

# Improving End-To-End Autonomous Driving with Synthetic Data from Latent Diffusion Models

## Supplementary Material

### 001 1. Organization

002 We outline the organization of the supplementary section  
003 here as follows:

- 004 1. We outline the dataset distribution across all subgroups  
005 for the datasets used in this paper in Section 2.
- 006 2. Section 3 discusses the impact of Caption Generation on  
007 the quality of synthetic data.
- 008 3. Section 4 discusses the subgroup or condition-specific  
009 performance of both semantic segmentation models and  
010 autonomous driving models fine-tuned on original and  
011 augmented datasets.
- 012 4. The performance of Autonomous Driving models (AD)  
013 over different subgroups is elaborated in Section 5.
- 014 5. Finally, we provide qualitative visualization for both seg-  
015 mentation and driving tasks in Section 6.

### 016 2. Dataset Analysis

017 For semantic segmentation tasks, the operation design do-  
018 main  $\mathcal{Z}$  and its corresponding semantic dimensions  $\mathcal{Z}_{[0,1]}$   
019 are based on weather  $\in \mathcal{Z}_0 = [\text{Rainy}, \text{Clear}, \text{Cloudy}]$  and  
020 time of day  $\in \mathcal{Z}_1 = [\text{Dawn/Dusk}, \text{Day}, \text{Night}]$ . We present  
021 the data distribution for the reader as a reference for both  
022 the BDD datasets and the Waymo Datasets.

023 For autonomous driving tasks, the operation design do-  
024 main  $\mathcal{Z}$  and its corresponding semantic dimensions  $\mathcal{Z}_{[0,1]}$   
025 are based on weather  $\in \mathcal{Z}_0 = [\text{Rainy}, \text{Clear}, \text{Cloudy}]$  and  
026 time of day  $\in \mathcal{Z}_1 = [\text{Twilight}, \text{Morning}, \text{Night}]$ . We present  
027 the data distribution for the reader as a reference for the ex-  
028 pert driving data compiled through CARLA.

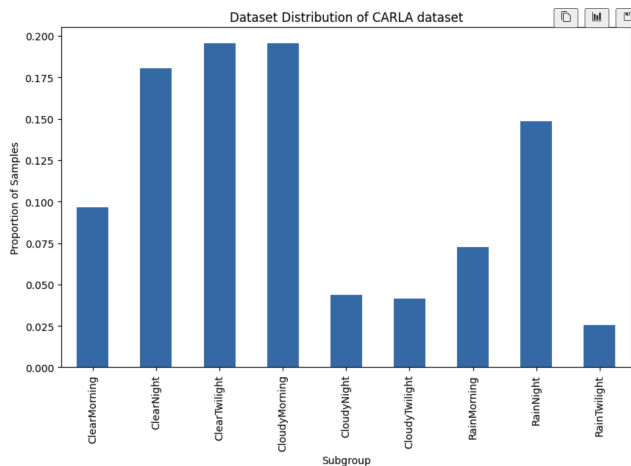


Figure 3. The distribution of autonomous driving AD for all identified subgroups in the training dataset.

Distribution	BDD100K		Waymo	
	CaG	no CaG	CaG	no CaG
Clear, Day	<b>162.16</b>	202.02	152.74	<b>146.99</b>
Clear, Dawn/Dusk	<b>66.47</b>	67.05	<b>150.92</b>	160.57
Clear, Night	<b>211.45</b>	273.28	<b>46.77</b>	80.84
Cloudy, Day	<b>134.24</b>	148.49	118.51	<b>114.93</b>
Cloudy, Dawn/Dusk	<b>144.94</b>	199.48	<b>214.55</b>	224.94
Cloudy, Night	<b>152.65</b>	246.72	<b>58.12</b>	107.57
Rainy, Day	<b>133.96</b>	154.72	121.69	<b>102.79</b>
Rainy, Dawn/Dusk	<b>199.83</b>	229.45	<b>124.35</b>	129.68
Rainy, Night	<b>291.66</b>	349.22	<b>62.21</b>	112.75

Table 1. Comparison of FD with and without caption

**generation for both datasets.** We show comprehensively that the caption generation reduces the FD score on CLIP-VIT-L16 features between the generated and the ground truth images.

