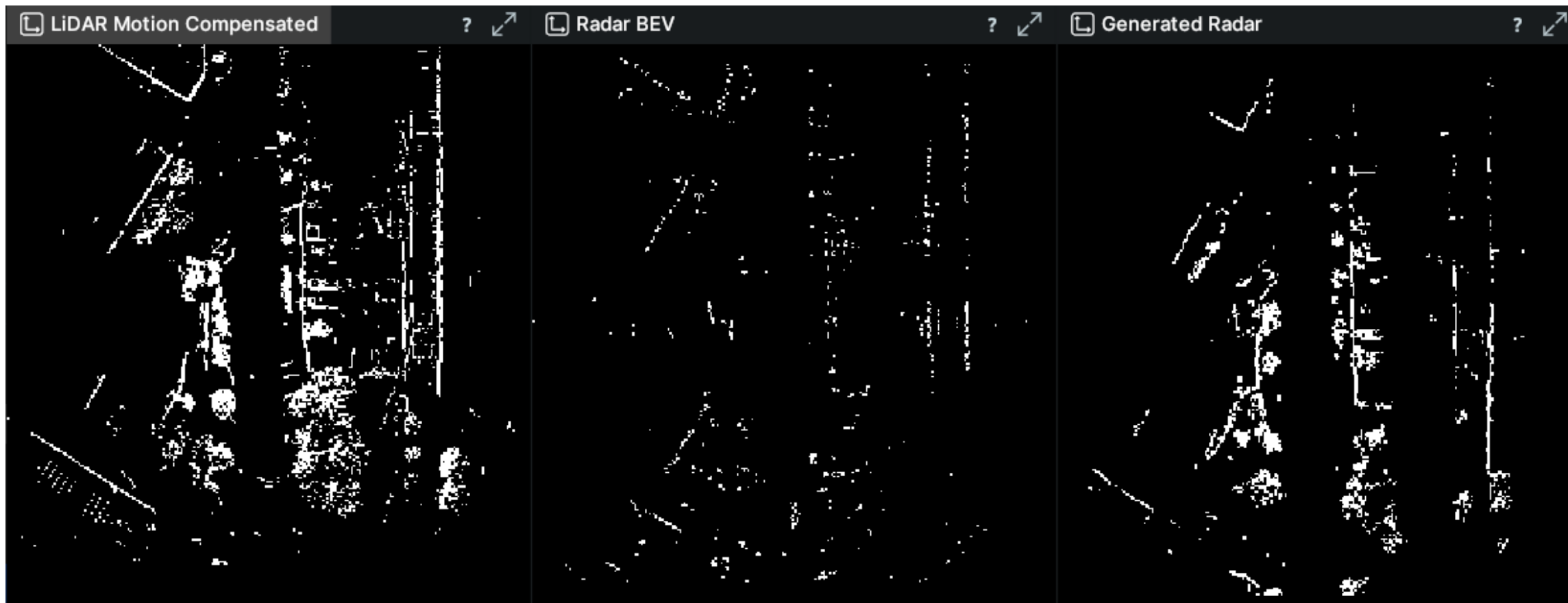


RCDIFF/x\_8



Approach	Occupancy IoU (%) Train-Val	Height MSE Train-Val	mAP (%)
<ul style="list-style-type: none"> <li>• End-To-End: Train Occupancy &amp; Height Map Denoise + CenterPoint               <ul style="list-style-type: none"> <li>– Loss: <math>MSE + L_{cls} + L_{reg}</math></li> <li>– Conditioning Signal: Radar Occupancy + Height Map</li> <li>– Augmentation: None</li> </ul> </li> </ul>	19-12	15-11	2
<ul style="list-style-type: none"> <li>• Stage 1: Train Occupancy &amp; Height Denoiser               <ul style="list-style-type: none"> <li>– Loss: MSE + BCE</li> <li>– Conditioning Signal: Radar Occupancy + <b>Height Map</b></li> <li>– Augmentation: Cut &amp; Mix</li> </ul> </li> <li>• Stage 2: Train CenterPoint on Generated Points</li> </ul>	37-18	6-8	5
<ul style="list-style-type: none"> <li>• Stage 1: Train Occupancy &amp; Height Denoiser               <ul style="list-style-type: none"> <li>– Loss: MSE</li> <li>– Conditioning Signal: Radar Occupancy Map</li> <li>– Augmentation: Cut &amp; Mix</li> </ul> </li> <li>• Stage 2: Train CenterPoint on Generated Points</li> </ul>	38-21	6-9	5
<ul style="list-style-type: none"> <li>• Stage 1: Train Occupancy &amp; Height Denoiser               <ul style="list-style-type: none"> <li>– Loss: MSE + BCE</li> <li>– Conditioning Signal: Radar Occupancy Map</li> <li>– Augmentation: Cut &amp; Mix</li> </ul> </li> <li>• Stage 2: Train CenterPoint on Generated Points</li> </ul>	34–15	7–11	3
<ul style="list-style-type: none"> <li>• Stage 1: Train Occupancy Denoiser               <ul style="list-style-type: none"> <li>– Loss: MSE + LPIPS</li> <li>– Conditioning Signal: Radar Occupancy Map</li> <li>– Augmentation: Cut &amp; Mix</li> </ul> </li> <li>• Stage 2: Train CenterPoint on Generated Points</li> </ul>	46–28	-	6
<ul style="list-style-type: none"> <li>• Stage 1: Train Occupancy Denoiser               <ul style="list-style-type: none"> <li>– Conditioning Signal: Radar Occupancy Map</li> <li>– Loss: MSE + LPIPS</li> <li>– Augmentation: None</li> </ul> </li> <li>• Stage 2: Train CenterPoint on Generated Points</li> </ul>	39–22	-	4

