



MOSS: Mask-Oriented Open-Set for 3D Scene Segmentation using Superpoint

Nov 2024 - ?

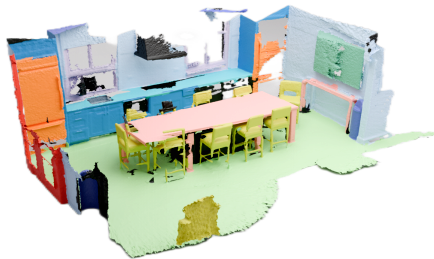
Hongrui Wu
Tongji University

Task Background

- Traditional methods rely on densely annotated 3D scenes.
- Have to utilize supervision from ground truth labels.
- 3D Annotation is **Time-consuming** and **Expensive!**



Input 3D
Geometry



Annotated
3D scenes

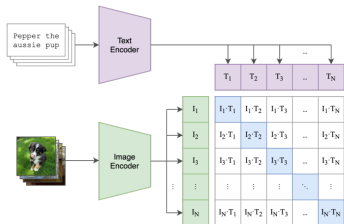
Task Background

Meanwhile computer vision is going through a transition from the previous **closed-set** perception to **open-set** perception:

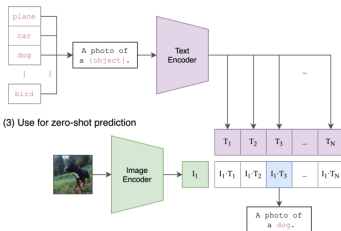
Closed-set: only handles predefined classes during training and has limited capability in dynamic world

Open-set: understands unseen, diverse and free-flowing language, mimicking how humans naturally interact with the world and each other

(1) Contrastive pre-training



(2) Create dataset classifier from label text

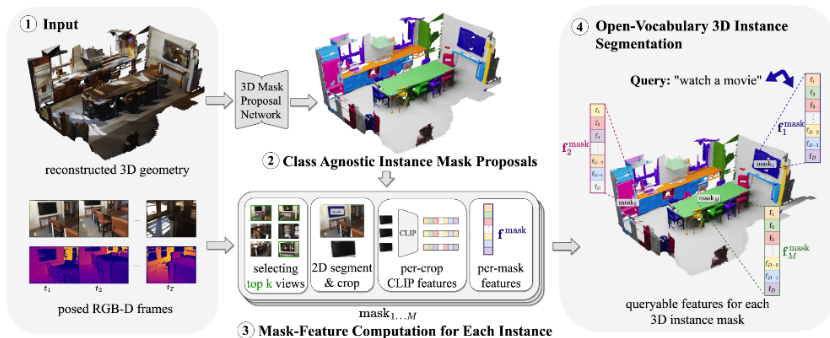


(3) Use for zero-shot prediction



2D open-set tasks can now understand new concepts, perform accurate segmentation and detection, and handle complex tasks requiring reasoning.

Example Pipeline

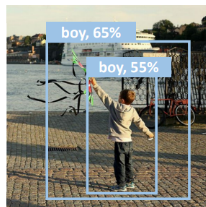


- ① Input RGB-D images of a 3D indoor scene and its point cloud.
- ② Generate class-agnostic instance mask proposals using the point cloud.
- ③ Compute feature representations for each mask.
- ④ Obtain an open-vocabulary 3D instance segmentation, enabling retrieval of objects related to queried concepts in the CLIP

Motivation

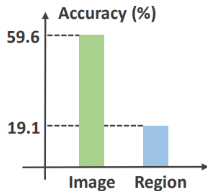
- Previous methods simply use the original CLIP model without addressing its regional limitations.
- Directly applying it for object detection leads to poor performance due to domain shift, as CLIP was trained to match whole images to text descriptions, without capturing fine-grained alignment between image regions and text spans.

Cropped image regions recognized by CLIP

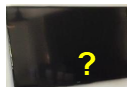


a

Image classification (ImageNet)
Region classification (LVIS)

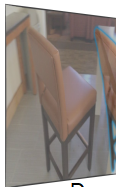


b



Motivation

- 2D instance segmentation achieves impressive results and is widely used in 3D open-set tasks.
- However, it is important to note that model sensitivity can cause mismatches between 2D segmented regions and 3D mask proposals, which can lead to mismatched inference results in subsequent stages.