### 029 3. Impact of Caption Generation(CaG)

030 To assess the impact of the proposed caption generation  
031 scheme we evaluate the quality of the synthetic images  
032 against the original ground truth images. As such we use  
033 Frechet Distance (FD) [?] scores as a suitable benchmark  
034 for the evaluation. FD score computes the distance between  
035 the feature distributions of synthetic and original images.  
036 We compute the FD scores between the data subgroup-  
037 specific distributions for both synthetic and ground truth  
038 images. Our computation of the FD is done over the  
039 image features extracted by CLIP-VIT-L16 which has a fea-  
040 ture dimension of size 768. Given that our caption genera-  
041 tion scheme using VLM improves zero shot synthetic data  
042 generation with lower FD scores, we, therefore justify the  
043 use of LLaVA to caption images for text descriptions that  
044 are used in the downstream fine-tuning of ControlNet with  
045 frozen Stable Diffusion weights for semantic segmentation  
046 and AD tasks.

### 047 4. Effect of Fine-tuning on condition-specific 048 performance

049 This set of experiments compares the effect of fine-tuning  
050 over synthetic data generated for various under-represented  
051 data subgroups. We refer you to Table 2 and Table 3 for  
052 these results.

053 For semantic segmentation specifically, we conduct  
054 an ablation study over two components of our proposed  
055 pipeline. First, we conduct an ablation over the effect

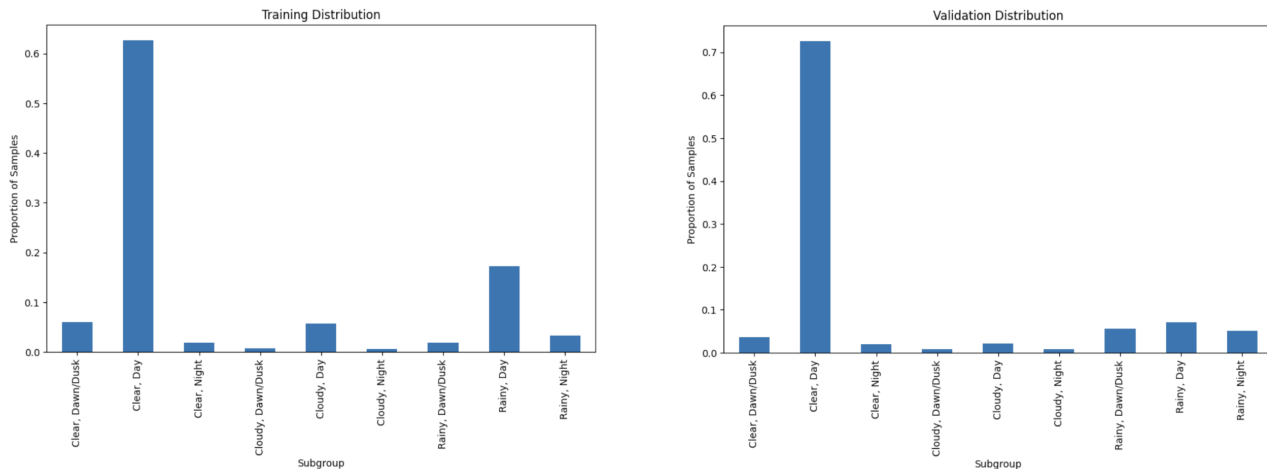


Figure 1. The distribution of image-segmentation mask pairs over all identified subgroups in the **Waymo** training and validation dataset.

