

FOMO-3D: Using Vision Foundation Models for Long-Tailed 3D Object Detection

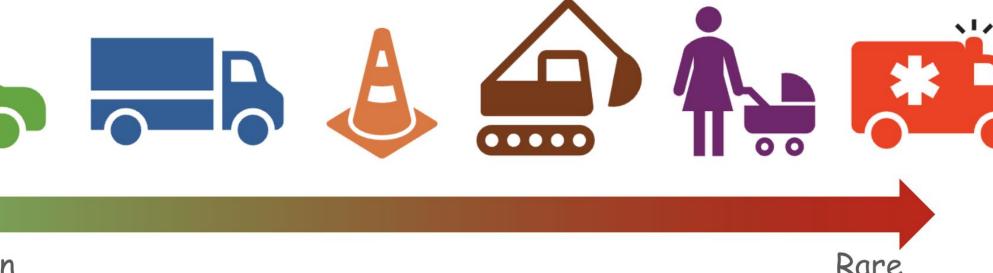
UNIVERSITY OF

detection.

Anqi Joyce Yang*, James Tu*, Nikita Dvornik, Thomas Li, Raquel Urtasun

* Equal Contribution

Long-Tailed 3D Detection

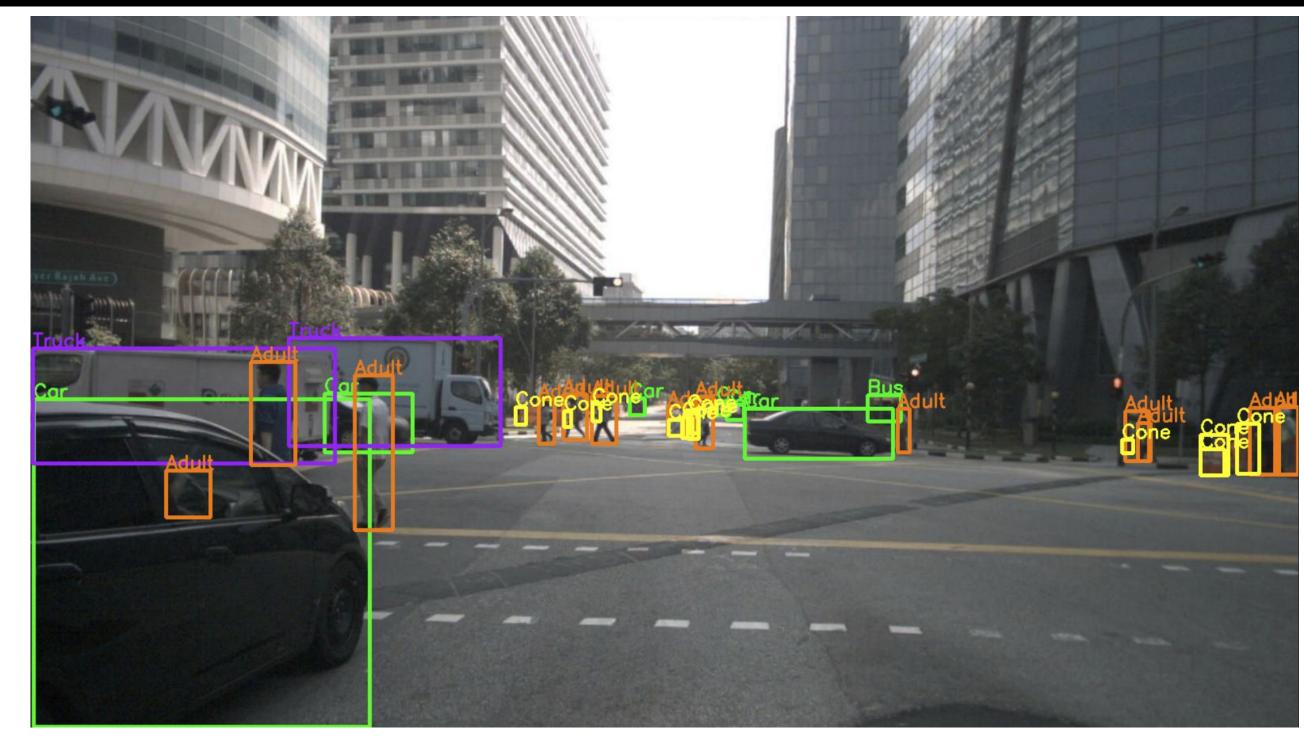


- To deploy a self-driving vehicle safely, it is crucial to detect objects