universitätfreiburg





Label-Efficient LiDAR Scene Understanding with 2D-3D Vision Transformer Adapters

Julia Hindel*, Rohit Mohan*, Jelena Bratulic, Daniele Cattaneo, Thomas Brox, and Abhinav Valada * Equal Contribution

Motivation

Foundation models have advanced vision and language models, but LiDAR still lacks large-scale pretraining and robust foundation models.



Results

Benchmark

BALVIT consistently outperforms baseline models on the SemanticKITTI and nuScenes datasets under low-label settings.



Existing methods rely on paired camera-LiDAR data and still require training 3D encoders from scratch to leverage foundation models.

We propose BALVIT, a universal foundation model for LiDAR segmentation that leverages a **novel 2D-3D adapter** to enable efficient reuse of vision model features through structured 2D representations.

Method

BALViT Architecture

Dual-View Encoding: LiDAR point clouds are encoded using

Method	Seman 0.1%	ticKITTI 1%	nuSo 0.1%	cenes 1%	
SR-Unet18 FRNet SphereFormer RangeViT	- 30.09 29.21 28.74	39.50 40.78 42.81 43.53	- 28.03 30.42 27.79	30.30 48.98 50.06 52.88	Fully Supervis Vision Distilla
SLidR ST-SLidR SEAL CLIP2Scene	-	44.60 44.72 46.63 42.60	-	38.30 40.75 45.84 56.30	Parameter Efficient Fine-Tuning
Frozen ViT backbone Bias tuning LoRA VPT ViT Adapter BALVIT (Ours)	29.97 30.86 31.65 31.07 29.55 32.85	45.91 45.63 46.27 46.08 45.01 51.80	28.72 28.15 28.27 29.68 27.50 31.86	54.70 56.05 57.57 55.67 56.06 59.27	

Effects of Vision Backbone Initialization

54 1

Cityscapes pretraining performs (1% SemanticKITTI), highlighting the value of domain-aligned vision backbones.

- Range View (RV) and Bird's-Eye View (BEV) branches.
- Vision Backbone on RV: A frozen vision transformer (ViT) backbone processes RV features with a learnable patch embedding, leveraging rich pre-trained visual representations.
- 2D-3D Adapter: Enables bidirectional injection of BEV and ViT-processed RV features via stacked parallel cross-attention for mutual refinement.
- 3D Positional Embeddings: Provide spatial alignment between RV and BEV, enabling seamless cross-view attention in the adapter.
- Parallel Decoding: RV and BEV use independent decoders to predict semantic labels.





Ablation Study on BALVIT Components

The 2D-3D adapter and BEV decoder contribute most to gains over the frozen ViT baseline (1% SemanticKITTI).





Inference Fusion

- The output is selected based on the highest logit confidence \hat{y} from RV and BEV predictions, using a fixed threshold s.
- This simple uncertainty-aware selection favors fine-grained RV when confident and falls back to BEV for globally consistent context.

$$Output = \begin{cases} \hat{y}_{\text{RV}}, & \text{if } \hat{y}_{\text{RV}} > s \\ \hat{y}_{\text{RV}} & \text{if } \hat{y}_{\text{RV}} \leq s \text{ and } \hat{y}_{\text{BEV}} < s , \\ \hat{y}_{\text{BEV}}, & \text{if } \hat{y}_{\text{RV}} \leq s \text{ and } \hat{y}_{\text{BEV}} > s \end{cases}$$

Qualitative



Ground Truth



Semantic KITTI