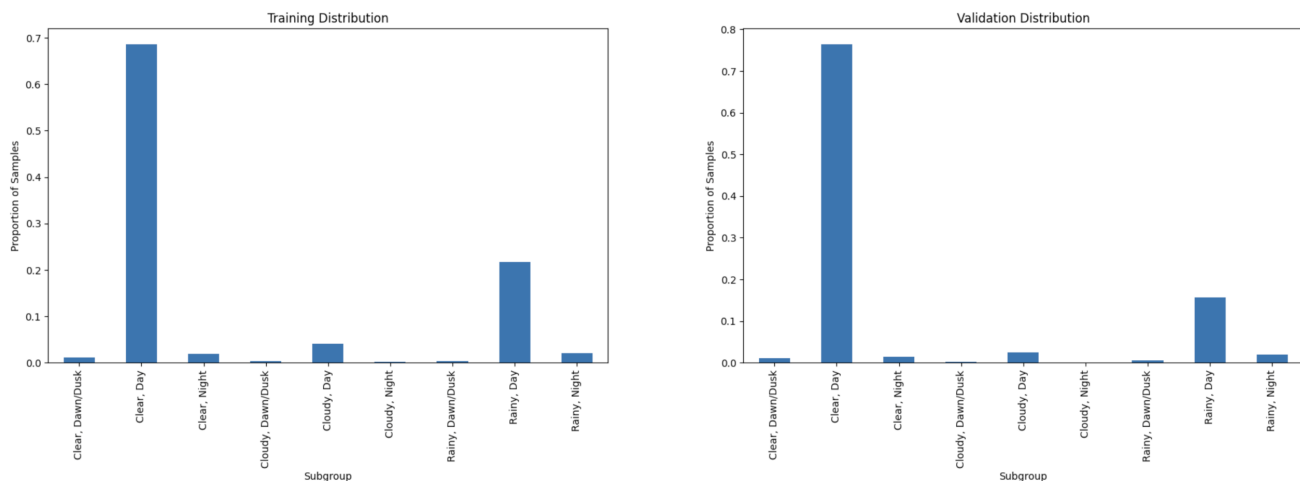


Figure 2. The distribution of image-segmentation mask pairs for all identified subgroups in the **BDD** training and validation dataset.

056 of fine-tuning ControlNet with data sampled over all sub-  
 057 groups equivalently. Thus, image-caption-segmentation  
 058 mask tuples that are from rarer subgroups are sampled more  
 059 selectively during fine-tuning. Synthetic data generated  
 060 from this variant is referred to as Synthetic-RST (Rare Sub-  
 061 Group Training). Second, we modify the method by which  
 062 we sample source images for which we want synthetic data  
 063 variants. Here synthetic images are sampled such that all  
 064 semantic categories are equivalently present. This results  
 065 in synthetic data with equivalent semantic class distribu-  
 066 tions that would enable selective training over rare seman-  
 067 tic classes. This was shown to improve the performance of  
 068 semantic segmentation models in prior work [? ]. We re-  
 069 port the results of the ablation study for the best-performing  
 070 model i.e. Mask2Former over all the synthetic datasets.  
 071 We see that across different subgroups, the best-performing  
 072 models are obtained by fine-tuning over datasets augmented  
 073 with synthetic data.

We report the per subgroup performance of various AD  
 models for our tests on Autonomous Driving. In the case  
 of Autonomous Driving, synthetic data is generated for all  
 camera views across an entire route. Hence, the ablations  
 proposed for semantic segmentation don't extend to AD in  
 our setup. The averaged driving scores are reported for all  
 9 data subgroups for all models fine-tuned on both origi-  
 nal dataset and augmented datasets. We see noticeable im-  
 provements in the driving score of AD models AIM-2D and  
 AIM-BEV when trained with synthetic data augmentations  
 using SynDiff-AD on all data subgroups. In contrast, syn-  
 thetic data augmentations degraded NEAT's performance  
 due to reasons mentioned in Section ??.

## 5. Performance of Autonomous Driving Mod- els

We present a detailed breakdown of the Driving Scores(DS)  
 as referenced in the main paper. Here we present the Route

091 Completion (RC) scores and the Infraction Scores (IS) of  
 092 the learned AD policies for each data subgroup and model.  
 093 In the following tables, we report the above metrics for each  
 094 AD model trained on the original and synthetic data. We  
 095 highlight the best performing models for each metric across  
 096 each subgroup. We refer you to Tables 4, 5 and 6

097 **6. Qualitative Visualizations**

098 We attempt to provide qualitative visualizations of the ob-  
 099 tained synthetic images for different tasks and datasets.  
 100 Here we sample an image and semantic mask pair and  
 101 showcase its variants across different data subgroups.

102 **6.1. Waymo**

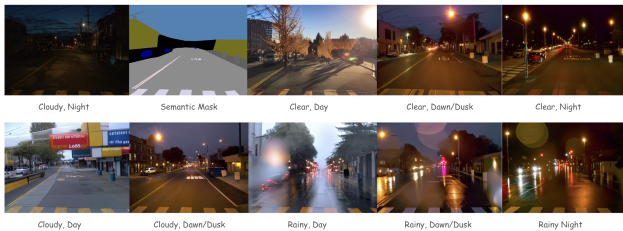


Figure 4. Sample visualization of synthetic images for a source image and mask taken in cloudy weather and night time



Figure 5. Sample visualization of synthetic images for a source image and mask taken in clear weather and day time

