

IDEAL-M3D: Instance Diversity-Enriched Active Learning for Monocular 3D Detection

Johannes Meier^{1,2,3,4,†,*}Florian Günther^{3,†}
Jacques Kaiser¹Riccardo Marin^{3,4}
Daniel Cremers^{3,4}Oussema Dhaouadi^{1,3,4}¹ DeepScenario² ETH Zurich³ TU Munich⁴ MCML

Figure 1. IDEAL-M3D is the first instance-based active learning method for monocular 3D detection. **Left:** While previous active learning approaches select entire images for labeling, we identify the most informative object instances (difference is highlighted in rounded boxes). **Right:** Our approach achieves full supervised performance using only 50-60% of the labeled boxes, significantly outperforming existing active learning methods across the majority of object categories (KITTI [12] validation set, $AP_{3D|R_{40}}^{0.7}$ Mod.).

Abstract

Monocular 3D detection relies on just a single camera and is therefore easy to deploy. Yet, achieving reliable 3D understanding from monocular images requires substantial annotation, and 3D labels are especially costly. To maximize performance under constrained labeling budgets, it is essential to prioritize annotating samples expected to deliver the largest performance gains. This prioritization is the focus of active learning. Curiously, we observed two significant limitations in active learning algorithms for 3D monocular object detection. First, previous approaches select entire images, which is inefficient, as non-informative instances contained in the same image also need to be labeled. Secondly, existing methods rely on uncertainty-based selection, which in monocular 3D object detection creates a bias toward depth ambiguity. Consequently, distant objects are selected, while nearby objects are overlooked.

To address these limitations, we propose IDEAL-M3D, the first instance-level pipeline for monocular 3D detection. For the first time, we demonstrate that an explicitly diverse, fast-to-train ensemble improves diversity-driven active learning for monocular 3D. We induce diversity with heterogeneous backbones and task-agnostic features, loss weight perturbation, and time-dependent bagging. IDEAL-M3D shows

superior performance and significant resource savings: with just 60% of the annotations, we achieve similar or better AP_{3D} on KITTI validation and test set results compared to training the same detector on the whole dataset.

1. Introduction

3D object detection estimates 3D bounding boxes of objects visible in 2D images by predicting their position, orientation, and dimensions. Identifying objects and their geometrical properties is a fundamental step toward understanding the environment, and with the progress in the autonomous driving industry [5, 7], it has gained significant traction. The monocular setting is particularly compelling, since relying on a single RGB camera is flexible, cheap, and simple to set up. Although estimating 3D properties such as depth from a monocular view is theoretically impossible, as discussed in [21, 24, 34, 42, 57], strong data priors can address most of these ambiguities in practice. This leads to reliable and practically useful predictions [15].

As a consequence, annotated data has become a pivotal element in fostering research in the field, and the research community has put tremendous effort into curating large-scale datasets [22, 51, 60, 73]. However, 3D data annotations are labor-intensive, costly, and difficult to scale. To limit such effort, Active Learning (AL) offers a compelling alterna-

*j.meier@tum.de † Equal contribution

tive by identifying the most informative unlabeled samples to annotate. Such selection reduces the redundant annotations [22, 67] and is particularly valuable for Monocular 3D Detection (M3D), where unlabeled data is abundant [60, 73].

Existing AL methods for M3D, such as Efficient AL [16] and MonoLiG [15], operate at the image level. They request entire images rather than specific objects. This induces unnecessary annotation overhead because every object in a selected image must be labeled, including trivial or already well-modeled instances (see Fig. 1). It also obscures the true annotation cost, which scales with the number of instances to be annotated [33]. As shown by Lyu et al. [33], labeling capacity is better utilized when annotating instances rather than images. In addition, existing AL methods for M3D [15, 16] are primarily uncertainty-based and rank samples by low predictive confidence or high disagreement (*e.g.*, ensemble variance). At the instance level, such criteria over-select distant objects, as they exhibit high aleatoric and epistemic uncertainty but deliver limited accuracy gains. We hypothesize that prior work [15, 16] overlooked this bias because image-level selection compels annotating all objects in each chosen image. This mixes informative and uninformative instances and masks the inefficiency.

We revisit AL for M3D from a diversity perspective. First, we formalize an instance-based AL pipeline and make the instance the unit of annotation. Then, we introduce **Instance Diversity-Enriched Active Learning for Monocular 3D Detection (IDEAL-M3D)**: The first instance-based AL method for M3D (*cf.* Tab. 1). A key limitation of prior diversity-based approaches is that they rely on features from a single detector, which introduces model- and task-specific biases. Motivated by information theory, our core idea is to use ensembles for diversity-based selection. We estimate each instance’s diversity jointly over a heterogeneous ensemble and combine complementary representations to reduce single-model bias. While ensembles are common for uncertainty estimation, to the best of our knowledge this is the first use of ensembles to drive diversity-based AL.

Our contributions include several measures to increase the heterogeneity of the ensemble. This involves random loss weight sampling, data sampling and the utilization of different backbones. As an additional benefit, these choices cut the additional training time by over 50% compared to a vanilla ensemble. Furthermore, we add visual features from a pre-trained image-based autoencoder to the ensemble. This adds no extra detector-training time and lets the selection exploit both task-related and task-independent cues.

We validate our approach on KITTI [12], Waymo [51], and additionally on Rope3D [60] for cross-perspective robustness. IDEAL-M3D achieves state-of-the-art results across all datasets. On the KITTI [12] validation and test set, it reaches 100% of the full-data performance using only 60% of the labels. In some cases, it even surpasses the

Table 1. Comparative analysis of key related research and methodologies. Our AL approach is the first instance-based method for Monocular 3D Detection (M3D).

Methods	Boxes	Modality	Instance-based	Diversity-based
ComPAS [33]	2D	Image	✓	✗
QBox [52]	2D	Image	✓	✗
DDFH [4]	3D	LiDAR	✗	✓
Efficient AL [16]	3D	Image	✗	✓
MonoLiG [15]	3D	Image & LiDAR	✗	✗
Ours	3D	Image	✓	✓

fully supervised counterpart. On Rope3D [60] IDEAL-M3D achieves over 97% accuracy with just 25% of data, while on Waymo over 95% with just 25% of data. Finally, we propose **Normalized Area under the Requested Curve (NAURC)**, a novel AL metric that allows for direct comparison across instance- and image-based methods. The contributions of this work are:

- We design the first instance-based AL pipeline for M3D and show that uncertainty-based selection underperforms at the instance level, where it over-selects distant objects with limited benefit.
- We introduce IDEAL-M3D, a diversity-based selector with highly diverse ensemble features that reduces training overhead by over 50% compared to the vanilla ensemble. Furthermore, we show that it leads to a complementary performance gain when these task-related detector features are combined with task-independent visual features.
- Our approach is simple to implement yet achieves state-of-the-art performance on KITTI [12], Waymo [51], and Rope3D [60].

2. Related Work

2.1. Monocular 3D Object Detection (M3D)

Monocular 3D Detection (M3D) aims to recover 3D information of objects from single 2D images [24, 30, 57, 71]. The research community has made a sustained effort to curate datasets and benchmarks for it, such as KITTI [12], Waymo [51], CDrone [36], Rope3D [60] and V2X-I [62, 63]. Such data often comes with additional LiDAR information, useful for providing additional depth supervision during training (*e.g.*, RD3D [69], CaDDN [45], MonoNeRD [56], OccupancyM3D [39], MonoTAKD [26]). However, LiDAR is not always available [9, 73], limiting the scalability of the approach. Methods that rely solely on 3D bounding box annotations such as MonoCon [27], MonoDETR [70],

MonoCD [57], MonoMAE [19] and MonoUNI [18] are compelling since they use only RGB information. To compensate for missing 3D information, these methods require a large number of annotated images. Among the most recent approaches, MonoDiff [43] frames the task as a denoising diffusion process, and MonoLSS [23] employs Gumbel Softmax to identify depth-relevant features. MonoLSS has particularly strong performance and independence from LiDAR supervision. We primarily evaluate on the standard benchmarks, KITTI [12] and Waymo [51], and additionally include Rope3D [60] to assess generalization to traffic-view scenarios.

2.2. Active Learning

AL reduces annotation costs by selecting the most informative samples for labeling. Traditional AL methods can be divided into uncertainty-based approaches (e.g., Maximum Entropy [48], BALD [11]), diversity-based approaches (e.g., Core-Set [47], DiscAL [13]), and hybrid methods [14, 67]. For an exhaustive survey on the topic, we point to [67]. However, these methods are primarily designed for classification tasks, whereas M3D is a regression problem.

Active Learning for 2D Object Detection. PPAL [58] combines weighted entropy scores with diversity filtering, while MI-AOD [64] first performs adversarial training and later measures uncertainty via image-instance inconsistencies. During our investigation, we found that uncertainty-based methods are not suitable for instance-based AL in 3D object detection since they tend to be biased toward the farthest objects in the scene.

Active Learning for 3D Object Detection. In LiDAR-based 3D detection, AL methods such as STONE [35] and CRB [32] identify representative prototypes in gradient space. STONE [35] incorporates uncertainty estimation via Monte Carlo dropout, while CRB [32] optimizes for uniform point density. DDFH [4] compresses high-dimensional features and bounding box information using t-SNE [53] and employs Gaussian mixture models for diverse sample selection. KECOR [31] leverages neural tangent kernels to quantify sample uniqueness. These methods address LiDAR-specific challenges such as point density diversity [4, 32]. In contrast, monocular image-based detection requires different handling strategies as it confronts different challenges, like the estimation of depth from a single image.

For M3D, AL remains underexplored. Efficient Active Learning [16] estimates epistemic uncertainty using heatmaps and maximizes diversity using a 2D detector trained on MS COCO [25] but fails to integrate highly task-related features. MonoLiG [15] is an uncertainty-based approach that relies on a LiDAR-trained teacher and an ensemble of five student models. This design imposes high training costs, yet it lacks mechanisms to reduce the ensemble overhead or to explicitly encourage ensemble diver-

sity. In our experiments, it still underperforms our approach. In contrast, we eliminate LiDAR dependency and adopt a diversity-driven selection tailored to M3D.

Instance-based Active Learning. All mentioned AL methods select entire scenes or images. This requires annotating all the objects in the image, even when they are not relevant. On the contrary, instance-based approaches aim to identify the individual objects which are the most informative to annotate. Didari et al. [10] and ViewAL [49] design specific techniques for segmentation by selecting pixels or superpixels in the images. In 2D object detection, methods like ComPAS [33] and QBox [52] rely on uncertainty, for instance selection, which, as already mentioned, is suboptimal for M3D. Instead, we use a representative-based selection.

3. IDEAL-M3D

We provide a high-level overview of IDEAL-M3D (Fig. 2). In Sec. 3.1, we present our adaptation of the AL pipeline to instance-based selection. The subsequent sections detail our instance selection strategy: Sec. 3.2 adapts our diversity-based baseline Core-Set [47] to M3D, while Sec. 3.3 and Sec. 3.4 show how diverse ensembles and task-agnostic features synergistically broaden feature-space coverage for more effective selection.

3.1. Instance-Based AL for M3D

Initialization. We begin by labeling a small random subset of images to initialize the detector. Then we train our model over multiple AL rounds. In each round, the network runs inference on all unlabeled images. We then use the predictions (in particular, the 2D bounding box center) to propose candidates for labeling until the labeling budget of the current round is exhausted.

Processing of Unlabeled Instances. During selection, 3D box locations are an unreliable matching criterion as depth can be very imprecise, especially at early iterations or at high distances. Instead, we identify the object to annotate as the one whose 2D center is closest to the requested 2D center. To avoid requiring the annotator to scan the entire image, we restrict the search to a depth-dependent window of size (r_x, r_y) around the predicted center; see the supplementary for the precise definition of this radius. We then select instances to label in the amounting to 5–10% of the global budget per round. This selection lies at the core of our approach and is detailed in the next section.

Training. After selecting training instances, we update the model with the loss used by the baseline detector (MonoLSS [23], MonoCon [27]). Because only a small set of instances is labeled, the classification head tends to assign unlabeled objects to the background, penalizing rare classes. Annotating all objects per image (as in image-based AL) mitigates this but is costly; in instance-based AL, rare classes may be missed.

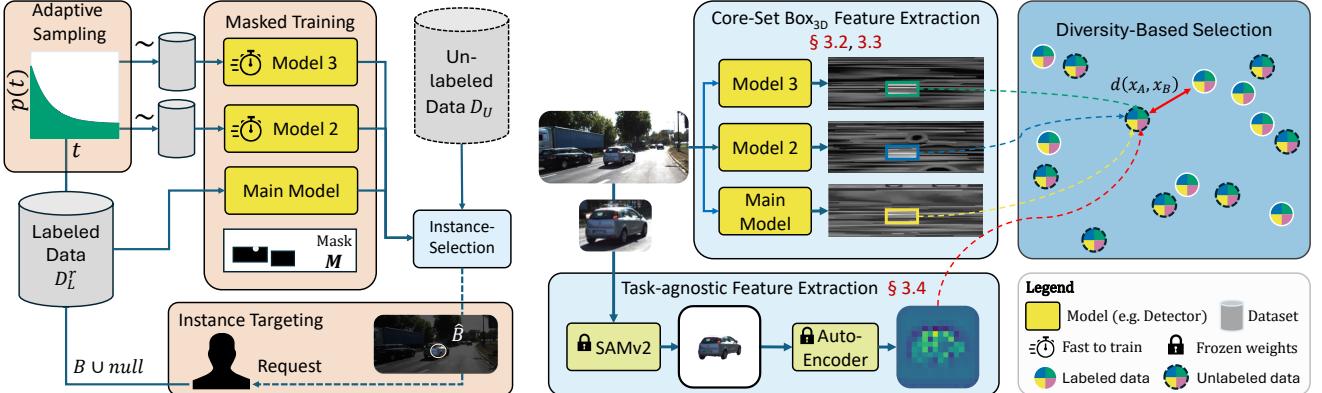


Figure 2. **Overview of IDEAL-M3D.** **Left:** Our instance-based AL pipeline couples precise instance targeting with time-adaptive sampling, minimizing expert effort while remaining training-time efficient (Sec. 3.1). **Right:** We maximize feature-space coverage by fusing Core-Set selection with an explicitly diverse, fast-to-train ensemble and task-agnostic visual embeddings, yielding robust geometry-aware selection under modest compute (Secs. 3.2 to 3.4). IDEAL-M3D uniquely integrates diversity-based selection with an ensemble purpose-built for representational diversity in M3D, delivering label efficiency without the cost of conventional ensembles.

