

CLIP-BEVFormer: Enhancing Multi-View Image-Based BEV Detector with Ground Truth Flow

Supplementary Material

001	1. Experiment Setup	
002	In our experiments, we adopt ResNet50 and ResNet101	
003	[2] as backbones for the BEVFormer-tiny and BEVFormer-	
004	base models, respectively. These backbones are initialized	
005	from the FCOS3D [6] checkpoint, following the configura-	
006	tion in BEVFormer [3]. We leverage the output multi-scale	
007	features from the Feature Pyramid Network (FPN) [4], with	
008	sizes of 1/16, 1/32, and 1/64, and the dimension of 256.	
009	During the training phase for CLIP-BEVFormer, we lever-	
010	age the pretrained language model in CLIP-RN101 [5] as	
011	our off-the-shelf language model.	
012	The BEV size for the tiny and base variants is set to $50 \times$	
013	50 and 200×200 , respectively, while the perception ranges	
014	span from -51.2m to 51.2m along the X and Y axes. The	
015	resolution of the BEV grid is set to 0.512m . We incorporate	
016	learnable positional embeddings for BEV queries to enrich	
017	the spatial representation.	
018	The BEV encoder comprises 6 encoder layers, consis-	
019	tently refining BEV queries in each layer. During the spatial	
020	cross-attention module, implemented using the deformable	
021	attention mechanism, each local query corresponds to four	
022	target points with different heights in 3D space. The prede-	
023	efined height anchors are uniformly sampled from -5 meters	
024	to 3 meters.	
025	For each reference point on 2D view features, we utilize	
026	four sampling points around this reference point for each	
027	head. During training, we use a 2-frame history BEV for the	
028	tiny variant and a 3-frame history BEV for the base variant.	
029	We train our models for 24 epochs with a learning rate of	
030	2×10^{-4} [3].	
031	2. 3D Object Detection Results with Various	
032	Baselines	
033	We have conducted experiments with various detection	
034	baselines, BEVformer [3], BEVformerV2 [7], and BEVerse	
035	[8]. We evaluate our model on both validation and test sets	
036	of nuScenes. The results presented in Tab. 1 show that our	
037	method consistently improves the perception capabilities of	
038	various baselines by significant margins on both sets, indi-	
039	cating its flexibility and model-agnostic nature.	
040	3. 3D Object Detection Metrics	
041	We adhere to standard evaluation metrics for 3D detection	
042	on the nuScenes dataset [1], encompassing metrics such as	
043	mean Average Precision (mAP), Average Translation Er-	
044	ror (ATE), Average Scale Error (ASE), Average Orientation	
	Error (AOE), Average Velocity Error (AVE), Average At-	045
	tribute Error (AAE), and nuScenes detection score (NDS).	046
	Mean Average Precision (mAP). For mAP, we utilize the	047
	Average Precision metric, modifying the definition of a	048
	match by considering the 2D center distance on the ground	049
	plane instead of intersection over union-based affinities.	050
	Specifically, we match predictions with ground truth ob-	051
	jects based on the smallest center distance within a certain	052
	threshold. Average precision (AP) is calculated by integrat-	053
	ing the recall vs precision curve for recalls and precisions	054
	> 0.1 . We then average over match thresholds of 0.5, 1, 2,	055
	4 meters and compute the mean across classes.	056
	True Positives (TP). TP metrics are designed to measure	057
	translation, scale, orientation, velocity, and attribute errors.	058
	These are calculated using a threshold of 2m center distance	059
	during matching and are positive scalars. Metrics are de-	060
	efined per class, and we then take the mean over classes to	061
	calculate mATE, mASE, mAOE, mAVE, and mAAE.	062
	<ul style="list-style-type: none">Average Translation Error (ATE). Euclidean center dis-	063
	tance in 2D in meters.	064
	<ul style="list-style-type: none">Average Scale Error (ASE). Calculated as $1 - \text{IOU}$ after	065
	aligning centers and orientation.	066
	<ul style="list-style-type: none">Average Orientation Error (AOE). Smallest yaw angle	067
	difference between prediction and ground truth in radi-	068
	ans. Orientation error is evaluated at 360 degrees for most	069
	classes, except barriers, where it is evaluated at 180 de-	070
	grees. Orientation errors for cones are ignored.	071
	<ul style="list-style-type: none">Average Velocity Error (AVE). Absolute velocity error	072
	in m/s. Velocity error for barriers and cones is ignored.	073
	<ul style="list-style-type: none">Average Attribute Error (AAE). Calculated as $1 - \text{acc}$,	074
	where acc is the attribute classification accuracy. At-	075
	tribute error for barriers and cones is ignored.	076
	nuScenes Detection Score (NDS). We consolidate the	077
	above metrics by computing a weighted sum: mAP, mATE,	078
	mASE, mAOE, mAVE, and mAAE. As a first step, we	079
	convert TP errors to TP scores using $\text{TP_score} = \max(1 -$	080
	$\text{TP_error}, 0.0)$. We then assign a weight of 5 to mAP and 1	081
	to each of the 5 TP scores, calculating the normalized sum.	082
	4. Training and Inference Efficiency	083
	Our model is trained with 4 A100 80GB GPUs. Our method	084
	does not introduce any additional parameters and computa-	085
	tions during the inference stage, which means that it allows	086
	for enhanced performance without sacrificing real-time pro-	087
	cessing capabilities. We provide detailed information on	088
	memory, training time, number of parameters and FPS in	089
	Tab. 2.	090