103 **6.2. BDD100K**

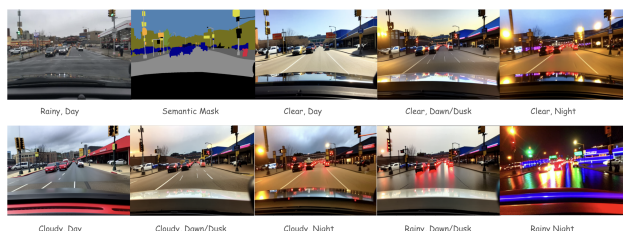


Figure 6. Sample visualization of synthetic images for a source image and mask taken in rainy weather and day time

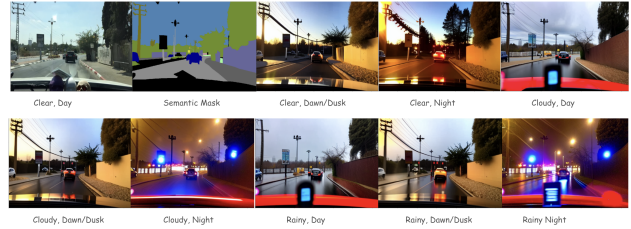


Figure 7. Sample visualization of synthetic images for a source image and mask taken in clear weather and day time

**6.3. Autonomous Driving Carla**

104

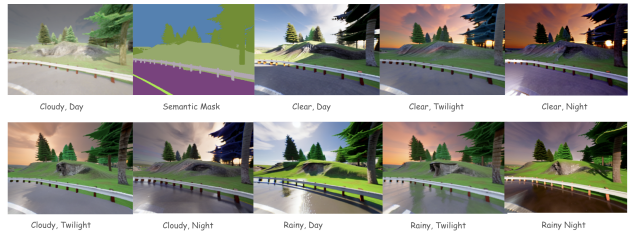


Figure 8. Sample visualization of synthetic images for a source image and mask taken in cloudy weather and day time

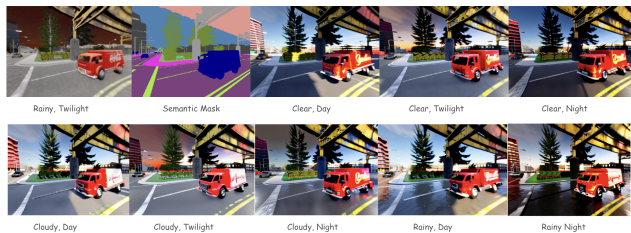


Figure 9. Sample visualization of synthetic images for a source image and mask taken in rainy weather and twilight time

Dataset	Sub-Group	Original	Synthetic	Synthetic RST	Synthetic CEQ	Synthetic RST-CEQ
Waymo	Clear, Dawn/Dusk	42.2	44.5	43.1	43.8	<b>46.6</b>
	Clear, Day	51.8	<b>52.9</b>	51.0	51.6	52.1
	Clear, Night	33.6	36.4	<b>39.3</b>	34.2	33.9
	Cloudy, Dawn/Dusk	48.2	48.8	47.7	<b>49.5</b>	49.2
	Cloudy, Day	<b>56.3</b>	55.8	<b>56.3</b>	52.5	55.4
	Cloudy, Night	38.3	38.1	35.9	37.5	<b>39.9</b>
	Rain, Dawn/Dusk	39.8	<b>41.4</b>	41.4	40.8	41.0
	Rain, Day	50.7	50.7	<b>52.7</b>	51.5	50.8
	Rain, Night	35.1	34.8	35.1	<b>36.5</b>	36.1
BDD100K	Clear, Dawn/Dusk	42.3	51.1	<b>52.3</b>	50.6	47.8
	Clear, Day	56.7	57.5	57.2	55.9	<b>57.8</b>
	Clear, Night	42.0	<b>51.3</b>	42.7	49.0	45.8
	Cloudy, Dawn/Dusk	40.8	42.4	42.6	35.8	<b>44.3</b>
	Cloudy, Day	52.2	<b>60.4</b>	55.8	57.0	59.6
	Cloudy, Night	35.4	49.3	47.8	49.0	<b>52.3</b>
	Rain, Dawn/Dusk	49.9	<b>57.6</b>	52.0	50.1	49.3
	Rain, Day	54.8	56.6	56.0	<b>57.7</b>	55.2
	Rain, Night	30.5	35.8	35.7	35.2	<b>36.4</b>

Table 2. **Improved performance over different data subgroups with synthetic data augmentation.** We conduct an ablation study that constructs synthetic datasets using three approaches - RST, CEQ, and RST - CEQ. RST datasets comprise images from a fine-tuned ControlNet that equally samples rare subgroups during training. CEQ datasets are sampled so that the synthetic dataset’s semantic class distribution is uniform. RST-CEQ incorporates both these strategies.