Hence, we introduce a binary mask $M_{i,j}$ for the classification loss to ignore uncertain regions. Let the predicted probability at pixel (i, j) for class c be $\hat{p}_{c,i,j}$ and the ground-truth be $p_{c,i,j}$. We define the binary Mask M as:

$$M_{i,j} = \begin{cases} 0 & \text{if } \sum_c p_{c,i,j} = 0 \text{ or } (i, j) \text{ falls inside the} \\ & \text{2D box of a previously predicted object,} \\ & \text{that is still unlabeled} \\ 1 & \text{otherwise.} \end{cases} \quad (1)$$

We compute the masked objectness loss, where ℓ denotes the weighted focal loss [72], as:

$$L_{cls}^{\text{masked}} = \sum_c \sum_{i,j} M_{i,j} \cdot \ell(\hat{p}_{c,i,j}, p_{c,i,j}). \quad (2)$$

3.2. Core-Set Box_{3D}

In our instance-based regime, popular uncertainty-based methods underperform (*cf.* Tab. 2): they often select far-away objects with high epistemic or aleatoric uncertainty, which are less useful for M3D than closer, well-resolved instances (*cf.* Sec. 7.5). These observations motivate a diversity-based baseline. We find that Core-Set [47], a simple diversity-based approach, performs surprisingly well in this setting, so we adopt it as our baseline.

Core-Set selects the unlabeled instance farthest from the labeled set in feature space. At labeling round r , with features $f(x)$, labeled set \mathcal{D}_L^r , and unlabeled set \mathcal{D}_U^r , it scores an unlabeled instance x by

$$\text{score}(x) = \min_{z \in \mathcal{D}_L^r} \|f(x) - f(z)\|_2, \quad (3)$$

and choose the next query

$$x^* = \arg \max_{x \in \mathcal{D}_U^r} \text{score}(x). \quad (4)$$

This greedy step repeats until the round budget is exhausted.

A first weakness of the classic variant is that $f(x)$ is a penultimate classification feature. This is suboptimal for M3D because they are trained to be invariant to geometric factors (size, depth, pose) that dominate 3D detection quality. We therefore extract features immediately before the 3D detection heads, flatten them, and use these as $f(x)$ in the Core-Set distance; we refer to this instance-level adaptation as *Core-Set Box_{3D}*. These pre-head features retain the 3D cues required for 3D box prediction, thereby aligning the distance metric with the information actually used for M3D.

3.3. Diverse ensembles

Motivation. One important limitation of Core-Set [47] (and inherited by *Core-Set Box_{3D}*) is the reliance on a single representation $f(\cdot)$. A single model's inductive biases and training dynamics (e.g., feature geometry shaped by loss composition, augmentations, and optimization) can dominate the notion of “distance,” risking that other task-relevant modes are overlooked. We remove this single-model bias by leveraging an ensemble. Prior work has used ensembles for uncertainty-based scoring; to our knowledge, we are the first to deploy an ensemble explicitly for diversity-based selection in M3D. Concretely, we replace $f(\cdot)$ with 3 detector-specific embeddings of the same instance x , denoted $f_1(x), f_2(x), f_3(x)$, and measure distances with a weighted cosine across views:

$$d(x_A, x_B) = \sum_{i \in \{1, 2, 3\}} \lambda_i \cdot d_{\cos}(f_i(x_A), f_i(x_B)), \quad (5)$$

where $\lambda_i \geq 0$ and $d_{\cos}(u, v) = 1 - \frac{u^\top v}{\|u\|_2 \|v\|_2}$. We then apply the Core-Set [47] selection in this fused space:

$$D_{\text{ens}}(x | \mathcal{D}_L^r) = \min_{z \in \mathcal{D}_L^r} d(x, z). \quad (6)$$

$$x^* = \arg \max_{x \in \mathcal{D}_U^r} D_{\text{ens}}(x \mid \mathcal{D}_L^r). \quad (7)$$

Information Theory Intuition. For intuition, we analyze the ensemble representation as a concatenation $F_{\text{ens}}(x) = [f_1(x), f_2(x), f_3(x)]$ and omit x for readability. By the chain rule of mutual information (MI),

$$\begin{aligned} I(F_{\text{ens}}; \mathcal{D}_L^r) &= I(f_1; \mathcal{D}_L^r) + \sum_{m=2}^3 I(f_m; \mathcal{D}_L^r \mid f_{<m}) \\ &\geq I(f_1; \mathcal{D}_L^r). \end{aligned} \quad (8)$$

Because the concatenation order is arbitrary, we can permute the views so that any f_m appears first, yielding:

$$I(F_{\text{ens}}; \mathcal{D}_L^r) \geq \max_{m \in \{1, 2, 3\}} I(f_m; \mathcal{D}_L^r). \quad (9)$$

The result in (8) and (9) gives strict gains when added views contribute a nonredundant signal.

Efficiency. Ensembles are often weakly diverse, since their members often differ mainly by random initialization and data-loading order. Also, they are slow to train. To address both issues, we train one main model in standard fashion (also used for evaluation) and add two auxiliary models that are faster and intentionally diversified. We introduce the following measures to increase representational diversity while reducing training time:

- **Heterogeneous, smaller backbones.** For the auxiliary models, we choose different backbones to induce diversity via both distinct architectures and initializations from their respective pretrained weights. For fairness, all detectors (main and auxiliary) are initialized from ImageNet [8] pretraining. In our experiments, we use RepViT-M1.0 [55] and MobileNet4-Conv-M [41], selected for (i) reasonable accuracy and (ii) faster training than the main model (see Fig. 12).
- **Shorter schedules.** Each auxiliary model trains for one third of the main model’s epochs, retaining most of the diversity signal while substantially reducing compute.
- **Multi-task loss weighting perturbations.** M3D comprises several subtasks (2D box, 3D offset, dimensions, depth, orientation, confidence, classification). At the start of each round r , for each auxiliary model m we sample per-subtask Ψ loss multipliers $w_{r,\Psi}^m \sim \text{Uniform}(1 - \delta, 1 + \delta)$ to induce mild specialization.
- **Time-dependent bagging.** Later rounds dominate training cost because the labeled set has grown, whereas early rounds are label-sparse and benefit from using a larger fraction of available labels; accordingly, we adopt a decreasing subsampling schedule over training progress. At training round r with progress $t \in [0, 1]$

(where $t \triangleq \frac{|\mathcal{D}_L^r|}{|\mathcal{D}_L^r| + |\mathcal{D}_U^r|}$), we subsample the labeled data with:

$$s(t) = 0.5 + 0.4 \exp(-\alpha t), \quad \alpha > 0, \quad (10)$$

which monotonically decreases the sampled fraction over time, using more data early and less later.

3.4. Task-agnostic features

Even with a multimodel ensemble, relying solely on task-specific features can bias the metric toward detector idiosyncrasies and under-represent appearance variation. To complement these signals, we augment the ensemble feature space with task-agnostic visual features that capture object appearance (e.g., texture, color, shape, local context).

Concretely, for each candidate instance we segment the object with SAMv2 [44] and encode the masked crop using a pretrained Stable Diffusion autoencoder [46] to obtain a compact flattened visual embedding $f_{\text{vis}}(x)$. Since the SAM mask is only used for feature selection and not for the final selection itself, a slight imperfection in the mask is not critical. We then integrate these features into our cosine-based distance and select the highest-scoring instances for labeling:

$$d(x_A, x_B) = \sum_{i \in \{1, 2, 3, \text{vis}\}} \lambda_i \cdot d_{\text{cos}}(f_i(x_A), f_i(x_B)). \quad (11)$$

Because the visual embeddings are obtained from frozen models (no additional training), the added computation is negligible relative to detector training. Overall, our diversity-based selection is straightforward to implement, shown in experiments to be effective, and markedly more training-time efficient than conventional ensembles.

4. Experiments

4.1. Datasets, Metrics, and Active Learning Setting

We evaluate on three benchmarks for monocular 3D detection: KITTI [12] and Waymo monocular [45, 51] (both vehicle-mounted), and Rope3D [60] (traffic scenes with diverse camera perspectives [36]). Waymo and Rope3D are large-scale datasets. MonoLSS [23] on KITTI [12] and MonoCon [27] on Waymo [51] and Rope3D [60] are used as baseline detectors.

- **KITTI [12].** Following [6, 23], we split the training set into 3,712 images for training, 3,769 for validation, and 7,518 for testing. We report $AP_{3D|R_{40}}$ [50] (3D Average Precision at 40 recall positions) with the standard Easy/Moderate/Hard difficulty levels, using Moderate as the primary benchmark. For active learning, we start from an initial 10% labeled training pool and increase the labeled set through iterations to 15%, 20%, 30%, 40%, 50%, and 60%.

Table 2. AL performance on KITTI [12] validation, Waymo [45, 51] validation and Rope3D [60] validation dataset. Results are averaged over three rounds, each initialized from the same checkpoint. **KITTI**: We report NAURC_{60%}. **Waymo**: We report NAURC_{25%}. **Rope3D**: We report NAURC_{35%}. The appendix provides a description of each method. **Type***: U=Uncertainty-based, D=Diversity-based, H=Hybrid.

Method	Type*	KITTI [12] Car			Waymo [51] Vehicle				Rope3D [60]				
		$AP_{3D R_{40}}^{0.7}$			IoU=0.5		IoU=0.7		Car		Big Vehicle		
		Easy	Mod.	Hard	AP_{3D}	APH_{3D}	AP_{3D}	APH_{3D}	$AP_{3D}^{0.5}$	Rope	$AP_{3D}^{0.5}$	Rope	
Image-based	Rand	-	18.01	12.97	10.77	8.07	8.00	1.72	1.70	32.35	44.84	13.65	28.27
	Conf	U	19.72	<u>14.20</u>	<u>11.80</u>	8.03	7.96	1.74	1.72	31.80	44.43	14.05	28.77
	Ens Depth Var	U	18.19	13.00	10.77	7.68	7.61	1.62	1.61	32.89	45.29	14.44	28.92
	Augm Depth Var	U	18.24	13.16	10.90	-	-	-	-	-	-	-	-
	Core-Set [47]	D	18.80	13.50	11.34	8.10	8.03	1.76	1.75	32.21	44.73	14.48	29.09
	BADGE [3]	H	19.02	13.53	11.39	8.14	8.08	1.76	1.75	32.83	45.22	14.05	28.61
	DDFH [4]	H	17.94	12.72	10.51	4.50	4.46	0.97	0.96	20.06	27.55	10.11	18.67
	CDAL [1]	H	18.38	12.99	10.82	7.80	7.73	1.67	1.66	30.96	42.43	12.18	25.97
Instance-based	Rand	-	18.69	13.53	11.21	8.19	8.12	1.85	1.84	30.71	43.55	10.75	26.14
	Conf	U	12.10	8.93	7.60	8.06	8.00	1.83	1.82	32.50	44.97	11.77	26.87
	Ens Depth Var	U	11.88	9.07	7.82	8.79	8.72	2.05	2.04	32.43	44.90	13.56	28.32
	Augm Depth Var	U	12.25	8.96	7.64	-	-	-	-	-	-	-	-
	ComPAS [33]	U	17.31	12.31	10.48	8.63	8.56	1.96	1.94	32.66	45.12	16.95	26.87
	Core-Set [47] Box _{3D}	D	18.83	13.82	11.68	8.84	8.77	2.02	2.01	32.66	45.09	<u>16.95</u>	<u>31.12</u>
	BADGE [3]	H	17.39	12.63	10.74	<u>8.86</u>	<u>8.79</u>	<u>2.07</u>	<u>2.06</u>	<u>32.92</u>	<u>45.32</u>	13.76	28.58
	IDEAL-M3D (Ours)	D	22.74	16.18	13.57	9.02	8.91	2.11	2.10	34.27	46.41	18.45	32.34

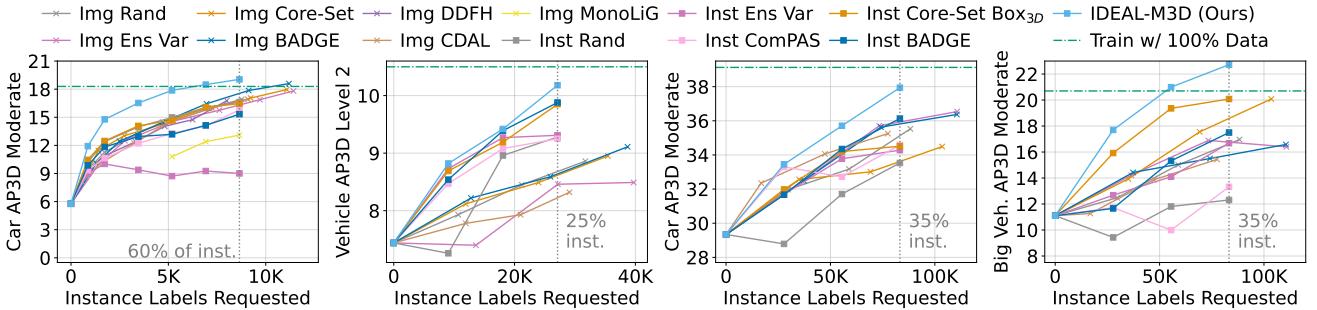


Figure 3. AL training curves. **Plot 1**: We report the KITTI validation [12] performance for cars $AP_{3D|R_{40}}^{0.7}$ Moderate. **Plot 2**: We report the Waymo validation [51] performance $AP^{0.5}$ for vehicles. **Plots 3-4**: We report the Rope3D validation performance $AP_{3D|R_{40}}^{0.5}$ for Cars and Big Vehicles. All results show mean performance across three rounds with identical initialization.