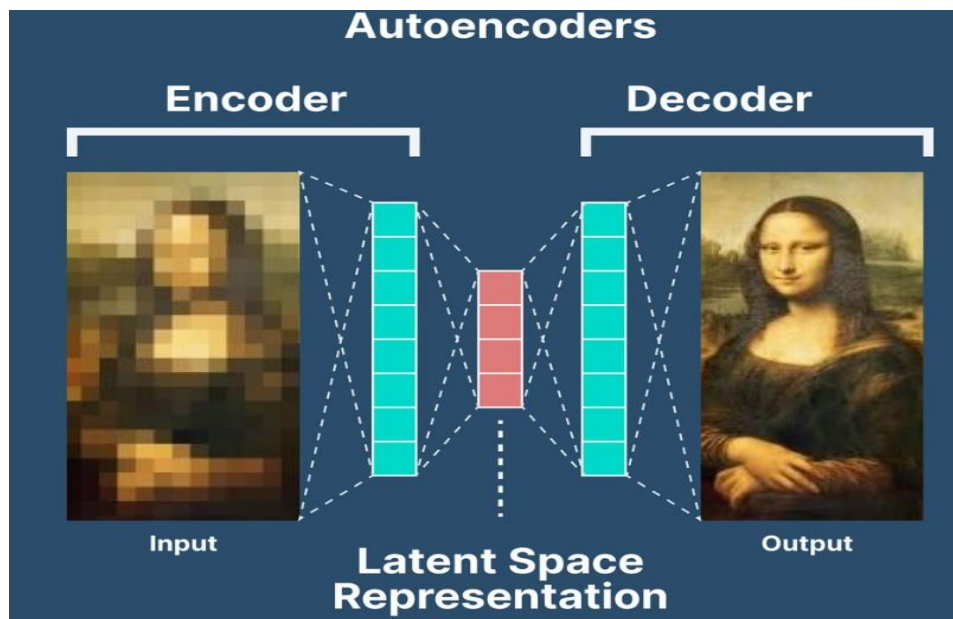
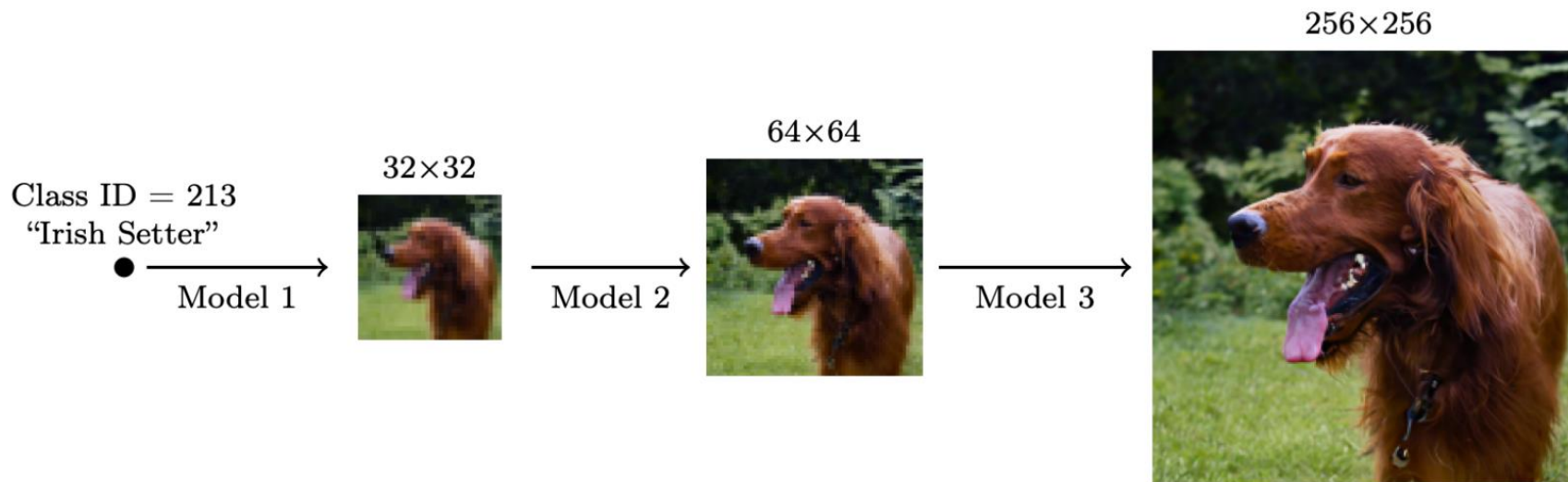




Model	Preprocessing Pipeline	Pointcloud Density(Relative to Original PC)	mAP
CenterPoint	<ul style="list-style-type: none"><li>• Points From LiDAR MultiSweeps</li><li>• Ego + Object Motion Compensation</li><li>• Ground Point Removal</li><li>• Points to BEV (Image Encoded Height)</li><li>• BEV to Points</li></ul>	70% @ 1024×1024 32% @ 512×512	42% @ 1024×1024 23% @ 512×512

Table 7.1: Comparison of preprocessing pipelines and resulting mAP.

# Future Work - Upsampling



# Future Work – Diffusion in Voxel Space