091 **5. Visualization**

092 In Fig. 1 and Fig. 2, we present a comprehensive visualiza-
093 tion of the qualitative detection performance achieved by
094 CLIP-BEVFormer. The images provide insights into both
095 camera and Bird’s Eye View (BEV) perspectives, offering
096 a nuanced understanding of the model’s predictions. No-
097 tably, these visualizations highlight the enhanced alignment
098 between CLIP-BEVFormer’s predictions and ground truth
099 detections in both camera and BEV views, underscoring the
100 model’s proficiency in accurately capturing the 3D environ-
101 ment.

102 **References**

- 103 [1] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora,
104 Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Gi-
105 ancarmo Baldan, and Oscar Beijbom. nuscnescenes: A multi-
106 modal dataset for autonomous driving. In *Proceedings of the*
107 *IEEE/CVF conference on computer vision and pattern recog-*
108 *niton*, pages 11621–11631, 2020. 1
- 109 [2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.
110 Deep residual learning for image recognition. In *Proceedings*
111 *of the IEEE conference on computer vision and pattern recog-*
112 *niton*, pages 770–778, 2016. 1
- 113 [3] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao
114 Sima, Tong Lu, Yu Qiao, and Jifeng Dai. Bevformer: Learn-
115 ing bird’s-eye-view representation from multi-camera images
116 via spatiotemporal transformers. In *European conference on*
117 *computer vision*, pages 1–18. Springer, 2022. 1
- 118 [4] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He,
119 Bharath Hariharan, and Serge Belongie. Feature pyramid net-
120 works for object detection. In *Proceedings of the IEEE confer-*
121 *ence on computer vision and pattern recognition*, pages 2117–
122 2125, 2017. 1
- 123 [5] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
124 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
125 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learn-
126 ing transferable visual models from natural language supervi-
127 sion. In *International conference on machine learning*, pages
128 8748–8763. PMLR, 2021. 1
- 129 [6] Tai Wang, Xinge Zhu, Jiangmiao Pang, and Dahua Lin.
130 Fcos3d: Fully convolutional one-stage monocular 3d object
131 detection. In *Proceedings of the IEEE/CVF International*
132 *Conference on Computer Vision*, pages 913–922, 2021. 1
- 133 [7] Chenyu Yang, Yuntao Chen, Hao Tian, Chenxin Tao, Xizhou
134 Zhu, Zhaoxiang Zhang, Gao Huang, Hongyang Li, Yu Qiao,
135 Lewei Lu, et al. Bevformer v2: Adapting modern image back-
136 bones to bird’s-eye-view recognition via perspective supervi-
137 sion. In *Proceedings of the IEEE/CVF Conference on Com-*
138 *puter Vision and Pattern Recognition*, pages 17830–17839,
139 2023. 1
- 140 [8] Yunpeng Zhang, Zheng Zhu, Wenzhao Zheng, Junjie Huang,
141 Guan Huang, Jie Zhou, and Jiwen Lu. Beverse: Unified per-
142 ception and prediction in birds-eye-view for vision-centric au-
143 tonomous driving. *arXiv preprint arXiv:2205.09743*, 2022. 1