from both common and rare classes, which could be challenging due to limited training data
- Foundation models show promise of bringing useful prior knowledge to tackle long-tailed 3D detection

Vision Foundation Models

OWLv2

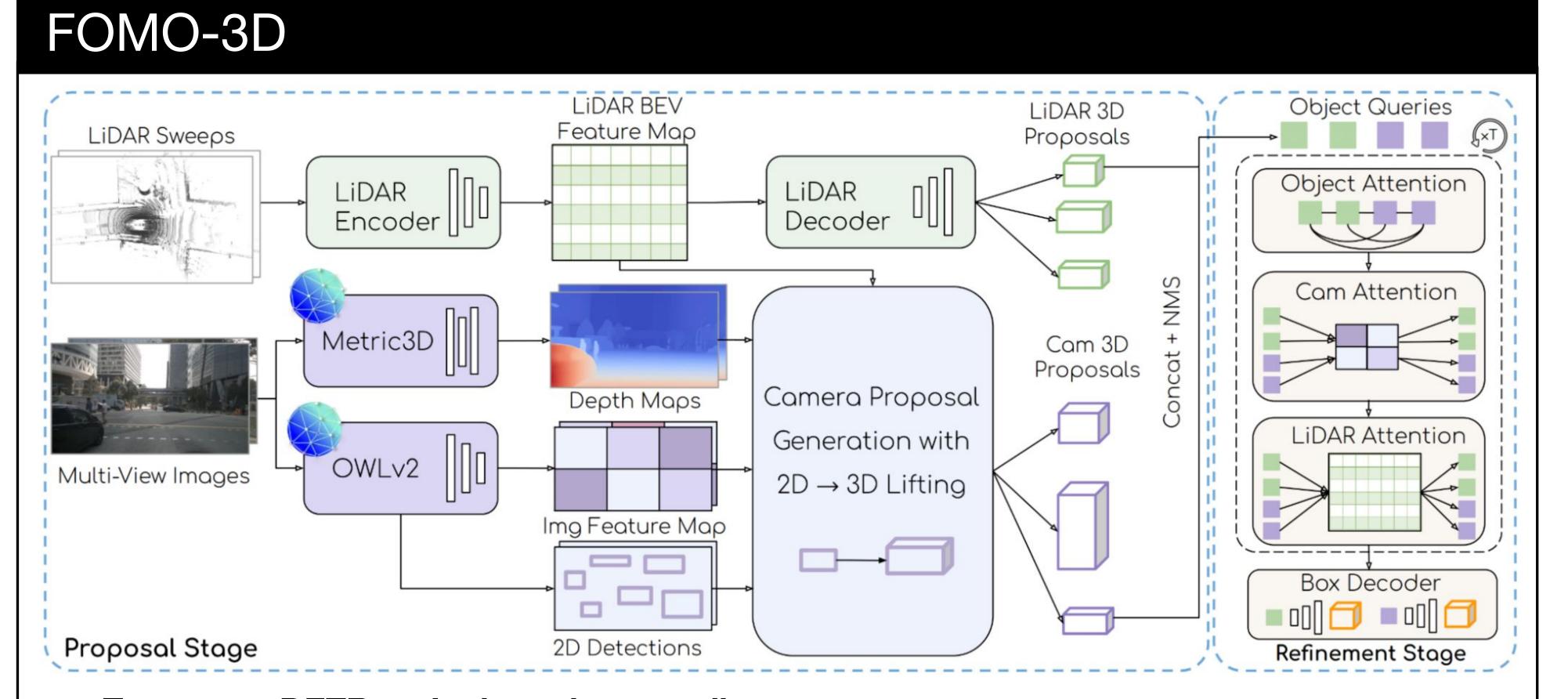
- Open-vocabulary 2D object detector
- Provides image embeddings
- Provides 2D detections from text prompts





