

POP-3D: OPEN-VOCABULARY 3D OCCUPANCY PREDICTION FROM IMAGES

Antonin Vobecky^{1,2}, O. Siméoni², D. Hurych², S. Gidaris², A. Bursuc²,
P. Pérez², J. Sivic¹



²valeo.ai

Motivation and Goal

Task: open-vocabulary 3D occupancy prediction

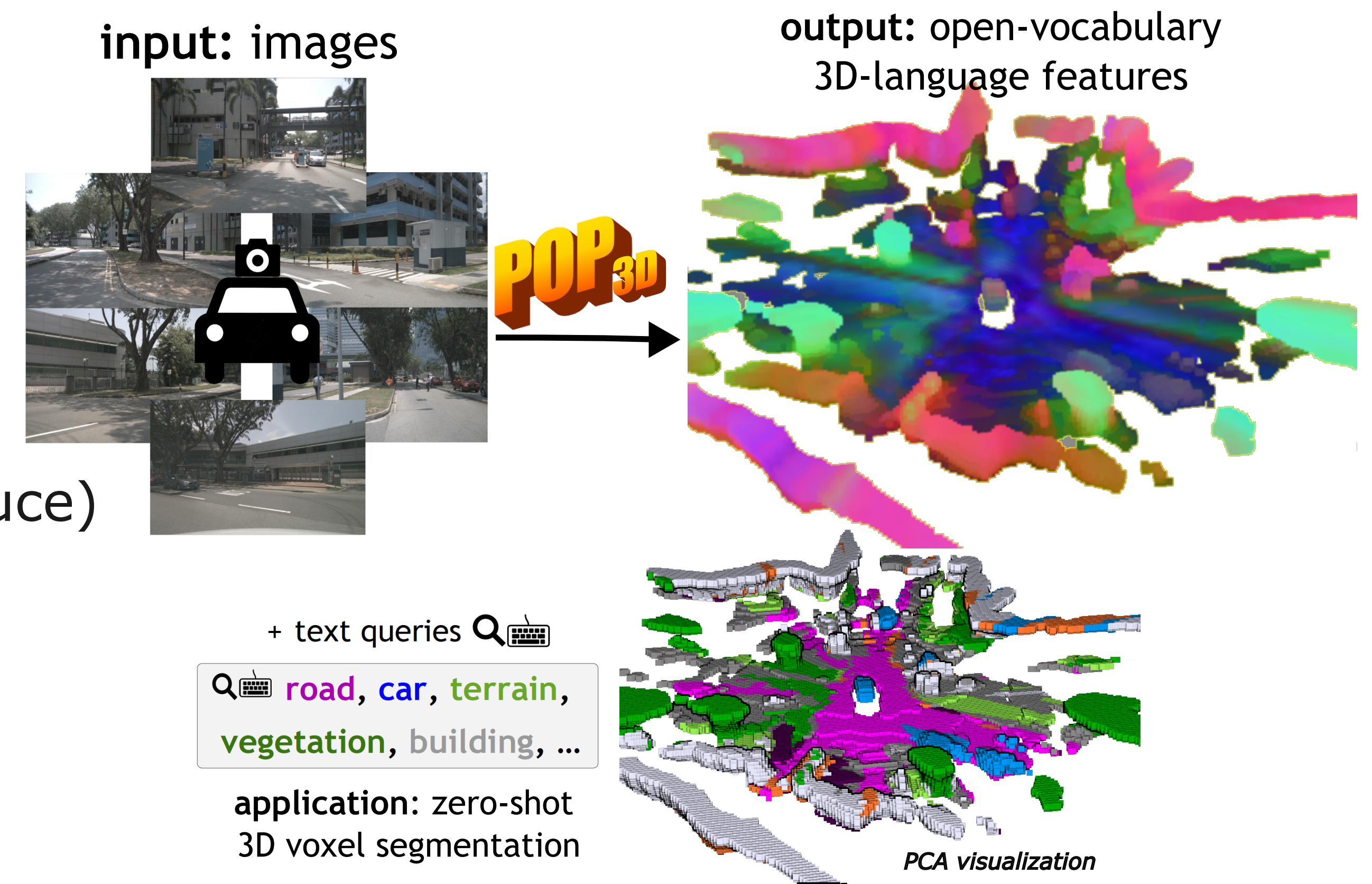
- generate 3D occupancy maps from surrounding images, labeling occupied voxels with semantic classes defined using open-form text