- **Waymo** [51]. Following the monocular setting [45], we use front-camera images with 52,386 training and 39,848 validation images. We report AP_{3D} and APH_{3D} at IoU thresholds 0.5 and 0.7, and follow MonoLiG [15] in using difficulty level 2. We begin with 10% labeled images and acquire an additional 5% per iteration up to a 25% maximum budget.
- **Rope3D** [60]. We use the heterologous split (changing camera views), comprising 40,333 training and 4,676 validation images. Performance is evaluated using $AP_{3D|R_{40}}$ [50] and the Rope Score [60] at IoU 0.5 and Moderate difficulty. We begin with 20% labeled images and acquire an additional 5% per iteration up to a 35% maximum budget.

4.2. NAURC: A budget-fair evaluation metric

The goal of AL is to maximize accuracy under a limited labeling budget. For comparability with prior work, we report our main results using the two standard strategies: training curves over requested instances [13, 15, 33, 47, 58, 64, 67] and performance at fixed labeled-data percentages [16, 33, 35]. These approaches, however, have notable drawbacks: training curves can be ambiguous because leadership changes across the budget, and fixed-percentage snapshots provide only a single number that depends on the chosen percentage. They also do not resolve the budget-unit mismatch between image-based methods (spending in images) and instance-based methods (spending in instances).

To address these issues, we propose a new metric,

Normalized Area under the Requested Curve (NAURC). NAURC is a single-scalar metric that integrates performance over the requested-instance budget, normalizes by a target budget, and enables direct comparison between image- and instance-based methods. It enables a common instance-based accounting for all methods: If an image-based method overshoots the target budget, we interpolate back to the budget; if it undershoots, we keep the last observed performance. We refer to the supplementary for a formal definition and more detailed motivation, derivation, and visualizations

4.3. Comparison with AL methods

In Tab. 2 and Fig. 3, we report average accuracy over three runs initialized from the same checkpoint, comparing IDEAL-M3D with other AL methods. Interestingly, methods based on uncertainty significantly drop their KITTI [12] and Rope3D [60] performance when they move from image-based to instance-based approaches. Intuitively, this is a result of their bias toward distant objects (Sec. 7.5), which can be mitigated by annotating all the objects in the image. Remarkably, with only 60% of instances, we surpass the 100% KITTI baseline, and on Rope3D we reach 97% (car) and 107% (big vehicle) of fully supervised AP performance using just 35% of labels. Results for the pedestrian and cyclist class are shown in the supplementary material.

Waymo [51] shows the opposite trend: all instance-based methods improve. We hypothesize that this stems from strong redundancy in the dataset, as Waymo provides video frames at 300 ms intervals [45], whereas KITTI and Rope3D are image-based. On Waymo, we match 97% ($AP^{0.5}$) and 95% ($AP^{0.7}$) supervised performance with 75% fewer labels. These quantitative results are supported by qualitative evidence in Fig. 4. We observe that IDEAL-M3D tends to prioritize nearby objects at the beginning of AL and gradually moves to more distant ones. Additionally, the detection quality generally shows consistent improvement over time.

4.4. Comparison with fully supervised methods

To assess generalizability, we run active learning up to 60% of the KITTI trainval instances and evaluate on the KITTI test set [12]. Test labels are not public, so evaluation requires submission to the official server. Competing fully supervised baselines are trained on 100% of trainval.

In Tab. 3, we compare IDEAL-M3D with recent fully supervised and semi-supervised approaches. Relative to our fully trained baseline detector MonoLSS [23], IDEAL-M3D achieves slightly lower AP on Moderate (-0.28) and Hard (-0.21). Yet, it attains higher AP on Easy ($+0.95$) while using 40% fewer labels. IDEAL-M3D is competitive with state-of-the-art under substantially reduced annotation. Several entries are semi-supervised: Mix-Teaching [59], DPL_{FLEX} [68], and MonoLiG [15] leverage additional unlabeled data during

Table 3. Comparison with state-of-the-art (SOTA) monocular methods on the KITTI [12] test set for the car category. **SSL** denotes methods that additionally require unlabeled data to perform semi-supervised learning. **AL** denotes methods that use active learning.

Method	SSL /AL	Train data	Car Test AP ⁷⁰ _{3D R40}		
			Easy	Mod	Hard
Mix-Teaching, [59], TCSV 23	✓/✗	100%	26.89	18.54	15.79
DPL _{FLEX} [68], CVPR 24	✓/✗	100%	24.19	16.67	13.83
MonoLiG [15], WACV 24	✓/✓	100%	24.90	18.86	16.79
MonoCD [57], CVPR 24	✗/✗	100%	25.53	16.59	14.53
MonoMAE [19], NIPS 24	✗/✗	100%	25.60	18.84	16.78
MonoDGP [40], CVPR 25	✗/✗	100%	26.35	18.72	15.97
GATE3D [17], CVPRW 25	✗/✗	100%	26.07	18.85	16.76
MonoLSS, 3DV 24 (Baseline)	✗/✗	100%	26.11	19.15	16.94
IDEAL-M3D 60% (Ours)	✗/✓	60%	27.06	<u>18.87</u>	16.73

training. MonoLiG [15] further uses LiDAR during training. In contrast, IDEAL-M3D uses RGB images only.

4.5. Ablation Study

We perform our main ablation study in Tab. 4 using NAURC_{60%} AP Mod. as the primary metric. Surprisingly, moving from image-based Core-Set [47] to a naïve instance-based Core-Set lowers performance from 13.50 to 13.16 AP Mod. We hypothesize two causes. First, monocular 3D detection is dominated by 3D box estimation (especially depth), while penultimate classification features over-emphasize semantics and are invariant to geometry. Second, image-based selection may benefit from an averaging effect: labeling one image annotates all objects, so even suboptimal choices often include informative instances. Aligning the selection space with the detection task, i.e. replacing classification features with pre-head detection features (*Core-Set Box*_{3D}), more than recovers this loss.

Building on this representation, the diverse ensemble (Sec. 3.3) contributes +0.72 AP Mod. Adding task-agnostic visual embeddings (Sec. 3.4) yields a larger +1.27 AP Mod. This shows that appearance cues complement task-conditioned geometry. The visual pathway adds negligible overhead because embeddings are computed with frozen encoders (Tab. 5). This preserves our favorable runtime profile. Using both sources together, detector-specific geometry and detector-agnostic appearance, adds a further +1.64 and +1.09 AP Mod. This aligns with our mutual-information view of non-redundant signals. Notice, that in an image-based selection setting, our approach still exceeds the classical Core-Set baseline by +1.33 AP Mod.

4.6. Training-time efficiency

In Tab. 5, we report the training-time ablation on the KITTI validation set [12]. The *Core-Set Box*_{3D} baseline trains for 18.2h. IDEAL-M3D requires 33.7h in total. Using an ensemble of three models therefore increases training time by only 85% relative to *Core-Set Box*_{3D}. In contrast, a

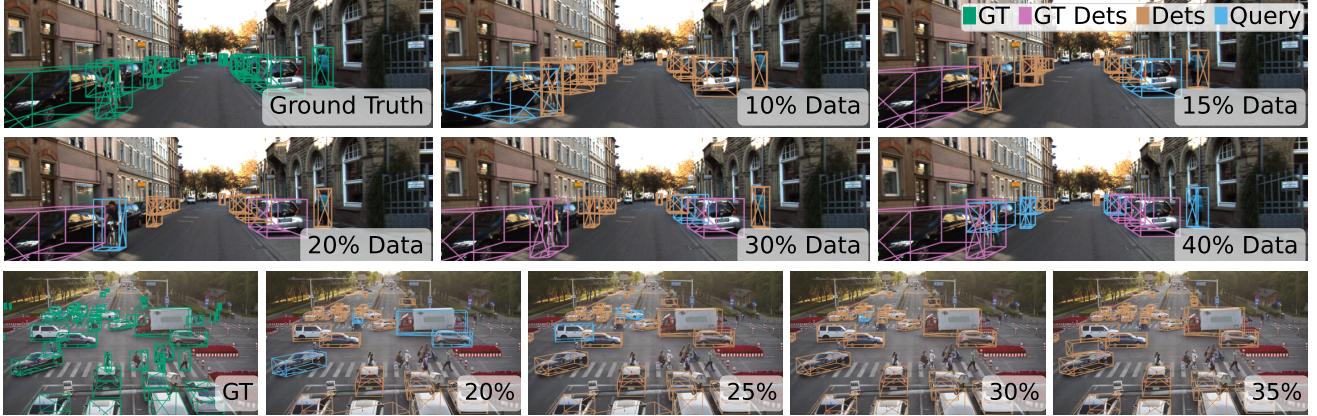


Figure 4. **Qualitative results of IDEAL-M3D on the KITTI [12] (first/second row), and Rope3D [60] (third row) datasets.** The results demonstrate prediction evolution and label selection strategy across time steps. Color coding: green boxes represent ground truth annotations, pink boxes indicate predictions on previously labeled objects, cyan boxes highlight predictions selected for the next labeling round, and orange boxes show predictions that remain unlabeled (best viewed in color with zoom).

Table 4. **Main ablation study on the KITTI [12] validation set (Cars, $AP_{3D|R_{40}}$, $\text{IoU}=0.7$).** **Inst:** Instance-based AL. **Box_{3D}:** Usage of backbone instead of classification features. **DE:** Diverse Ensemble. **VD:** Visual diversity. **Easy/Moderate/Hard:** We report NAURC_{60%}. **Final AP:** Moderate AP after training on 60% of the data. For image-based methods we report the interpolated result.

Inst	Box _{3D}	DE	VD	Easy	Moderate	Hard	Final AP
✓				18.80	13.50	11.40	16.75
✓				18.28	13.16	10.94	16.37
✓	✓			18.83	13.82	11.68	16.45
✓	✓	✓		20.55	14.54	12.28	18.30
✓	✓		✓	20.79	15.09	12.54	18.51
	✓	✓	✓	20.77	14.83	12.34	17.96
✓		✓	✓	20.21	14.96	12.53	17.82
✓	✓	✓	✓	22.74	16.18	13.57	19.04

Table 5. Runtime ablation study on the KITTI [12] validation set.

Method	Total training time
Inst. <i>Core-Set Box_{3D}</i>	18.2h
Ours w/o ensemble	18.7h
Ours w/o visual diversity	33.3h
Ours w/o SAMv2	33.4h
IDEAL-M3D (Ours)	33.7h
Ours w/o diverse backbones	37.0h
Ours w/o data sampling	40.9h
Ours w/ full epochs	41.3h
Inst Ens Depth Var	54.7h

vanilla ensemble like *Inst Ens Depth Var* would naively incur a 200% increase (approximately 3× the baseline). Among our efficiency measures, shortening the training schedule (fewer epochs) is the most effective. Time-dependent bagging (subsampling more aggressively later) provides the next-largest savings. Smaller, heterogeneous backbones further reduce cost. Together, these measures offset most of the

multi-model overhead. Finally, adding task-agnostic visual diversity increases total time by only 1% because the visual embeddings are computed with frozen models and require no additional detector training.

5. Conclusion

We introduced IDEAL-M3D, the first instance-based active learning framework for monocular 3D object detection. We show that diversity-driven ensembles are highly effective for instance selection when representational diversity is explicitly maximized. Our design relies on simple mechanisms that increase diversity without heavy training overhead. We also integrate appearance coverage via task-agnostic visual embeddings at negligible cost. Overall, our approach improves label efficiency while keeping computation modest.

Extensive experiments support these claims across three datasets: KITTI, Waymo, and Rope3D, including results on the private KITTI test set. On KITTI, we match fully supervised accuracy using only 60% of the labeled data. The pipeline is simple to implement and robust in practice.

A limitation is that our approach does not explicitly address extreme class imbalance, where ultra-rare categories may be undersampled. Future work will incorporate multi-modal cues (e.g., LiDAR-derived multi-view images, video) to further diversify signals, and explore semi-supervised learning where feature similarity can guide high-quality pseudo-label selection.

Acknowledgments This work is a result of the joint research project STADT:up. The project is supported by the German Federal Ministry for Economic Affairs and Climate Action (BMWK), based on a decision of the German Bundestag. The author is solely responsible for the content of this publication. This work was also supported by the ERC Advanced Grant SIMULACRON, the Georg Nemetschek Institute project AI4TWINNING and the DFG project 4D-YouTube CR 250/26-1.

