## Appendices

## A BACKGROUND OF TRANSFORMER ARCHITECTURE AND SELF-ATTENTION

This section outlines the Transformer architecture and its self-attention mechanism, foundational to Vaswani et al.] (2017). Designed initially for NLP tasks, the Transformer employs an encoder-decoder structure to process sequences. The encoder converts an input sequence into continuous representations, which the decoder uses to synthesize an output sequence. Each consists of repeated layers that include self-attention and position-wise feedforward networks. The core of the Transformer's innovative approach lies in its self-attention mechanism, which updates the representation of each sequence element based on the context provided by the entire sequence:

$$Attention(Q, K, V) = softmax(\frac{QK^{T}}{\sqrt{d_{k}}})V$$
(4)

where Q (Query), K (Key), and V (Value) matrices are generated through linear transformations of the input's hidden states, with  $d_k$  representing the dimensionality of these states for appropriate scaling. This self-attention operation allows each element to dynamically adjust its representation by integrating information weighted by attention scores from the entire sequence, thereby enhancing the model's ability to capture contextual relationships within the data.

775 776 777

756

757 758 759

760 761

762

763

764

765

766

767 768

769

## B KITTI PEDESTRIAN ANNOTATION WITH MMCAT

KITTI Pedestrian Category. To further demonstrate the efficacy and generalizability of our MMCAT
 multimodal framework, we conduct experiments on the Pedestrian class of KITTI. It's important
 to note that, for a predicted 3D box to be considered acceptable in the Pedestrian category, it must
 achieve an Intersection over Union (IoU) greater than 0.5 with the ground truth, contrasting with the
 higher threshold of 0.7 sets for the Vehicle category.

Data Preparation. Of the 3,712 training scenes in KITTI, only 951 scenes contain pedestrian labels.
Considering the small number of training samples, we randomly choose 25% (515 samples) of 2,257 samples in those scenes. Compared to previously fully supervised PointRCNN, which leverages all 951 exhaustively annotated scenes with 2,257 pedestrian samples, we use far fewer and weak 2D supervisions.

788 **Implementation Details.** As for sparser Pedestrian point clouds, we adapted from previous studies, 789 normalizing point counts per sample within a batch to the median (around 450) across the batch, 790 addressing point density variability. We adjust the MMCAT encoder for the Pedestrian category to 791 avoid overfitting on these sparser samples. The image encoder also has two blocks ( $L_2 = 2$ ), the 2D 792 box encoder two ( $L_3 = 2$ ), and each multimodal encoder comprises two blocks ( $L_4 = L_5 = 2$ ), with 793 SA and Batch-SA configured to a hidden size of 512 and eight attention heads. During inference, 794 similar to the Vehicle class, we apply MMCAT to our frustum point cloud and use IoU = 0.3. The training process takes 1000 epochs with a batch size of 800 on four Nvidia A100 GPUs. We used 795 the Adam optimizer with a starting learning rate of  $1 \times 10^{-4}$ , a cosine annealing scheduler for 796 adjustments, and a weight decay of 0.05. 797

798 799

B.1 EXPERIMENT ANALYSIS ON KITTI VAL (PEDESTRIAN) SET.

Our experiments show that our MMCAT model extends effectively to additional categories, such as
 Pedestrians. As evidenced in Table 7. MMCAT indicates very promising results in the Pedestrian
 class, surpassing the majority of existing fully supervised models while utilizing only 25% labeled
 data under 2D weak supervision. This verifies MMCAT's good generalizability across diverse object
 classes and efficiency in handling smaller objects.

806 B.2 QUALITATIVE ANALYSIS ON KITTI VAL (PEDESTRIAN) SET.807

In Figure 4, we present visual examples of MMCAT's performance relabeling the KITTI training set for the Pedestrian category. The first two columns (easy samples) illustrate the model's precision in generating 3D bounding boxes for clear, non-occluded pedestrians at a moderate distance

811	Table 7: Results of KITTI Val set (Pedestrian), compared to the fully supervised PointRCNN and
812	other autolabeler.