Model	Aug	Clear			Cloudy			Rain		
		Tw	Day	Night	Tw	Day	Night	Tw	Day	Night
NEAT	No	35.86	5.09	21.87	17.71	36.66	17.46	3.49	23.10	21.72
NEAT	Yes	12.14	16.73	8.86	15.95	14.18	20.24	6.86	14.48	13.32
AIM-2D	No	19.69	23.20	6.11	<b>37.25</b>	18.77	43.72	14.68	23.93	3.42
AIM-2D	Yes	39.04	40.30	<b>29.02</b>	19.11	33.44	<b>46.32</b>	16.68	50.94	<b>34.23</b>
AIM-BEV	No	39.78	31.42	2.73	29.88	<b>44.68</b>	43.22	18.07	43.73	3.97
AIM-BEV	Yes	<b>58.37</b>	<b>47.94</b>	25.42	14.93	<b>44.22</b>	35.03	<b>27.52</b>	<b>53.67</b>	15.64

Table 3. **AD models trained on augmented datasets exhibit improved driving performance** We show that AD models fine-tuned on augmented datasets (indicated by Aug) have improved performance, especially over rare subgroups where the models trained on the original dataset underperform.

Metric	Aug	Clear			Cloudy			Rain		
		Tw	Day	Night	Tw	Day	Night	Tw	Day	Night
RC		53.08	43.29	28.44	35.48	53.75	53.41	9.22	53.24	33.81
IS	No	<b>0.828</b>	0.386	0.747	0.629	<b>0.829</b>	0.499	<b>0.595</b>	<b>0.645</b>	<b>0.667</b>
DS		35.86	5.09	21.87	17.71	36.66	17.46	3.49	23.10	21.72
RC		33.56	33.47	52.21	32.97	33.22	50.73	31.10	33.63	30.03
IS	Yes	0.597	0.644	0.441	<b>0.654</b>	0.631	0.667	0.453	0.633	0.649
DS		12.14	16.73	8.86	15.95	14.18	20.24	6.86	14.48	13.32

Table 4. **Performance of NEAT across different data sub-groups**

Metric	Aug	Clear			Cloudy			Rain		
		Twi	Day	Night	Twi	Day	Night	Twi	Day	Night
<b>RC</b>		76.04	55.22	84.91	54.05	76.05	55.02	48.65	84.61	45.97
<b>IS</b>	No	0.224	0.483	0.073	0.729	0.204	0.727	0.259	0.244	0.352
<b>DS</b>		19.69	23.20	6.11	<b>37.25</b>	18.77	43.72	14.68	23.93	3.42
<b>RC</b>		84.81	<b>55.30</b>	84.05	54.98	83.59	55.02	82.58	77.02	<b>69.38</b>
<b>IS</b>	Yes	0.392	0.631	0.312	0.339	0.373	<b>0.791</b>	0.266	0.551	0.472
<b>DS</b>		39.04	40.30	<b>29.02</b>	19.11	33.44	<b>46.32</b>	16.68	50.94	<b>34.23</b>

Table 5. Performance of AIM-2D across different data-subgroups

Metric	Aug	Clear			Cloudy			Rain		
		Twi	Day	Night	Twi	Day	Night	Twi	Day	Night
<b>RC</b>		83.48	55.18	76.31	<b>55.18</b>	<b>100.0</b>	<b>55.18</b>	64.51	83.14	69.21
<b>IS</b>	No	0.436	0.589	0.038	0.573	0.447	0.706	0.269	0.462	0.255
<b>DS</b>		39.78	31.42	2.73	29.88	<b>44.68</b>	43.22	18.07	43.73	3.97
<b>RC</b>		<b>100.0</b>	<b>55.28</b>	<b>100.0</b>	23.66	90.86	36.99	<b>85.92</b>	<b>100.0</b>	<b>69.38</b>
<b>IS</b>	Yes	0.584	<b>0.706</b>	0.254	0.438	0.452	0.677	0.374	0.536	0.285
<b>DS</b>		<b>58.37</b>	<b>47.94</b>	25.42	14.93	<b>44.22</b>	35.03	<b>27.52</b>	<b>53.67</b>	15.64

Table 6. Performance of AIM-BEV across different data-subgroups