References

- [1] Sharat Agarwal, Himanshu Arora, Saket Anand, and Chetan Arora. Contextual diversity for active learning. In *ECCV*, volume 12361, pages 137–153, 2020. [6](#), [15](#), [18](#)
- [2] David Arthur and Sergei Vassilvitskii. k-means++: the advantages of careful seeding. In Nikhil Bansal, Kirk Pruhs, and Clifford Stein, editors, *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2007, New Orleans, Louisiana, USA, January 7-9, 2007*, pages 1027–1035, 2007. [17](#)
- [3] Jordan T. Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. In *ICLR*, 2020. [6](#), [17](#), [18](#)
- [4] Huang-Yu Chen, Jia-Fong Yeh, Jiawei, Pin-Hsuan Peng, and Winston H. Hsu. Distribution discrepancy and feature heterogeneity for active 3d object detection. In *CoRL*, 2024. [2](#), [3](#), [6](#), [15](#), [17](#), [18](#)
- [5] Li Chen, Penghao Wu, Kashyap Chitta, Bernhard Jaeger, Andreas Geiger, and Hongyang Li. End-to-end autonomous driving: Challenges and frontiers. *IEEE Trans. Pattern Anal. Mach. Intell.*, 46(12):10164–10183, 2024. [1](#)
- [6] Xiaozhi Chen, Kaustav Kundu, Yukun Zhu, Andrew G Berneeshawi, Huimin Ma, Sanja Fidler, and Raquel Urtasun. 3d object proposals for accurate object class detection. *Advances in neural information processing systems*, 28, 2015. [5](#)
- [7] Pranav Singh Chib and Pravendra Singh. Recent advancements in end-to-end autonomous driving using deep learning: A survey. *IEEE Trans. Veh.*, 9(1):103–118, 2024. [1](#)
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009. [5](#)
- [9] Oussema Dhaouadi, Johannes Meier, Luca Wahl, Jacques Kaiser, Luca Scalerandi, Nick Wandelburg, Zhuolun Zhou, Nijanthan Berinpanathan, Holger Banzhaf, and Daniel Cremers. Highly accurate and diverse traffic data: The deepscenario open 3D dataset. In *IV*, pages 377–384, 2025. [2](#)
- [10] Sima Didari, Wenjun Hu, Jae Oh Woo, Heng Hao, Hankyu Moon, and Seungjai Min. Bayesian active learning for semantic segmentation. *CoRR*, 2024. [3](#)
- [11] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *ICML*, volume 70, pages 1183–1192, 2017. [3](#)
- [12] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? The KITTI vision benchmark suite. In *CVPR*, pages 3354–3361, 2012. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [8](#), [13](#), [14](#), [15](#), [16](#), [17](#), [18](#), [19](#), [20](#), [22](#)
- [13] Daniel Gissin and Shai Shalev-Shwartz. Discriminative active learning. *CoRR*, abs/1907.06347, 2019. [3](#), [6](#), [13](#)
- [14] Yinan He, Lile Cai, Jingyi Liao, and Chuan-Sheng Foo. Hybrid active learning with uncertainty-weighted embeddings. *Trans. Mach. Learn. Res.*, 2024, 2024. [3](#)
- [15] Aral Hekimoglu, Michael Schmidt, and Alvaro Marcos-Ramiro. Monocular 3d object detection with lidar guided semi supervised active learning. In *WACV*, pages 2335–2344, 2024. [1](#), [2](#), [3](#), [6](#), [7](#), [13](#), [17](#), [19](#)
- [16] Aral Hekimoglu, Michael Schmidt, Alvaro Marcos-Ramiro, and Gerhard Rigoll. Efficient active learning strategies for monocular 3d object detection. In *IV*, pages 295–302, 2022. [2](#), [3](#), [6](#), [13](#), [17](#)
- [17] Eunsoo Im, Jung Kwon Lee, and Changhyun Jee. Gate3d: Generalized attention-based task-synergized estimation in 3d*. In *CVPRW*, 2025. [7](#)
- [18] Jinrang Jia, Zhenjia Li, and Yifeng Shi. Monouni: A unified vehicle and infrastructure-side monocular 3d object detection network with sufficient depth clues. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *NeurIPS*, 2023. [3](#)
- [19] Xueying Jiang, Sheng Jin, Xiaoqin Zhang, Ling Shao, and Shijian Lu. Monomae: Enhancing monocular 3d detection through depth-aware masked autoencoders. In *NeurIPS*, 2024. [3](#), [7](#)
- [20] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. [15](#)
- [21] Abhinav Kumar, Garrick Brazil, Enrique Corona, Armin Parhami, and Xiaoming Liu. DEVIANT: depth equivariant network for monocular 3d object detection. In *ECCV*, volume 13669, pages 664–683, 2022. [1](#)
- [22] Dongyuan Li, Zhen Wang, Yankai Chen, Renhe Jiang, Weiping Ding, and Manabu Okumura. A survey on deep active learning: Recent advances and new frontiers. *CoRR*, 2024. [1](#), [2](#)
- [23] Zhenjia Li, Jinrang Jia, and Yifeng Shi. Monolss: Learnable sample selection for monocular 3d detection. In *International Conference on 3D Vision, 3DV*, pages 1125–1135, 2024. [3](#), [5](#), [7](#), [12](#), [13](#), [15](#), [17](#), [19](#)
- [24] Zhuoling Li, Zhan Qu, Yang Zhou, Jianzhuang Liu, Haoqian Wang, and Lihui Jiang. Diversity matters: Fully exploiting depth clues for reliable monocular 3d object detection. In *CVPR*, pages 2781–2790, 2022. [1](#), [2](#)
- [25] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV*, volume 8693, pages 740–755, 2014. [3](#)
- [26] Hou-I Liu, Christine Wu, Jen-Hao Cheng, Wenhao Chai, Shian-Yun Wang, Gaowen Liu, Hugo Latapie, Jhih-Ciang Wu, Jenq-Neng Hwang, Hong-Han Shuai, and Wen-Huang Cheng. Monotakd: Teaching assistant knowledge distillation for monocular 3d object detection. In *CVPR*, pages 22266–22275, 2025. [2](#)
- [27] Xianpeng Liu, Nan Xue, and Tianfu Wu. Learning auxiliary monocular contexts helps monocular 3d object detection. In *AAAI*, pages 1810–1818, 2022. [2](#), [3](#), [5](#), [12](#), [13](#), [15](#)
- [28] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *CVPR*, pages 11966–11976, 2022. [19](#)
- [29] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. [15](#)

- [30] Yan Lu, Xinzhu Ma, Lei Yang, Tianzhu Zhang, Yating Liu, Qi Chu, Junjie Yan, and Wanli Ouyang. Geometry uncertainty projection network for monocular 3d object detection. In *ICCV*, pages 3091–3101, 2021. 2
- [31] Yadan Luo, Zhuoxiao Chen, Zhen Fang, Zheng Zhang, Mahsa Baktashmotagh, and Zi Huang. Kecor: Kernel coding rate maximization for active 3d object detection. In *ICCV*, pages 18233–18244, 2023. 3
- [32] Yadan Luo, Zhuoxiao Chen, Zijian Wang, Xin Yu, Zi Huang, and Mahsa Baktashmotagh. Exploring active 3d object detection from a generalization perspective. In *ICLR*, 2023. 3
- [33] Mengyao Lyu, Jundong Zhou, Hui Chen, Yijie Huang, Dongdong Yu, Yaqian Li, Yandong Guo, Yuchen Guo, Liyu Xiang, and Guiguang Ding. Box-level active detection. In *CVPR*, pages 23766–23775, 2023. 2, 3, 6, 13, 17, 18
- [34] Xinzhu Ma, Yinmin Zhang, Dan Xu, Dongzhan Zhou, Shuai Yi, Haojie Li, and Wanli Ouyang. Delving into localization errors for monocular 3d object detection. In *CVPR*, pages 4721–4730, 2021. 1
- [35] RUIYU MAO, Sarthak Kumar Maharana, Rishabh K Iyer, and Yunhui Guo. STONE: A submodular optimization framework for active 3d object detection. In *NeurIPS*, 2024. 3, 6, 13
- [36] Johannes Meier, Luca Scalerandi, Oussema Dhaouadi, Jacques Kaiser, Araslanov Nikita, and Daniel Cremers. CARLA Drone: monocular 3d object detection from a different perspective. In *German Conference on Pattern Recognition, GCPR*, 2024. 2, 5, 15
- [37] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Kosecka. 3d bounding box estimation using deep learning and geometry. In *CVPR*, pages 5632–5640, 2017. 13
- [38] Maxime Oquab, Timothée Darct, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jégou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision. *TMLR*, 2024, 2024. 20, 21
- [39] Liang Peng, Junkai Xu, Haoran Cheng, Zheng Yang, Xiaopei Wu, Wei Qian, Wenxiao Wang, Boxi Wu, and Deng Cai. Learning occupancy for monocular 3d object detection. In *CVPR*, pages 10281–10292, 2024. 2
- [40] Fanqi Pu, Yifan Wang, Jiru Deng, and Wenming Yang. Monodgp: Monocular 3d object detection with decoupled-query and geometry-error priors. In *CVPR*, 2025. 7
- [41] Danfeng Qin, Chas Leichner, Manolis Delakis, Marco Fornoni, Shixin Luo, Fan Yang, Weijun Wang, Colby R. Banbury, Chengxi Ye, Berkin Akin, Vaibhav Aggarwal, Tenghui Zhu, Daniele Moro, and Andrew G. Howard. Mobilenetv4: Universal models for the mobile ecosystem. In *ECCV*, volume 15098, pages 78–96. Springer, 2024. 5, 13, 19
- [42] Zequn Qin and Xi Li. Monoground: Detecting monocular 3d objects from the ground. In *CVPR*, pages 3783–3792, 2022. 1
- [43] Yasiru Ranasinghe, Deepti Hegde, and Vishal M. Patel. Monodiff: Monocular 3d object detection and pose estimation with diffusion models. In *CVPR*, pages 10659–10670, 2024. 3
- [44] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädl, Chloé Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross B. Girshick, Piotr Dollár, and Christoph Feichtenhofer. SAM 2: Segment anything in images and videos. *CoRR*, 2024. 5, 13, 20
- [45] Cody Reading, Ali Harakeh, Julia Chae, and Steven L. Waslander. Categorical depth distribution network for monocular 3d object detection. In *CVPR*, pages 8555–8564, 2021. 2, 5, 6, 7, 18
- [46] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10674–10685, 2022. 5, 13, 20
- [47] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In *ICLR*, 2018. 3, 4, 6, 7, 13, 18
- [48] Claude E. Shannon. A mathematical theory of communication. *ACM SIGMOBILE Mob. Comput. Commun. Rev.*, 5(1):3–55, 2001. 3
- [49] Yawar Siddiqui, Julien Valentin, and Matthias Nießner. Viewal: Active learning with viewpoint entropy for semantic segmentation. In *CVPR*, pages 9430–9440, 2020. 3
- [50] Andrea Simonelli, Samuel Rota Bulò, Lorenzo Porzi, Manuel Lopez-Antequera, and Peter Kontschieder. Disentangling monocular 3d object detection. In *ICCV*, pages 1991–1999, 2019. 5, 6
- [51] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo open dataset. In *CVPR*, pages 2443–2451, 2020. 1, 2, 3, 5, 6, 7, 13, 14, 16, 18, 23
- [52] Ying-Peng Tang, Xiu-Shen Wei, Borui Zhao, and Sheng-Jun Huang. Qbox: Partial transfer learning with active querying for object detection. *IEEE Trans. Neural Networks Learn. Syst.*, 34(6):3058–3070, 2023. 2, 3
- [53] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008. 3, 17, 20, 21
- [54] Stefan Sylvius Wagner and Stefan Harmeling. Object-aware DINO (oh-a-dino): Enhancing self-supervised representations for multi-object instance retrieval. *CoRR*, 2025. 21
- [55] Ao Wang, Hui Chen, Zijia Lin, Jungong Han, and Guiguang Ding. Repvit: Revisiting mobile cnn from vit perspective. In *CVPR*, pages 15909–15920, 2024. 5, 13, 19, 20
- [56] Junkai Xu, Liang Peng, Haoran Chen, Hao Li, Wei Qian, Ke Li, Wenxiao Wang, and Deng Cai. Mononerd: Nerf-like representations for monocular 3d object detection. In *ICCV*, pages 6791–6801, 2023. 2
- [57] Longfei Yan, Pei Yan, Shengzhou Xiong, Xuanyu Xiang, and Yihua Tan. Monocd: Monocular 3d object detection with

- complementary depths. In *CVPR*, pages 10248–10257, 2024. [1](#), [2](#), [3](#), [7](#)
- [58] Chenhongyi Yang, Lichao Huang, and Elliot J. Crowley. Plug and play active learning for object detection. In *CVPR*, pages 17784–17793, 2024. [3](#), [6](#), [13](#)
- [59] Lei Yang, Xinyu Zhang, Jun Li, Li Wang, Minghan Zhu, Chuang Zhang, and Huaping Liu. Mix-teaching: A simple, unified and effective semi-supervised learning framework for monocular 3d object detection. *IEEE Trans. Circuits Syst. Video Technol.*, 33(11):6832–6844, 2023. [7](#), [19](#)
- [60] Xiaoqing Ye, Mao Shu, Hanyu Li, Yifeng Shi, Yingying Li, Guangjie Wang, Xiao Tan, and Errui Ding. Rope3D: The roadside perception dataset for autonomous driving and monocular 3D object detection task. In *CVPR*, pages 21309–21318, 2022. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [8](#), [13](#), [14](#), [15](#), [17](#), [18](#), [24](#)
- [61] Fisher Yu, Dequan Wang, Evan Shelhamer, and Trevor Darrell. Deep layer aggregation. In *CVPR*, pages 2403–2412. Computer Vision Foundation / IEEE Computer Society, 2018. [13](#), [19](#)
- [62] Haibao Yu, Yizhen Luo, Mao Shu, Yiyi Huo, Zebang Yang, Yifeng Shi, Zhenglong Guo, Hanyu Li, Xing Hu, Jirui Yuan, and Zaiqing Nie. DAIR-V2X: A large-scale dataset for vehicle-infrastructure cooperative 3d object detection. In *CVPR*, pages 21329–21338, 2022. [2](#)
- [63] Haibao Yu, Wenxian Yang, Hongzhi Ruan, Zhenwei Yang, Yingjuan Tang, Xu Gao, Xin Hao, Yifeng Shi, Yifeng Pan, Ning Sun, Juan Song, Jirui Yuan, Ping Luo, and Zaiqing Nie. V2x-seq: A large-scale sequential dataset for vehicle-infrastructure cooperative perception and forecasting. *CoRR*, 2023. [2](#)
- [64] Tianning Yuan, Fang Wan, Mengying Fu, Jianzhuang Liu, Songcen Xu, Xiangyang Ji, and Qixiang Ye. Multiple instance active learning for object detection. In *CVPR*, pages 5330–5339, 2021. [3](#), [6](#), [13](#)
- [65] Xueying Zhan, Qing Li, and Antoni B. Chan. Multiple-criteria based active learning with fixed-size determinantal point processes. *CoRR*, 2021. [14](#)
- [66] Xueying Zhan, Huan Liu, Qing Li, and Antoni B. Chan. A comparative survey: Benchmarking for pool-based active learning. In Zhi-Hua Zhou, editor, *IJCAI*, pages 4679–4686, 8 2021. [14](#)
- [67] Xueying Zhan, Qingzhong Wang, Kuan-Hao Huang, Haoyi Xiong, Dejing Dou, and Antoni B. Chan. A comparative survey of deep active learning. *CoRR*, 2022. [2](#), [3](#), [6](#), [13](#), [14](#)
- [68] Jiacheng Zhang, Jiaming Li, Xiangru Lin, Wei Zhang, Xiao Tan, Junyu Han, Errui Ding, Jingdong Wang, and Guanbin Li. Decoupled pseudo-labeling for semi-supervised monocular 3d object detection. In *CVPR*, pages 16923–16932, 2024. [7](#)
- [69] Qiude Zhang, Chunyu Lin, Zhijie Shen, Nie Lang, and Yao Zhao. Revisiting monocular 3d object detection from scene-level depth retargeting to instance-level spatial refinement, 2024. [2](#)
- [70] Renrui Zhang, Han Qiu, Tai Wang, Ziyu Guo, Ziteng Cui, Yu Qiao, Hongsheng Li, and Peng Gao. Monodetr: Depth-guided transformer for monocular 3d object detection. In *ICCV*, pages 9121–9132, 2023. [2](#)
- [71] Yunpeng Zhang, Jiwen Lu, and Jie Zhou. Objects are different: Flexible monocular 3d object detection. In *CVPR*, pages 3289–3298, 2021. [2](#)
- [72] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *CoRR*, 2019. [4](#), [12](#)
- [73] Xiaosu Zhu, Hualian Sheng, Sijia Cai, Bing Deng, Shaopeng Yang, Qiao Liang, Ken Chen, Lianli Gao, Jingkuan Song, and Jieping Ye. Roscenes: A large-scale multi-view 3d dataset for roadside perception. In *ECCV*, volume 15099, pages 331–347, 2024. [1](#), [2](#)