Method	Modality Full Supervisi	Full Supervision	$AP_{3D}(IoU = 0.5)$			$AP_{BEV}(IoU = 0.5)$		
incurou	inodunty i un supervision		Easy	Moderate	Hard	Easy	Moderate	Hard
PointRCNN Shi et al. (2019)	LiDAR	1	63.70	69.43	58.13	68.89	63.54	57.63
PointPillarsLang et al. (2019)	Lidar		66.73	61.06	56.50	71.97	67.84	62.41
STDYang et al. (2019)	LIDAR	\$ \$	73.90	66.60	62.90	- 75.09	- 69.90	- 66.00
VoxelNetZhou & Tuzel (2018)	LiDAR	1	-	-	-	70.76	62.73	55.05
	Comparison with other Autolabeler							
WS3D Meng et al. (2020)	LiDAR	BEV Centroid	74.65	69.96	66.49	74.99	71.23	67.45
CAT Qian et al. (2023)	LiDAR	2D Box	75.15	70.06	67.09	74.79	71.27	66.75
MMCAT (ours)	LiDAR	2D Box	76.85	72.01	70.88	76.25	73.20	70.01
2D Box	CT	Orien			OT	0		
20 DOK	GI VS.	Jurs	2D B0	x	GIVS	. Ours		
the second second			and the second		<b>N</b>			
	-			W.				
	<b>4</b>				:	XII		
	<b>1</b>		-VAL					
	EE.					****		
	Least Contract of Charmer 1							
				Diant	N.			
		4					12	
						<i>(</i>		
			112285		Conceptor .			
				- Branks	1			
T TIL reformbour				The Cont	Į.			
						:		
				JISS N	- M			
			1					
	-		A.		.*			
The			- 20				i	
			a star	and the second second	1	<b>F</b>	1	
		==	12	No 19	1\ +	· · ·		

Figure 4: Qualitative Analysis with MMCAT performance on the KITTI Training Set (Pedestrian):
Demonstrates MMCAT's robustness and precision in generating 3D bounding boxes, especially under challenging conditions such as substantial truncation, heavy occlusion, or significant distance from the sensor (highlighted in the last two columns). MMCAT effectively generates comprehensive amodal 3D bounding boxes (shown in green) for pedestrian structures.

characterized by dense point clouds. As shown in the last two columns (hard samples), MMCAT
 also accurately identifies heavily occluded or very far pedestrian samples, significantly sparser than
 vehicles. This demonstrates MMCAT's capability to process vehicles and adapt to other categories
 within autonomous driving contexts, even with minimal and more readily available supervision,
 particularly for challenging samples.

865		Table 8: Annotation	n Speed and Tr	aining Time.		
866		Methods	Inference (s)	Training (h)	Parameters (M)	
867		Humman Annotation Song et al. (2015)	114	-	-	
868		WS3D Meng et al. (2021b)	2.5	10.0	-	
869		MTrans Liu et al. (2022a)	0.04	6.5	-	
870		CAT Qian et al. (2023)	0.02	6.5	4	
871		MMCAT (ours)	0.02	7.0	4.6	
872						
873	<b>B</b> .3	ANNOTATION SPEED.				
874	An ar	alysis of the annotation speed is detaile	d in Table 🛛 N	AMCAT was a	able to relabel the	KITTI
875	traini	ng set which contains 3 712 frames with 1	5654vehicle i	nstances in al	out 3 minutes and	otating
876	at 0.0	2 seconds per instance. In contrast, hum	an annotations	take around	114 seconds per in	istance.
877	or 30	seconds with the assistance of a 3D ob	ject detector,	as reported in	previous studies	Huang
878	et al.	(2019); Song et al. (2015); Meng et al.	(2021b). MM	CAT matches	CAT's annotation	1 speed
879	despi	te being a larger model. This trend is cor	nsistent on the	Waymo, with	MMCAT maintai	ning an
880	annot	ation speed of 0.02 seconds per car, simi	lar to CAT.			
881						
882						
883						
884						
885						
886						
887						
888						
889						
890						
891						
892						
893						
894						
895						
896						
897						
898						
899						
900						
901						
902						
903						
904						
905						
906						
907						
908						
909						
910						
911						
912						
913						
914						
915						
916						
917						