Metric3Dv2 (M3D)

- Monocular depth estimation model
- Provides dense depth images



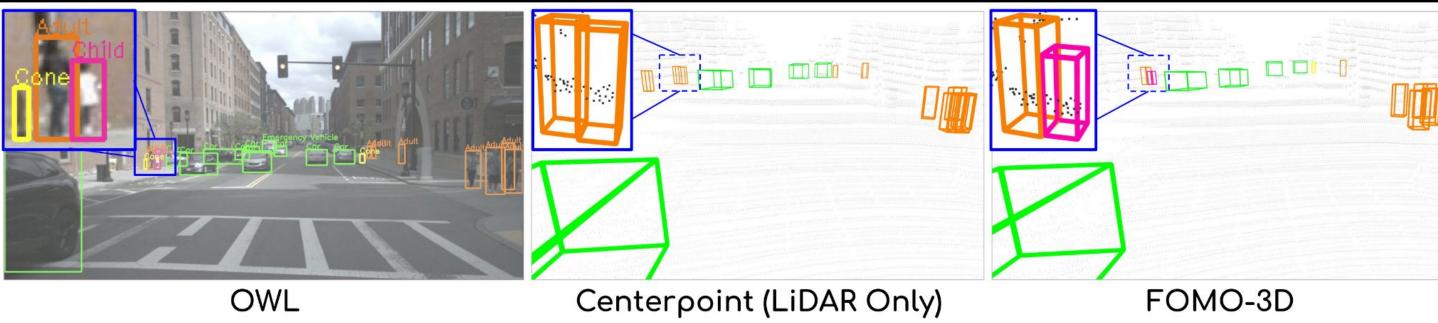
- Two-stage DETR-style detection paradigm
- First proposal stage LiDAR proposal branch (CenterPoint), novel camera proposal branch lifting 2D OWL boxes into 3D detections, using M3D depths and fusing with LiDAR features
- Second refinement stage object self-attention, object to LiDAR attention, object to **OWL** features attention

Camera Proposal IoU-Based Matching in Frustum

- . Obtain 2D detections and image features from OWL
- 2. Lift each 2D detection into 3D using M3D depths
- 3. Lift image features to 3D as a feature point cloud, and concat with LiDAR features as a fused BEV feature map
- 4. Refine 3D detection via attention to the BEV LiDAR and image features in the object frustum During training – frustum-aligned matching between proposals and ground truth for supervision

Quantitative Results Method BEVFormer [38] nuScenes: CenterPoint (Group-Free) [1, 5] BEVFusion-L [2] **FOMO-3D** outperforms SoTA detectors in all TransFusion [25] BEVFusion [2] aggregated object CMT [2] groups, esp. on *Few* $MMF^* (w/OWL + M3D) [5]$ MMLF [24] FOMO-3D MMLF Highway: FOMO-3D (no cam prop) FOMO-3D FOMO-3D shines for rare classes and long-range

Qualitative Results



d but has false positive c OWL detects child

Towed Object

- LiDAR-only misclassified child as adult
- FOMO-3D combines multimodal information to successfully detect the chi reject the false positive con
- OWL detects cone but duplicates person and cyclist
- LiDAR-only performs poorly on c and person
- FOMO-3D performs well on c e, person, without duplicating person/cyclist
- All share a failure mode confusing trailer and truck