Contents

1. Introduction	1
2. Related Work	2
2.1. Monocular 3D Object Detection (M3D)	2
2.2. Active Learning	3
3. IDEAL-M3D	3
3.1. Instance-Based AL for M3D	3
3.2. Core-Set Box _{3D}	4
3.3. Diverse ensembles	4
3.4. Task-agnostic features	5
4. Experiments	5
4.1. Datasets, Metrics, and Active Learning Setting	5
4.2. NAURC: A budget-fair evaluation metric	6
4.3. Comparison with AL methods	7
4.4. Comparison with fully supervised methods	7
4.5. Ablation Study	7
4.6. Training-time efficiency	7
5. Conclusion	8
6. IDEAL-M3D: Further Details	12
6.1. Problem Statement	12
6.1.1 Monocular 3D Object Detection (M3D)	12
6.1.2 Active Learning.	12
6.2. Loss Functions	12
6.2.1 Image-level Losses	12
6.2.2 Object-level Losses	13
6.3. Obtaining Diverse Features	13
7. Experiments	13
7.1. NAURC Evaluation Metric	13
7.2. Implementational Details	14
7.3. Implementation Details of Baseline Methods	15
7.4. Comparison with AL methods	17
7.5. Uncertainty vs. Diversity-based Methods	17
7.6. Training time comparison with fully supervised methods	19
7.7. Backbone ablation	19
7.8. Feature diversity ablation	20
7.9. Visual diversity ablation	20
7.10 Qualitative Results	21

6. IDEAL-M3D: Further Details

6.1. Problem Statement

6.1.1 Monocular 3D Object Detection (M3D).

Monocular 3D Detection (M3D) predicts categories and 3D bounding boxes \mathcal{B}_i for objects in an RGB image I with

intrinsics $K \in \mathbb{R}^{3 \times 4}$. Each \mathcal{B}_i is parameterized by position $(x_i, y_i, z_i) \in \mathbb{R}^3$, dimensions $(w_i, h_i, l_i) \in \mathbb{R}^3$, orientation $R_i \in SO(3)$, and class $c_i \in \mathbb{N}$. Given a dataset $\{(I_j, K_j, \mathcal{B}(I_j))\}_{j=1}^M$, with $\mathcal{B}(I_j)$ as ground-truth boxes for image I_j , the goal is to train a model capable of predicting $\mathcal{B}(I)$ for any given image I . Starting solely from a 2D RGB image poses a significant challenge due to depth ambiguity.

6.1.2 Active Learning.

Active Learning (AL) starts with a small labeled dataset $\mathcal{D}_L^0 = \{(I_i, K_i, \mathcal{B}(I_i))\}$ and a large unlabeled dataset $\mathcal{D}_U = \{(I_j, K_j)\}$. In each round r , we select $\mathcal{D}_r^* = \{(I_j, K_j, \mathcal{S}_{j,r}^*)\}_{j \in \mathcal{J}_r^*}$, where \mathcal{J}_r^* is the set of selected images and $\mathcal{S}_{j,r}^* \subseteq \hat{\mathcal{B}}(I_j)$ the subset of bounding boxes chosen for labeling. The oracle $\Omega : (I, \hat{\mathcal{B}}) \mapsto \mathcal{B} \cup \{\text{null}\}$, which in practice corresponds to a human annotator, refines each selected $\hat{\mathcal{B}}_{j,k} \in \mathcal{S}_{j,r}^*$, returning $\mathcal{B}_{j,k}$ or, in case it is not labeled, null. The labeled dataset is updated as $\mathcal{D}_L^r = \mathcal{D}_L^0 \cup \bigcup_{r'=1}^r \{(I_j, K_j, \{\mathcal{B}_{j,k}\}_{k \in \mathcal{S}_{j,r'}^*}) \mid \hat{\mathcal{B}}_{j,k} \neq \text{null}\}$, and the model is fine-tuned on \mathcal{D}_L^r . This repeats until the total labeled bounding boxes reach the budget $\mathcal{T} = \sum_{r=1}^R \sum_{j \in \mathcal{J}_r^*} |\mathcal{S}_{j,r}^*|$. In summary, AL pipelines can be seen as iterative cycles that repeat two steps: data selection for labeling, and training.

6.2. Loss Functions

We detail the loss computations for our baseline methods and omit the loss weights for clarity.

MonoLSS [23]:

$$L = L_{cls} + L_{c,o} + L_{h,w} + L_{S_{3d}} + L_\theta + L_{depth} \cdot \text{Sample } S \quad (12)$$

MonoCon [27]:

$$L = L_{cls} + L_{c,o} + L_{h,w} + L_{depth} + L_{S_{3d}} + L_\theta + L_{kp,h} + L_{kp,o} + L_{kp,co} \quad (13)$$

The individual losses can be categorized into image- and object-related losses.

6.2.1 Image-level Losses

The following losses are image-specific. Therefore, we apply the masking strategy as described in *cf.* Sec. 3.1 on both of these losses:

- L_{cls} : Gaussian kernel weighted focal loss for classification, following CenterNet [72].
- $L_{kp,h}$: (MonoCon only) Gaussian weighted focal loss for projected 3D keypoints as an auxiliary task.

6.2.2 Object-level Losses

These losses are specific to objects and do not require specialized masking:

- $L_{c,o}$: L1 Loss for offset from most confident foreground bin to precise projected 3D center
- $L_{h,w}$: L1 loss for 2D height and width
- $L_{S_{3d}}$: Dimension-loss. Dimension-aware L1 loss (L1 loss normalized by ground truth) in case of MonoCon and L1 loss in case MonoLSS
- L_{depth} : Laplacian aleatoric uncertainty loss
- L_θ : Multi-bin loss following Mousavian et al. [37]
- $L_{kp,o}$: L1 loss for keypoint offsets from keypoint heatmap
- $L_{kp,co}$: L1 loss for keypoint offsets from projected 3D center

6.3. Obtaining Diverse Features

Ensemble Features When extracting features for $Core\text{-}Set Box_{3D}$ we orientate on our baseline detectors. The idea is simple: We use the features that lead to a bounding box prediction. For MonoLSS [23] we use the region of interest (RoI) features of dimension $d \times 7 \times 7$, while for MonoCon [27] we use the features of size $d \times 3 \times 3$. Before applying Core-Set [47] selection, the tensors are flattened into a single dimensional vector.

For our main model we employ the standard DLA-34 [61] backbone ($d = 64$). For our auxiliary models we replace the DLA-34 with the backbones of RepVIT M [55] ($d = 56$) and MobileNetv4 M [41] ($d = 48$). These models are very lightweight in terms of parameters and still offer an acceptable 2D and 3D detection performance, while being easily interchangeable with minimal code modifications (*cf.* Fig. 12).

Visual Features The extraction of visual features follows a two step process. First we mask the image, then we encode it using an off-the-shelf image autoencoder.

To ensure compatibility with Core-Set selection, all object features must share the same dimensionality, yet the pixel space of instances differs. For example, a car closer to the camera has a larger pixel height than a car further away. While resizing objects to a fixed size is a straightforward solution, we adopt a more effective strategy (*cf.* Tab. 6). Specifically, we crop each object to a fixed height and width of 320×320 pixels, centering the crop on the 2D center of the object. If the object lies near the image boundary, we apply padding using the background color. We resize the cropped region to 128×128 pixels before feeding it into the

Table 6. **Ablation study on the KITTI [12] validation set showing the effect of resizing.** Results are reported using $NAURC_{60\%}$ $AP_{3D|R_{40}}$ for cars (IoU 0.7) and pedestrians/cyclists (IoU 0.5). **Final AP:** Moderate AP after training on 60% of the data.

Method	Easy	Moderate	Hard	Final AP
Ours w/ object resizing	21.51	15.17	12.55	18.59
IDEAL-M3D (Ours)	22.74	16.18	13.57	19.04

autoencoder. We utilize the autoencoder of Stable Diffusion v2-base [46] for this purpose.

Our approach ensures that the visual features effectively capture both the object’s 2D size and depth through the fixed cropping strategy, which preserves relative scale and encourages depth-diversity. Additionally, the resized crop maintains the object’s visual characteristics, enabling the autoencoder to learn a rich, low-dimensional representation of the visual appearance.

For segmentation [44], we utilize the SAMv2 ViT-B [44] architecture on Rope3D [60] and Waymo [51] and the larger SAMv2 ViT-L [44] variant on KITTI [12].

7. Experiments

7.1. NAURC Evaluation Metric

A central challenge in AL evaluation is to obtain a scalar, budget-aware summary that is comparable across selection paradigms. Previous work typically use two approaches: (i) training curves that plot performance against the number of labeled instances on the x-axis [13, 15, 33, 47, 58, 64, 67]; while informative, rankings often swap across budgets/time, complicating objective comparison; and (ii) fixed-budget snapshots that report performance at selected percentages of labeled data [16, 33, 35]; these yield a single value but depend on the chosen percentage and ignore the rest of the trajectory.

We introduce the Normalized Area under the Requested Curve (NAURC) to provide a fair, single-scalar comparison across image- and instance-based methods. NAURC adopts a common instance-based accounting for all approaches: for image-based selection, each selected image contributes the number of labelable instances it contains (the budget accumulates the total instances from the chosen images); for instance-based selection, the budget counts all requested instances, including requests that are later deemed non-labelable, thereby reflecting annotator verification effort. NAURC normalizes by a target budget and handles budget mismatch robustly: if a method overshoots the budget, we linearly interpolate back to it; if it undershoots, we keep the last observed performance (no extrapolation).

Empirical motivation for instance-based accounting. We analyze the relationship between actually labeled instances and the allocated budget in Fig. 5. For image-based methods, we quantify the budget at each AL iteration in terms of actually labeled boxes, since these methods request whole images. Most image-based strategies, with the exception of *Img Random*, tend to select images containing above-average numbers of instances, revealing that an image-count budget unfairly mixes annotation efforts because images vary widely in object count. For instance-based methods, the labeled ratio captures the proportion of true positive requests. Because false positive requests still require annotator verification, they count against the budget. This slightly penalizes instance-based methods compared to image-based methods. However, it also considers the annotators effort to verify that a requested object is a false positive.

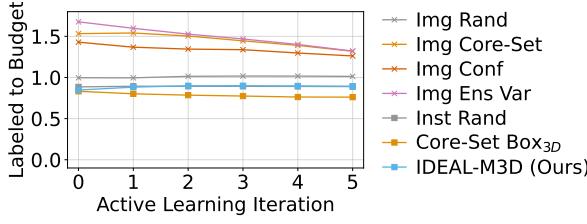


Figure 5. Ratio of labeled instances vs. instance-based budget on KITTI [12]. Most image-based methods request images with more than the average number of objects.

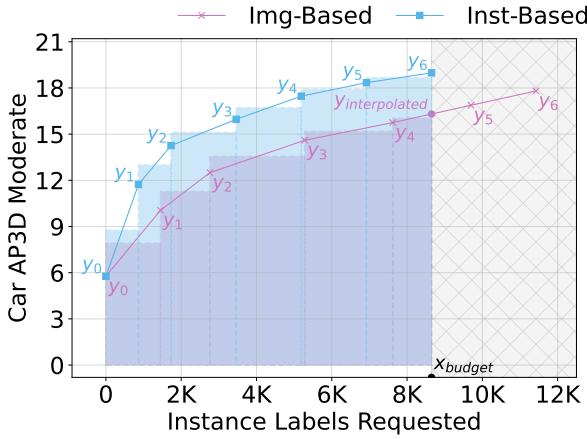


Figure 6. Visualization of the Normalized Area under the Requested Curve computation. For methods exceeding the label budget x_{budget} , the final performance metric is interpolated.

Formally, we compute the metric as the Area under the Requested Curve (AURC) up until the total requested label budget x_{budget} normalized by x_{budget} (refer to Fig. 6). The

NAURC is defined as:

$$\text{NAURC}_{x_{\text{budget}}} = \frac{1}{x_{\text{budget}}} \left(\text{AURC}_{\text{final}} + \sum_{i=0}^k \text{AURC}_i \right), \quad (14)$$

where k denotes the last AL iteration before reaching the requested instance label budget x_{budget} :

$$k = \max \{i \mid x_{i+1} \leq x_{\text{budget}}\}. \quad (15)$$

We compute the metric across iterations using the trapezoidal rule to calculate the AURC between consecutive data points.

$$\text{AURC}_i = \frac{(y_{i+1} + y_i)}{2} \cdot (x_{i+1} - x_i), \quad (16)$$

where y_i and y_{i+1} are the metric values at the i -th and $(i+1)$ -th AL iteration and x_i and x_{i+1} are their respective requested instance labels. The AURC of the final interval is computed as:

$$\text{AURC}_{\text{final}} = \frac{(y_{\text{interpolated}} + y_k)}{2} \cdot (x_{\text{budget}} - x_k), \quad (17)$$

where $y_{\text{interpolated}}$ is the metric value at the budget point. To handle both cases where methods exceed or fall short of the target budget, we define:

$$y_{\text{interpolated}} = \begin{cases} y_k + \Delta y \cdot \frac{x_{\text{budget}} - x_k}{x_{k+1} - x_k} & \text{if } x_{k+1} > x_{\text{budget}} \\ y_k & \text{otherwise} \end{cases} \quad (18)$$

where $\Delta y = y_{k+1} - y_k$.