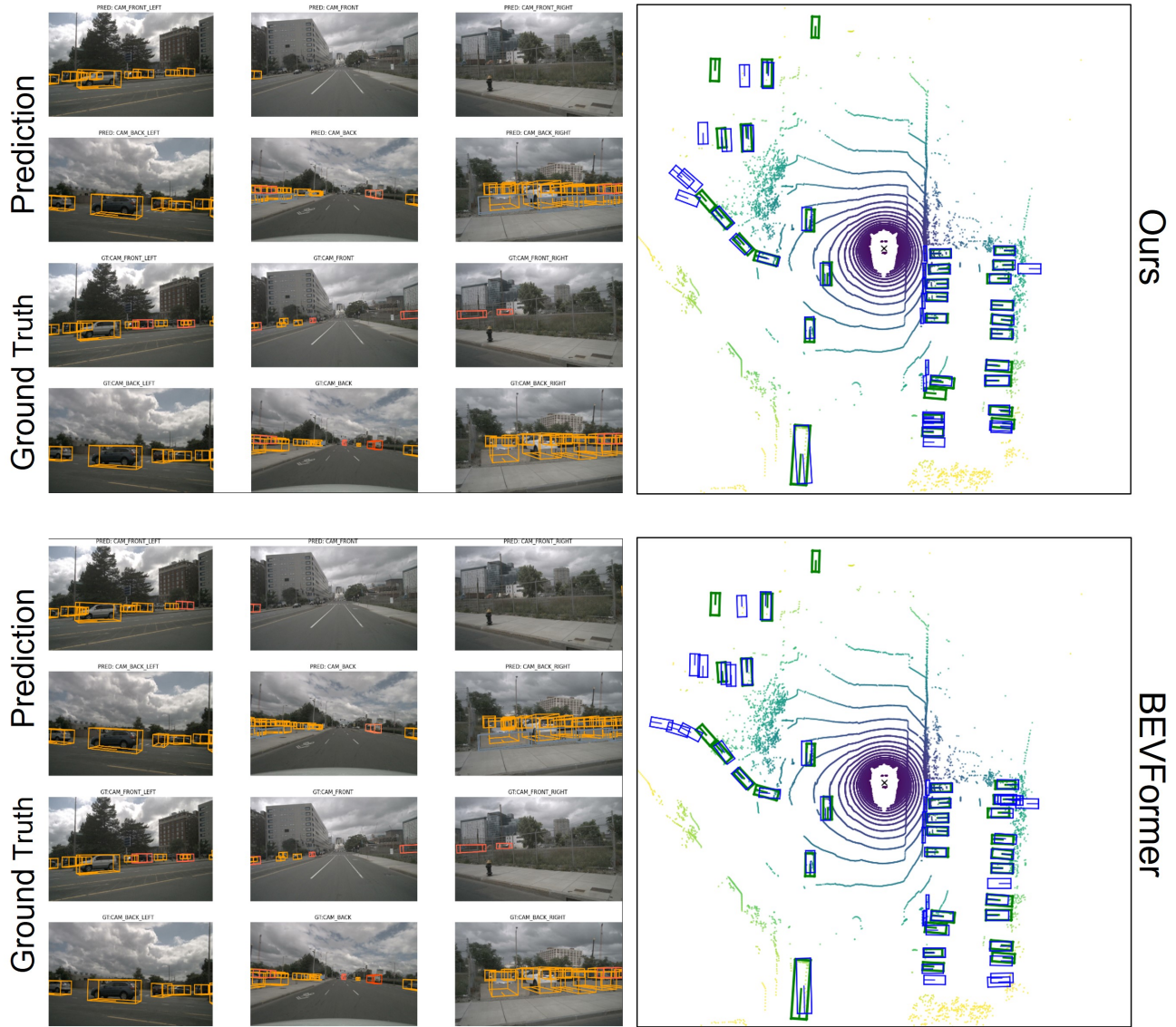


Figure 1. **Visualization results on nuScenes validation set.** We demonstrate qualitative detection performance on both camera and BEV images. As can be seen in BEV images, our CLIP-BEVFormer method demonstrates improved alignment with ground truth detections.

Model	Backbone	Validation Set								Test Set					
		NDS \uparrow	mAP \uparrow	mATE \downarrow	mASE \downarrow	mAOE \downarrow	mAVE \downarrow	mAAE \downarrow	NDS \uparrow	mAP \uparrow	mATE \downarrow	mASE \downarrow	mAOE \downarrow	mAVE \downarrow	mAAE \downarrow
BEVformer-tiny	R50	35.5	25.1	0.898	0.293	0.651	0.657	0.216	37.2	27.3	0.856	0.283	0.609	0.753	0.146
+Ours	R50	38.8	27.3	0.856	0.282	0.583	0.538	0.228	41.1	29.3	0.811	0.271	0.554	0.579	0.136
BEVformer-base	R101	51.7	41.6	0.673	0.274	0.372	0.394	0.198	53.5	44.5	0.631	0.257	0.405	0.435	0.143
+Ours	R101	55.1	44.1	0.641	0.253	0.319	0.307	0.172	54.7	44.7	0.591	0.257	0.417	0.371	0.128
BEVformerV2	R50	42.6	35.1	0.753	0.286	0.466	0.807	0.186	42.5	35.4	0.707	0.278	0.506	0.895	0.134
+Ours	R50	44.1	37.0	0.729	0.281	0.438	0.791	0.204	43.6	37.9	0.676	0.272	0.475	0.975	0.141
BEVerse	Swin	46.6	32.1	0.681	0.278	0.466	0.328	0.190	50.1	36.2	0.610	0.257	0.451	0.355	0.131
+Ours	Swin	48.3	34.2	0.665	0.270	0.456	0.318	0.170	52.2	37.4	0.556	0.247	0.413	0.301	0.129

Table 1. 3D object detection results on nuScenes validation and test sets.