Motivation: Tackle the following issues:

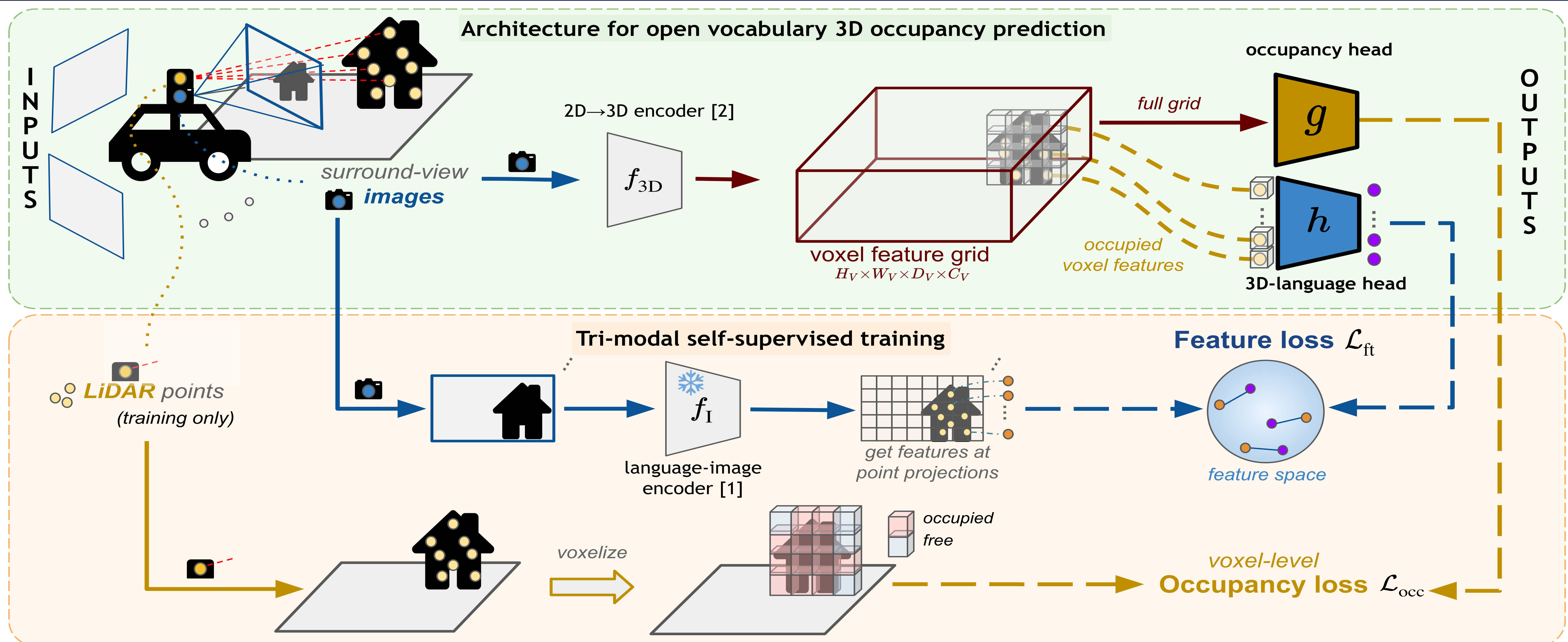
- annotations constrained by closed-vocabulary class definitions
- limited availability of sparse LiDAR annotations (expensive to produce)

Goal:

- open-vocabulary 3D semantic occupancy prediction
- *training*: unlabeled image-LiDAR data + pre-trained image-language model
- *inference*: images only



Training



Given surround-view images on the input, our method first produces voxel feature grid using 2D-3D encoder.

Two heads on top of the feature grid:

- occupancy head g** : Outputs **voxel grid occupancy**. Trained by occupancy from LiDAR points.
- 3D-language head h** : Outputs **open-vocabulary 3D-language features**. Trained by language-image features from a frozen image encoder [1].

Results