Compared to prior AUC-style metrics [65–67], NAURC enables fair cross-paradigm comparison across image- and instance-based methods: If a method overshoots the budget, NAURC linearly interpolates the terminal performance at x_{budget} ; if it undershoots, it holds the last observed value. This removes extrapolation and prevents artificial inflation.

7.2. Implementational Details

Active Learning Parameters. The distance weights (*cf.* Eq. (5)) are determined as $\lambda_1 = \lambda_2 = \lambda_3 = \frac{1}{6}$, with λ_{vis} set to $\frac{1}{2}$ to balance the contribution of visibility metrics. The Loss weight multiplier parameter is set to $\delta = 0.2$. For adaptive sampling, we configure the time-dependent decay parameter α to 3.0 for KITTI [12] and 30.0 for Rope3D [60] and Waymo [51] to account for the different dataset size. On KITTI [12], labeling proposals exclude objects with heights below 25 pixels, as ground truth instances are constrained to a minimum height of 25 pixels, which eliminates unsuitable candidates during this process. All experiments, including run time evaluations, are conducted on a single NVIDIA A40 GPU with 32GB RAM. For instance-based AL, we mitigate redundant labeling by skipping requests for objects that fall within 95% of the radius (r_x, r_y) from a previous request

and share the same predicted class. This prevents repeated annotations or requests regarding the same false-positive predictions. For masking and visual feature extraction, we leverage the 2D bounding box predictions of the main model. To accelerate computation, we apply principal component analysis (PCA) to the extracted features, compressing the dimensions while retaining at least 99% of the variance.

MonoLSS [23] Configuration. The training setup utilizes a batch size of 16 and optimizes the model using the Adam [20] optimizer with a weight decay of 1e-5. Starting from a pre-trained checkpoint, each AL iteration spans 150 epochs for the main model. During training, we initialize the learning rate at 1e-3 and decay it by a factor of 0.1 after 60% and 80% of the epochs. Additionally, the first cycle incorporates a 5-epoch cosine warmup to stabilize gradient updates. After the initial phase of training, the LSS module gets activated at the 50th epoch. To enhance model robustness, we apply comprehensive data augmentation techniques, including random horizontal flipping, shifting (W: ± 256 pixels, H: ± 77 pixels), scaling (0.6-1.4), and MixUp3D using fully labeled images.

MonoCon [27] configuration. For MonoCon, we adopt a batch size of 24 and train the model using the AdamW [29] optimizer with a weight decay of 1e-5 and a learning rate of 0.0011. The initial training phase encompasses 90 epochs, followed by an additional 30 epochs for every subsequent AL cycle. To maintain stable training, we implement gradient clipping with a norm of 35 and apply cosine learning rate scheduling. For computational efficiency, images are resized to 960×640 pixels, in line with the settings from [36]. The augmentation pipeline integrates a rich variety of transformations, including photometric distortion, random shift (± 32 pixels), horizontal flipping, and random cropping (900×550 pixels). Furthermore, for Rope3D [60], we learn the $SO(3)$ orientation matrix in alignment with the GroundMix [36] approach.

Labeling radius. Also, to simplify the task for the labeler, we expect that an object lies within a depth-dependent radius (r_x, r_y), letting the user focus on a small area. Such radius is defined as:

$$r_x = H \cdot \frac{f_x}{\hat{z}}, \quad r_y = H \cdot \frac{f_y}{\hat{z}}, \quad (19)$$

where H is a scaling factor, (f_x, f_y) are the camera focal lengths, and \hat{z} is the predicted depth. This ensures accurate targeting by accounting for the geometric effects of depth, as pixel-space errors decrease with distance. We define the labeling radius via $H = 2.0$, which corresponds to approximately 47 pixels for objects at a 30m distance on KITTI [12].

Instance matching between ensemble members. Using the predictions from the main model as a reference, we associate predictions from auxiliary models by calculating the 2D bounding box IoU. This computation uses a relaxed threshold of 0.5 to accommodate additional object matches. As auxiliary models are trained with fewer resources, we adjust detection thresholds to compensate for their lower capacity. Specifically, thresholds are reduced from 0.2 to 0.1 for MonoCon [27] and from 0.2 to 0.05 for MonoLSS [23]. For cases of multiple associations, we prioritize the object with the highest confidence score. This strategy significantly enhances instance coverage, particularly for rare and low-confidence objects, while maintaining diversity and minimizing discard rates.

7.3. Implementation Details of Baseline Methods

For completeness, we provide additional implementation details for a subset of our baseline methods. In our evaluation framework, we adapt image-based methods by selecting images based on their highest-scoring contained instance, while instance-based methods directly employ our labeling pipeline (Sec. 3.1) with their respective acquisition functions. To manage computational resources effectively, we implement dataset-specific sampling strategies: For datasets containing on average more than 10 objects per image (*e.g.* Rope3D [60]), we limit the maximum number of requested images to one-third of the number of instances, though this restriction is not applied to KITTI [12] due to its limited size. In the following, we provide further details on the individual baseline methods.

Augmentation Depth Variance (Augm Depth Var): This method provides a computationally efficient alternative to ensemble approaches by employing multiple augmented forward passes of a single model. The augmentation pipeline includes blur ($\sigma = 0.5$), brightness adjustments (factors: 0.5-1.5, probability: 0.2-0.5), and hue shifts (± 0.1). We exclude this method for MonoCon [27] evaluations due to its built-in color augmentation during training.

CDAL [1]: This approach implements Core-Set sampling by computing class-wise confusions from softmax probabilities of both labeled data heatmap indices and predictions. The sampling strategy employs pairwise and class-wise KL divergence for measuring image similarity.

CloseDepth: Instances with low depth values are preferred for labeling.

Confidence (Conf): This uncertainty-based approach combines depth and classification confidence scores from the detector, selecting instances with lower confidence for labeling prioritization.

DDFH [4]: DDFH is the current state-of-the-art method for active learning in LiDAR-based 3D detection. The approach optimizes three key aspects: balancing class distribution, maximizing frame-level heterogeneity, and selecting

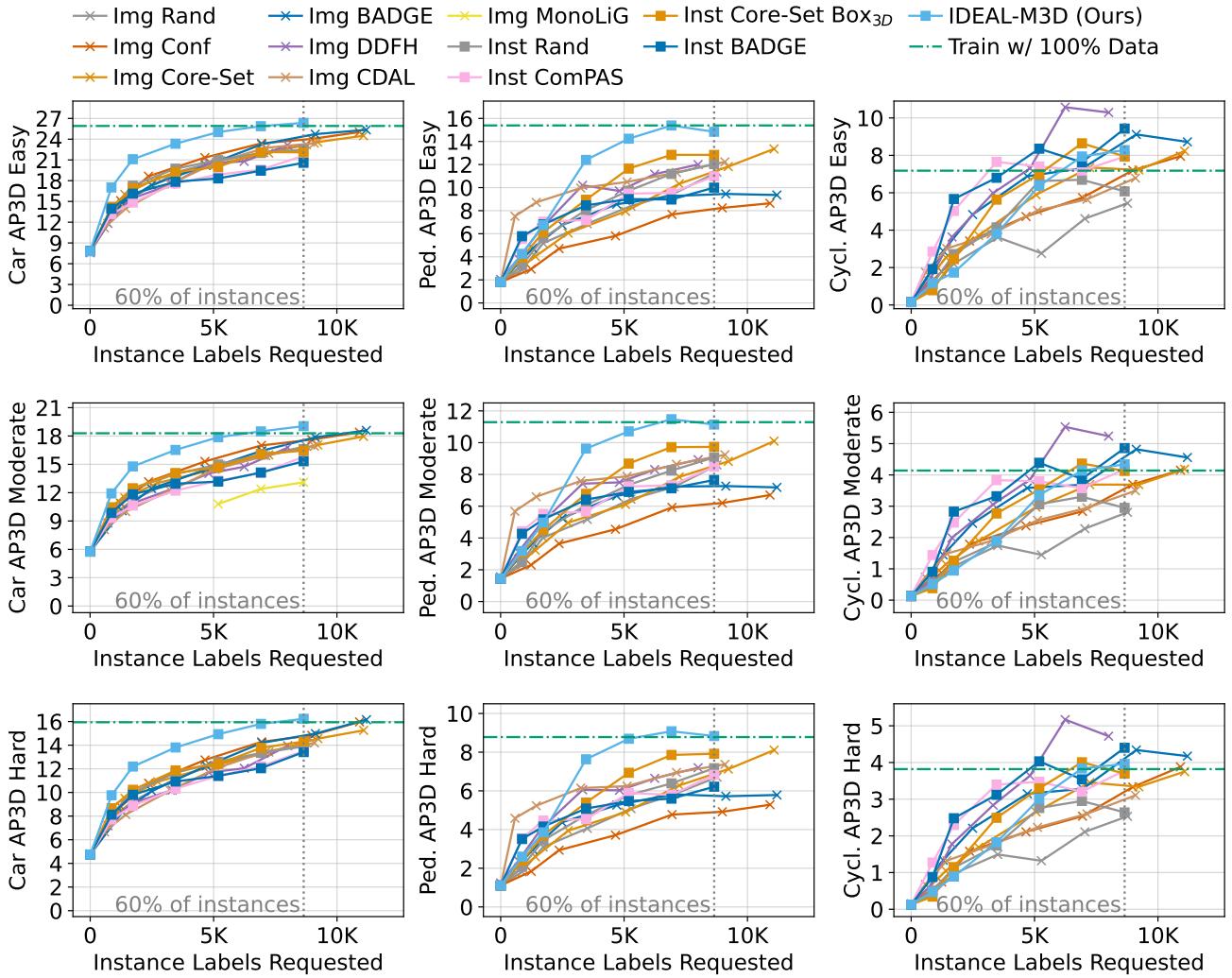


Figure 7. AL methods evaluated on the KITTI validation set [12] for Easy, Moderate and Hard on Car (IoU=0.7), Pedestrian (IoU=0.5) and Cyclist (IoU=0.5).

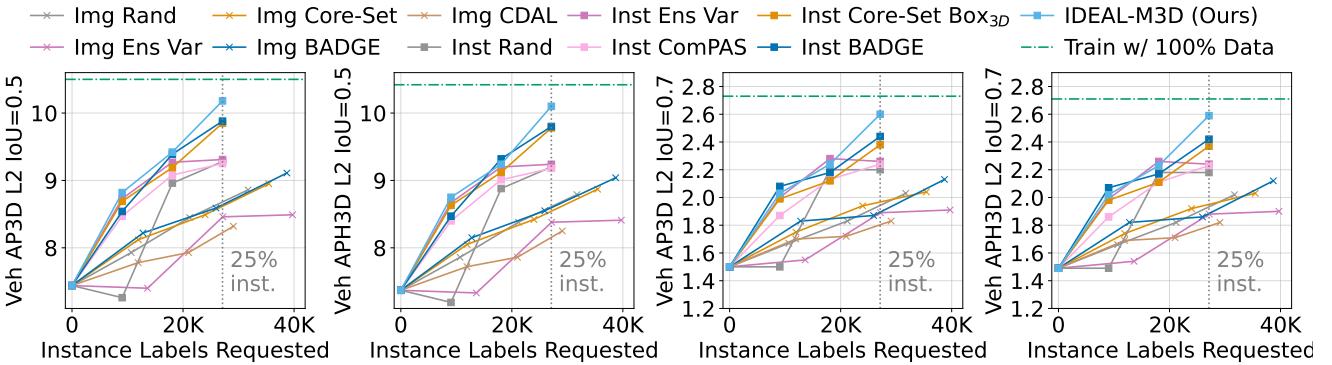


Figure 8. AL methods evaluated on the Waymo validation set [51].

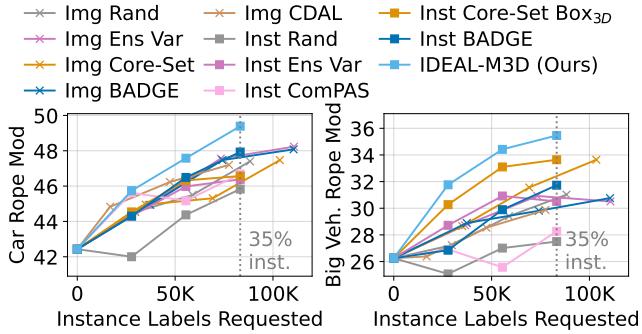


Figure 9. AL methods evaluated on the Rope3D validation set [60] with Rope Score at IoU 0.5.

diverse instances. For instance diversity, it fits a Gaussian mixture model to t-SNE [53] compressed model features and additional geometric features. We adapt this approach to monocular 3D detection by dropping the point density information from the feature vectors.

BADGE [3]: BADGE computes gradient embeddings derived from the penultimate classification layer for each data point and selects samples via the k-MEANS++ seeding algorithm [2]. We compute gradients based on the focal loss for the predicted heatmap locations on an instance basis.

ComPAS [33]: We adapted the active learning component of ComPAS [citation] to our setting, which selects instances based on their localization and classification disagreement between a chairman and its committee members. This disagreement is measured using multiple data augmentations, including scale, shift, contrast, solarize, saturation, sharpness, and brightness.

Dropout-based Methods: We investigated dropout-based uncertainty estimation but excluded it from our final evaluation due to significant performance degradation. Specifically, when applying dropout to backbone features, we observed substantial drops in accuracy: While the baseline achieves 18.29 for Car $AP_{3D|R_{40}}^{IoU=0.7}$ Moderate, introducing dropout rates of 1%, 5%, and 10% reduces performance to 17.73, 15.81, and 13.67 on the KITTI [12] validation set, respectively (baseline: MonoLSS [23]).

Efficient AL [16]: We excluded Efficient AL [16] from our comparison due to ambiguous evaluation metrics (unspecified IoU thresholds for mAP and no differentiation between KITTI [12]'s easy, moderate, and hard instances).

Ensemble Depth Variance (Ens Depth Var): Using an ensemble of three models, we associate predictions via a 2D IoU threshold of 0.5. The uncertainty metric is derived from the variance of depth predictions relative to their mean, prioritizing instances exhibiting higher variance.