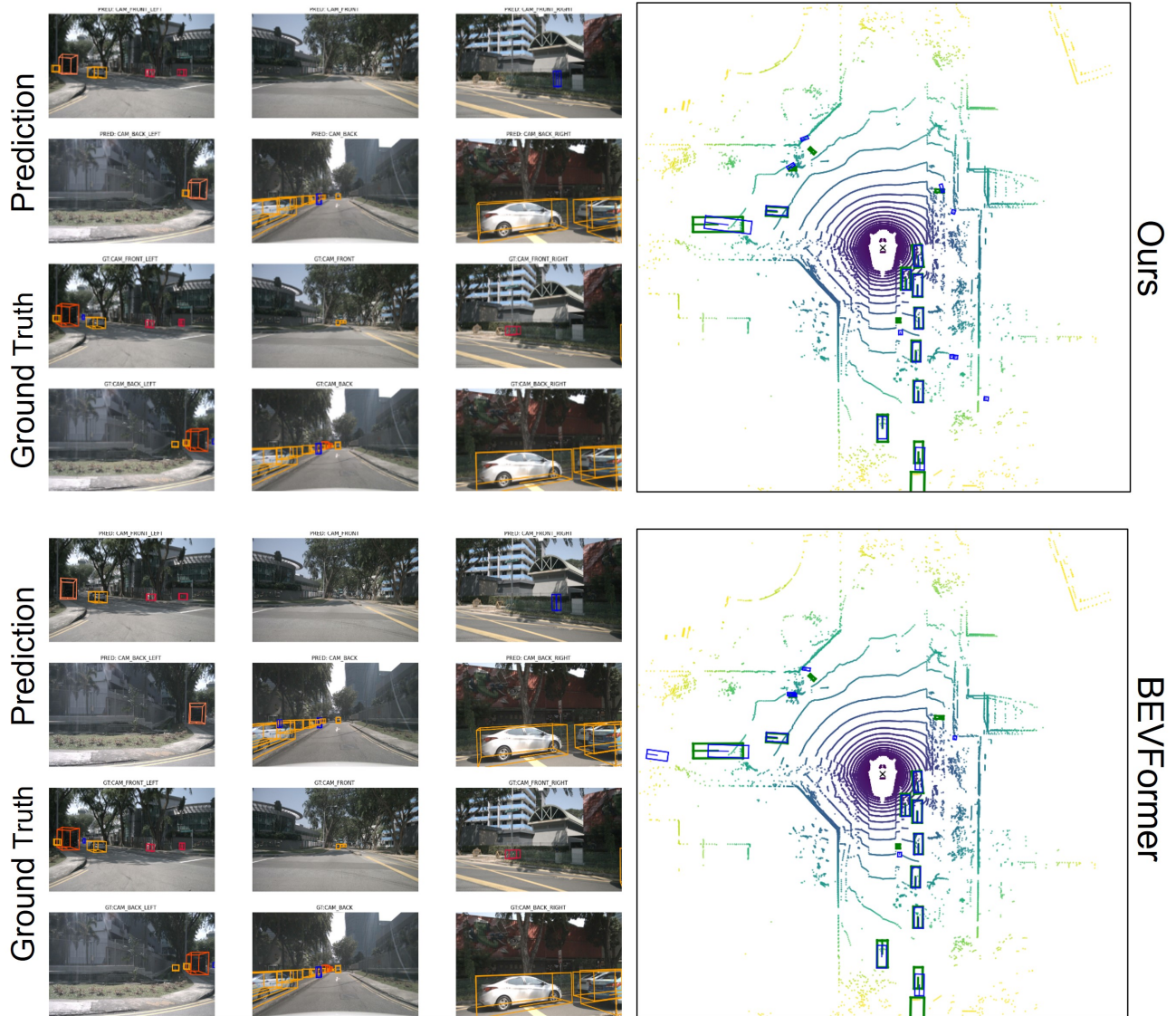


Figure 2. **Visualization results on nuScenes validation set.** Our CLIP-BEVFormer demonstrates improved alignment with ground truth detections on both camera and BEV images.

Model	Train Mem (GB)	Train Hrs	# Params (M)	FPS
BEVformer-tiny	~7	~46	33	5.1
+Ours	~7	~46	33	5.1
BEVformer-base	~25	~90	69	2.1
+Ours	~25	~90	69	2.1
BEVformerV2	~46	~38	56	2.3
+Ours	~46	~38	56	2.3
BEVerse	~48	~72	102.5	4.4
+Ours	~48	~72	102.5	4.4

Table 2. Efficiency details. FPS is tested on 1 V100 GPU.