Ensemble Relative Standard Deviation (Ens Rel Std): Similar to *Ens Depth Var*, but normalizes the depth standard deviation by the mean predicted depth.

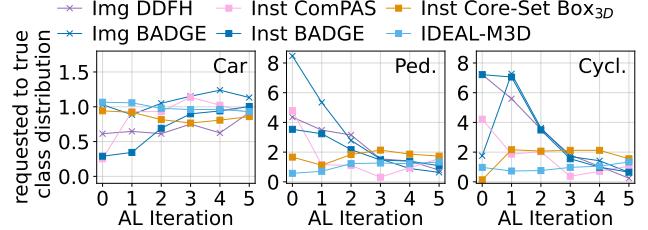


Figure 10. KITTI [12] ratio of requested AL instances to the true dataset distribution (<1: Undersampling, 1: Balanced, >1: Oversampling).

FarDepth: Instances with large depth values are preferred for labeling. Predicted instances need to have a 2D height of at least 25 pixels and a depth smaller 50m.

MonoLiG [15]: MonoLiG [15] combines AL with semi-supervised learning by leveraging an ensemble of five models alongside a LiDAR-based teacher model. The instance selection strategy integrates three uncertainty measures: the aleatoric uncertainty of the teacher model, the disagreement among ensemble members, and the teacher-student prediction inconsistency. These measures are aggregated into a unified scoring mechanism for ranking potential labeling candidates. Due to unavailability of the official implementation at submission time, we report the results from the paper. To establish a fair comparison, we estimate the number of labeled instances per AL iteration based on the dataset's average object count per image.

7.4. Comparison with AL methods

We provide additional comparisons to contextualize the main results in Tab. 2. Final accuracies at the end of AL training are summarized in Tab. 8, and the corresponding training curves are shown in Figs. 7 to 9, offering both endpoint and trajectory views of performance.

Per-class results on KITTI for pedestrian and cyclist are reported in Fig. 7 and Tab. 7. IDEAL-M3D significantly outperforms all baselines on Car and Pedestrian. While ComPAS [33], Core-Set Box_{3D} BADGE [3], and DDFH [4] report higher AP on Cyclist, this comes from a highly skewed budget allocation: they assign up to 7× more annotations to Cyclist at the expense of the other classes (*cf.* Fig. 10). In contrast, IDEAL-M3D maintains a balanced acquisition across categories, which we hypothesize emerges from the representational diversity of our ensembles. This yields stronger average performance and a more uniform gain among all classes.

7.5. Uncertainty vs. Diversity-based Methods

To understand why uncertainty-based methods are less effective for instance-based M3D (*cf.* Tab. 2), we conduct a detailed analysis comparing *Core-Set Box_{3D}* with three

Table 7. **AL performance on the KITTI [12] validation dataset.** Results are averaged over three rounds, each initialized from the same checkpoint. **KITTI:** We report NAURC_{60%} $AP_{3D|R_{40}}$ with IoU thresholds of 0.7 (cars) and 0.5 (pedestrians, bicyclists). **Type***: U=Uncertainty-based, D=Diversity-based, H=Hybrid.

		Car			Pedestrian			Cyclist			Average			
	Method	Type*	Easy	Mod.	Hard	Easy	Mod.	Hard	Easy	Mod.	Hard	Easy	Mod.	Hard
Image-based	Rand	-	18.01	12.97	10.77	6.31	4.83	3.87	2.88	1.44	1.29	9.07	6.41	5.31
	Conf	U	19.72	14.20	11.80	5.56	4.31	3.47	4.19	2.11	1.88	9.82	6.87	5.72
	Ens Depth Var	U	18.19	13.00	10.77	4.36	3.47	2.74	4.15	2.09	1.89	8.90	6.19	5.13
	Augm Depth Var	U	18.24	13.16	10.90	4.53	3.61	2.86	3.34	1.67	1.49	8.70	6.15	5.08
	Core-Set [47]	D	18.80	13.50	11.34	7.20	5.59	4.47	4.88	2.45	2.23	10.29	7.18	6.01
	BADGE [3]	H	19.02	13.53	11.39	7.40	5.76	4.62	5.68	2.91	2.60	10.70	7.40	6.20
	DDFH [4]	H	17.94	12.72	10.51	8.90	6.68	5.38	6.47	3.37	3.12	11.10	7.59	6.34
	CDAL [1]	H	18.38	12.99	10.82	9.88	7.49	6.00	4.27	2.16	1.93	10.84	7.55	6.25
Instance-based	Rand	-	18.69	13.53	11.21	8.03	6.05	4.72	4.49	2.10	1.91	10.40	7.23	5.95
	Conf	U	12.10	8.93	7.60	3.71	2.90	2.30	2.96	1.42	1.27	6.26	4.22	3.87
	Ens Depth Var	U	11.88	9.07	7.82	3.66	2.82	2.22	3.67	1.75	1.59	6.40	4.55	3.88
	Augm Depth Var	U	12.25	8.96	7.64	3.15	2.42	1.94	2.80	1.36	1.24	6.07	4.25	3.61
	ComPAS [33]	U	17.31	12.31	10.48	8.04	6.24	5.01	6.29	3.18	2.87	10.55	7.24	6.12
	Core-Set [47] Box _{3D}	D	18.83	13.82	11.68	9.38	7.07	5.68	5.50	2.78	2.54	11.24	7.89	6.63
	BADGE [3]	H	17.39	12.63	10.74	7.98	6.13	4.88	6.55	3.32	3.04	10.64	7.36	6.22
	IDEAL-M3D (Ours)	D	22.74	16.18	13.57	11.42	8.61	6.86	4.84	2.51	2.32	13.00	9.10	7.58

Table 8. **Final AL performance on KITTI [12] validation, Waymo [45, 51] validation and Rope3D [60] validation dataset.** Results are averaged over three rounds, each initialized from the same checkpoint. We report the final accuracy after training on 60% of data (KITTI) and 25% of data (Rope3D, Waymo). For image-based methods exceeding the budget we report the interpolated result. **Type***: U=Uncertainty-based, D=Diversity-based, H=Hybrid.

		KITTI [12] Car			Waymo [51] Vehicle				Rope3D [60]				
		$AP_{3D R_{40}}^{0.7}$			IoU=0.5		IoU=0.7		Car		Big Vehicle		
	Method	Type*	Easy	Mod.	Hard	AP _{3D}	APH _{3D}	AP _{3D}	APH _{3D}	AP _{3D} ^{0.5}	Rope	AP _{3D} ^{0.5}	Rope
Image-based	Rand	-	22.47	16.42	13.80	8.68	8.20	1.74	1.73	35.13	47.08	16.63	30.74
	Conf	U	23.89	17.51	14.72	8.12	8.05	1.66	1.65	34.04	46.24	17.61	31.63
	Ens Depth Var	U	22.81	16.28	13.66	8.46	8.37	8.37	1.87	35.93	47.74	16.76	30.81
	Augm Depth Var	U	22.20	16.15	13.72	-	-	-	-	-	-	-	-
	Core-Set [47]	D	23.10	16.73	14.19	8.32	8.25	1.90	1.88	34.11	46.25	18.59	32.41
	BADGE [3]	H	<u>24.41</u>	<u>17.55</u>	<u>14.81</u>	8.35	8.33	1.74	1.74	35.81	47.62	15.77	30.05
	DDFH [4]	H	23.13	16.77	13.96	8.22	8.15	1.85	1.84	35.70	47.54	19.02	32.71
	CDAL [1]	H	23.30	16.69	14.04	7.88	7.81	1.70	1.69	35.25	47.20	15.43	29.87
Instance-based	Rand	-	22.50	16.61	13.89	9.28	9.20	2.20	2.18	33.53	45.83	12.31	27.50
	Conf	U	13.00	10.00	8.67	9.05	8.97	2.32	2.31	34.85	46.86	14.38	29.05
	Ens Depth Var	U	11.29	9.01	8.43	9.31	9.24	2.26	2.24	34.28	46.38	16.68	30.95
	Augm Depth Var	U	9.69	7.38	6.58	-	-	-	-	-	-	-	-
	ComPAS [33]	U	21.59	15.79	13.52	9.25	9.18	2.24	2.23	34.58	46.68	13.32	28.27
	Core-Set [47] Box _{3D}	D	22.10	16.45	14.31	9.85	9.78	2.38	2.37	34.51	46.57	<u>20.08</u>	<u>33.64</u>
	BADGE [3]	H	20.59	15.33	13.42	<u>9.88</u>	<u>9.80</u>	<u>2.44</u>	<u>2.42</u>	<u>36.13</u>	<u>47.93</u>	17.50	31.73
	IDEAL-M3D (Ours)	D	26.35	19.04	16.23	10.18	10.10	2.60	2.59	37.94	49.39	22.21	35.46

uncertainty-based methods:

- *Conf*: Selects instances with lowest model confidence scores (aleatoric uncertainty)
- *Ens Rel Std*: Targets instances with highest relative

depth deviation from the ensemble mean (three models)

- *Ens Depth Var*: Prioritizes instances with highest absolute depth error relative to the ensemble mean

Our analysis of the selected instances in Fig. 11 reveals

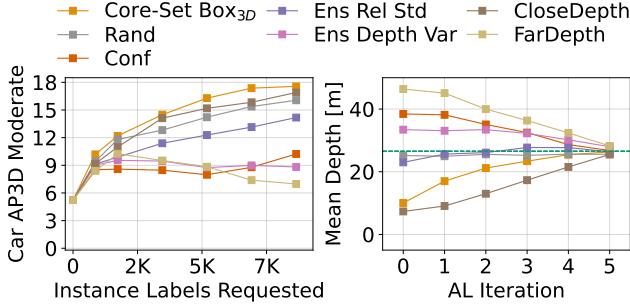


Figure 11. Comparison of instance-based active learning methods for monocular 3D detection. Uncertainty-based methods tend to select distant objects, which provide weaker training signals, while diversity-based approaches favor closer objects that are more informative. **Left:** The effectiveness of various selection strategies. **Right:** Distribution of selected instances, showing diversity-based methods’ preference for close objects compared to uncertainty-based methods’ bias toward moderately far and distant objects.

distinct selection patterns: Both *Conf* and *Ens Depth Var* demonstrate a clear bias towards distant objects, which is expected given that absolute depth errors typically increase with distance. In contrast, *Ens Rel Std* achieves more balanced sampling across different depths, while *Core-Set Box_{3D}* shows a preference for closer objects. We hypothesize that closer objects generally provide richer visual information due to higher pixel density and more distinct features, making them more conducive to accurate depth learning. Interestingly, diversity-based methods also tend to perform better for far away objects (hard category, *cf.* Tab. 2), even though they are undersampled during training time. To further investigate this, we introduce two simple AL strategies:

- *CloseDepth*: Prioritizes close instances for labeling
- *FarDepth*: Prioritizes most distant instances (filtered to instances ≤ 25 pixels and ≤ 50 m depth to avoid false positives)

Our experiments in Fig. 11 show that *CloseDepth* performs second-best after *Core-Set Box_{3D}*, leading to two key insights. First, closer instances provide more effective training signals, even for detecting distant objects. Second, uncertainty-based methods are suboptimal for instance-based M3D due to their inherent bias towards moderate far and distant objects

7.6. Training time comparison with fully supervised methods

We compare total training time in Tab. 9 against fully supervised and semi-supervised baselines. Relative to MonoLSS [23], IDEAL-M3D trains for roughly twice as long. It reaches (near) fully supervised accuracy using only 60%

of the labeled data on the validation and test sets. Since compute is typically far cheaper than human annotation, this is a favorable trade-off.

Compared to MonoLiG [15], IDEAL-M3D achieves higher accuracy on KITTI validation and test (*cf.* Tab. 3 and Fig. 3). It also reduces training time by about 3 \times . This efficiency stems from concrete design choices. MonoLiG trains an ensemble of five image-based student models and an additional sixth LiDAR-based model. In contrast, we train only three models. Our auxiliary ensemble components are lightweight and fast to train.

We also compare with the semi-supervised Mix-Teaching [59]. Its training time is more than 4 \times higher. The method trains a five-model ensemble across one supervised and three semi-supervised rounds, which introduces substantial overhead. In contrast, IDEAL-M3D attains higher KITTI test accuracy without using unlabeled data and with fewer labeled samples (*cf.* Tab. 3).

Table 9. Comparison of training time of selected methods on the KITTI [12] trainval set. **SSL** denotes that the method uses semi-supervised learning. **AL** denotes that the method uses active learning.

Method	SSL	AL	Total training time
MonolSS [23] (Baseline)	✗	✗	37h
IDEAL-M3D (Ours)	✗	✓	75h
MonoLiG [15]	✓	✓	240h
Mix-Teaching [59]	✓	✗	305h

7.7. Backbone ablation

To promote ensemble diversity while keeping compute modest, we equip the auxiliary models with distinct lightweight backbones (*cf.* Sec. 3.3). We evaluate several candidates on KITTI [12] validation set using 100% of the training data to isolate backbone effects. Throughput and accuracy are summarized in Fig. 12.

For the main model, we retain DLA-34 [61] due to its strong accuracy, ensuring a fair and comparable reference across experiments. This backbone remains fixed in all primary results.

For the auxiliary models, we select RepViT-M1.0 [55] and MobileNetV4-Conv-M [41]. MobileNetV4-Conv-M offers the best speed–accuracy trade-off among the tested lightweight backbones, being both faster and more accurate than alternatives. ConvNeXt-Pico [28] is marginally faster than RepViT-M1.0 but is approximately 25% less accurate; we therefore prefer RepViT-M1.0. Together, RepViT-M1.0 and MobileNetV4-Conv-M provide complementary inductive biases and increase architectural diversity at low training cost.

Table 10. Ensemble ablation study on the KITTI [12] validation set. Results are reported using NAURC_{60%} $AP_{3D|R_{40}}$ for cars (IoU=0.7) and pedestrians/cyclists (IoU=0.5). **Final AP:** Moderate AP after training on 60% of the data.

Method	Easy	Moderate	Hard	Final AP
Ours w/o diverse backbones	20.85	15.08	12.66	18.10
Ours w/o data sampling	21.21	15.33	12.61	18.78
Ours w/ full epochs	22.38	16.03	13.42	18.42
Ours w/o random loss	22.02	15.55	12.55	18.93
IDEAL-M3D (Ours)	22.74	16.18	13.57	19.04

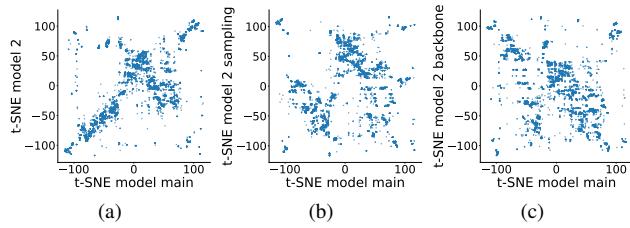


Figure 13. t-SNE [53] visualization of RoI. Considered features from ensemble models trained on 30% of the KITTI [12] dataset. **(a)** Identical data and backbones yield correlated features, limiting diversity. **(b)-(c)** Adaptive sampling and varied backbones (RepViT [55]) enhance feature diversity, improving ensemble effectiveness.

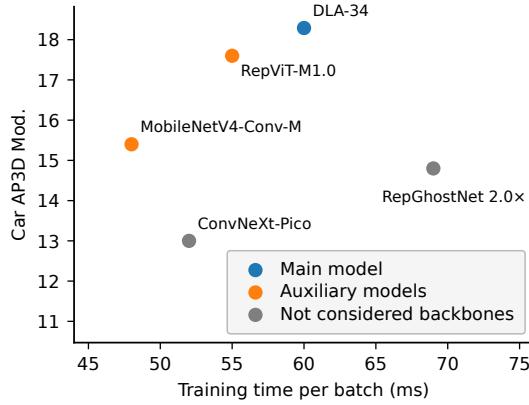


Figure 12. Comparing KITTI [12] validation set performance under diverse backbones (supervised raining with 100% of data).

7.8. Feature diversity ablation

A key component of our approach is instance selection driven by diverse, fast-to-train feature ensembles. We increase diversity and reduce compute through four complementary mechanisms: (i) heterogeneous lightweight backbones, (ii) time-adaptive data sampling that varies the training trajectory across ensemble members, (iii) fewer training epochs for auxiliary models, and (iv) random sampling of loss weights to perturb optimization and feature emphasis

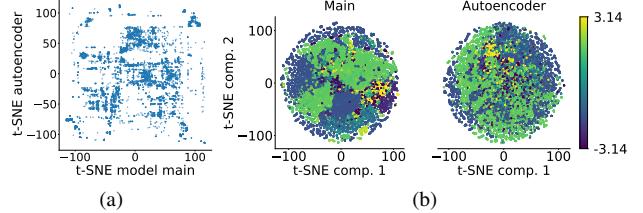


Figure 14. t-SNE [53] visualization of features trained on 30% of the KITTI [12] dataset, highlighting the complementarity of RoI and autoencoder features. **(a)** RoI and autoencoder features show low correlation, capturing distinct information. **(b)** Example: RoI features are dominated by the object orientation around the *y*-axis, while Stable Diffusion v2-base [46] autoencoder features are more orientation-invariant. The vehicle’s orientation is shown by color.

Table 11. Visual diversity ablation study on the KITTI [12] validation set. Results are reported using NAURC_{60%} $AP_{3D|R_{40}}$ for cars (IoU=0.7) and pedestrians/cyclists (IoU=0.5). **Final AP:** Moderate AP after training on 60% of the data.

Method	Easy	Moderate	Hard	Final AP
Ours w/o SAMv2 [44]	17.72	13.24	11.11	16.36
Ours w/ DINOv2 [38]	17.87	13.08	11.12	16.60
IDEAL-M3D (Ours)	22.74	16.18	13.57	19.04

across tasks/heads.

Together, these mechanisms reduce total training time by approximately 15 hours compared to a vanilla, homogeneous ensemble trained for full schedules, while improving selection quality through more diverse feature views.

The ablation in Tab. 10 supports this design. Reducing epochs on auxiliary models preserves performance to within noise levels, indicating that full schedules are unnecessary for effective feature-based selection. Adding backbone diversity, time-adaptive sampling, and random loss-weight sampling yields incremental gains in $AP_{3D|R_{40}}$ (Moderate) of +1.10, +0.85, and +0.63, respectively.

We further analyze feature diversity in Fig. 13 via t-SNE [53]. Ensembles using time-adaptive sampling and heterogeneous backbones exhibit markedly lower inter-model feature correlation than a vanilla ensemble, confirming that our mechanisms produce complementary representations that enhance the quality of selected instances.

7.9. Visual diversity ablation

In Tab. 11, we ablate components that promote visual diversity in the feature ensemble.

Removing Segment Anything Model 2 (SAMv2) reduces AP_{3D} (Moderate) by more than 2 points. We hypothesize two causes. First, without SAMv2, background content leaks into the features. The same object with different back-

grounds then appears artificially diverse, which dilutes foreground cues. Second, the foreground mask carries coarse geometric information. E.g, the silhouette of a car at 45° yaw differs from one at 0° , which is useful for 3D selection.

Replacing the autoencoder with DINOv2 [38] features causes a similar drop. DINOv2 emphasizes global scene semantics. It captures less local, instance-centric detail [54]. The autoencoder yields object-focused representations that better support instance selection.

Fig. 14 shows t-SNE [53] visualizations of autoencoder features and detector features. The two spaces exhibit minimal correlation. This indicates that the task-agnostic autoencoder provides complementary signals to the task-specific model features.

7.10. Qualitative Results

In the subsequent pages (*cf.* Fig. 15, Fig. 16, Fig. 17) we present further qualitative results on IDEAL-M3D highlighting the prediction accuracy and selection process over time.

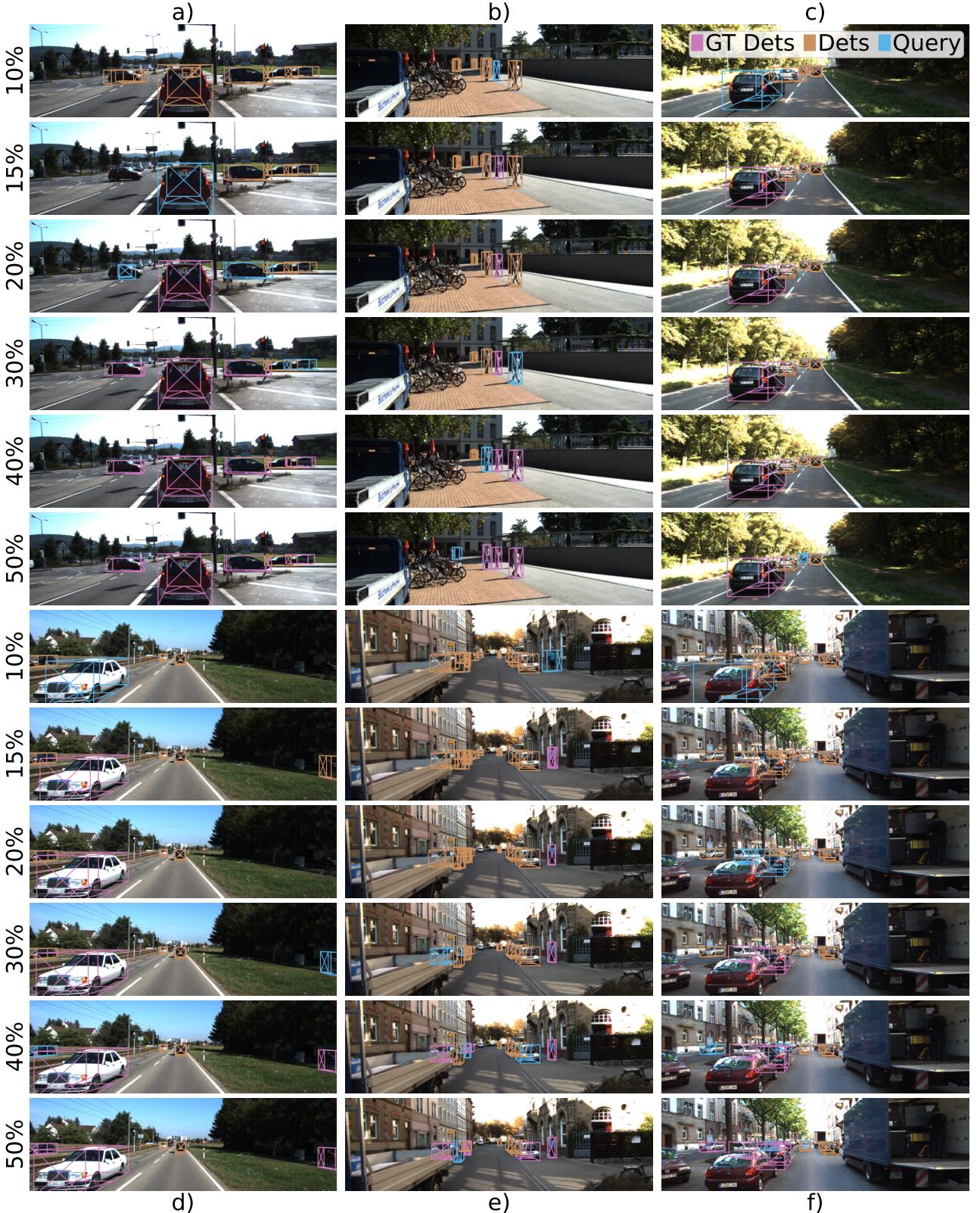


Figure 15. Qualitative results of IDEAL-M3D on the KITTI [12] dataset showing prediction evolution over time (top to bottom). Pink boxes represent previously labeled instances now in the training set, cyan boxes indicate predictions selected for current labeling, and orange boxes show predictions not chosen for annotation. The progression demonstrates the model’s improvement through strategic label acquisition



Figure 16. Qualitative results of IDEAL-M3D on the Waymo [51] dataset showing prediction evolution over time (top to bottom). Pink boxes represent previously labeled instances now in the training set, cyan boxes indicate predictions selected for current labeling, and orange boxes show predictions not chosen for annotation. The progression demonstrates the model’s improvement through strategic label acquisition

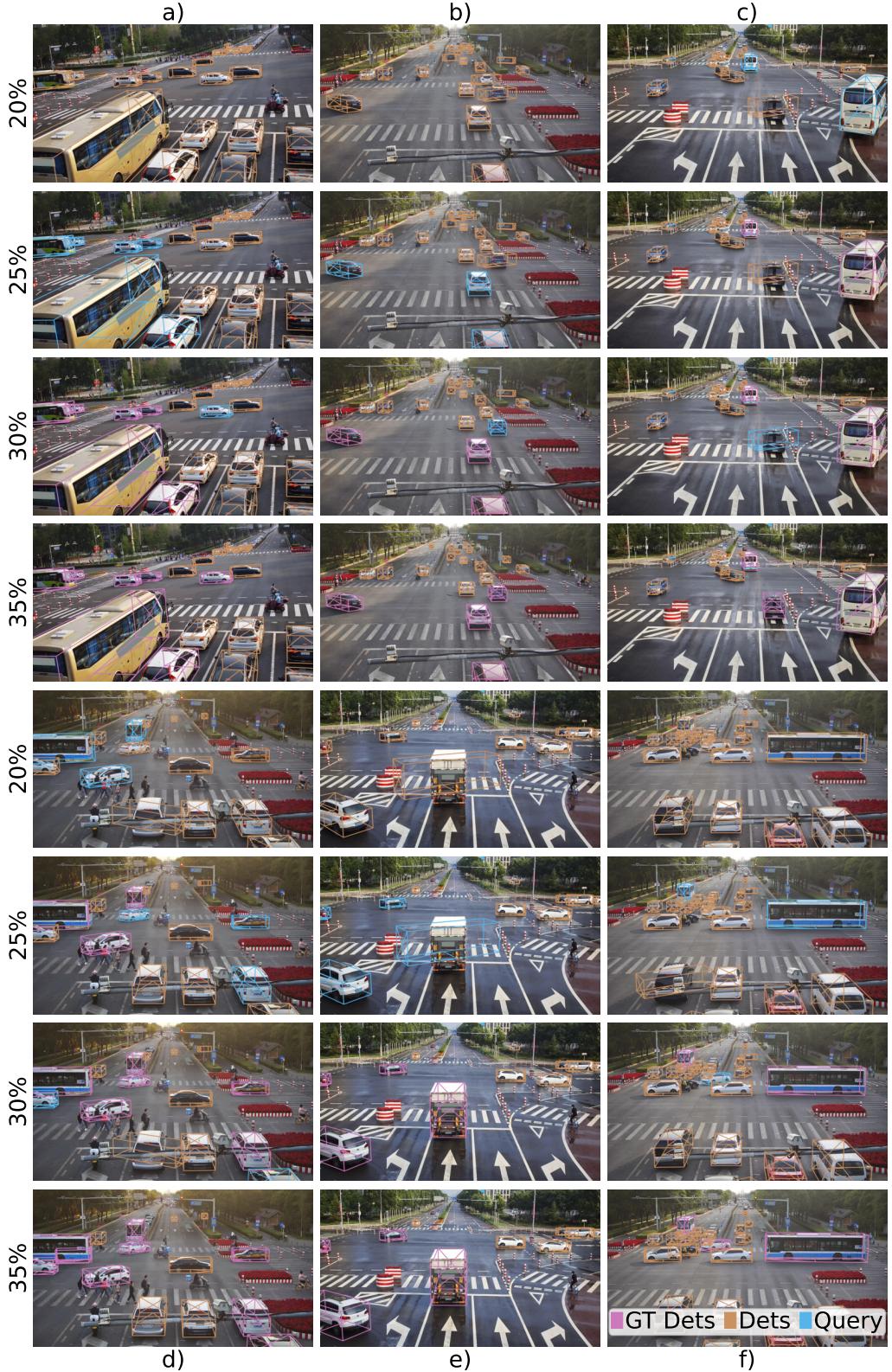


Figure 17. Qualitative results of IDEAL-M3D on the Rope3D [60] dataset showing prediction evolution over time (top to bottom). Pink boxes represent previously labeled instances now in the training set, cyan boxes indicate predictions selected for current labeling, and orange boxes show predictions not chosen for annotation. The progression demonstrates the model’s improvement through strategic label acquisition