Demo from SAM
official repo

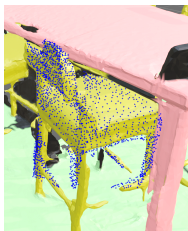


Motivation

- Most existing work heavily relies on the mask proposals generated by pre-trained 3D models (like Mask3D), where the quality of these masks directly affects the performance of instance segmentation.
- However, open-set tasks should not be constrained by closed-set models. Additionally, prior knowledge from 2D segmentation models can alleviate the limitations observed in current 3D class performance.



Mask3D



3D Points to 2D Pixels

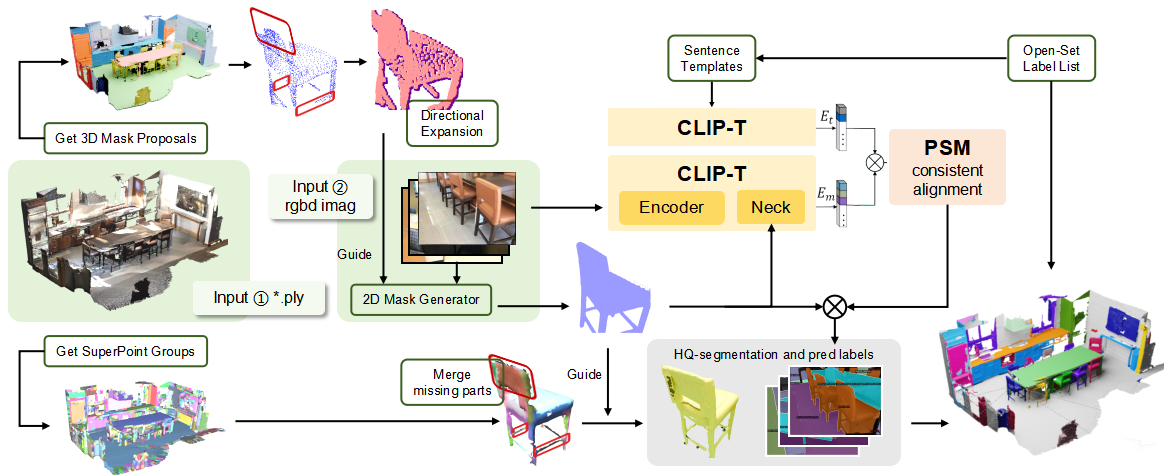


Ground truth

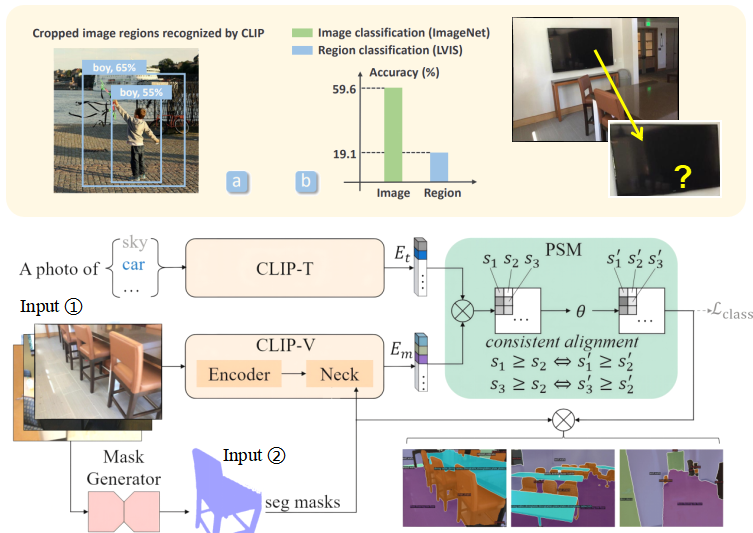
Contributions

- We proposed **MOSS**, a mask-based framework for open-set 3D scene semantic segmentation that enables efficient cross-dimensional feature transfer and inference.
- We enhanced the frame by implementing **global information input with mask** constraints to strengthen attention.
- We employed **a density-guided dilation algorithm** to optimize the matching precision between 2D and 3D masks.
- We also introduce a novel method to enhance 3D mask proposals, which leverages 2D prior knowledge to perform **back-projection** on a 3D pre-trained model. This approach guides **the capture of superpoint clusters** in the 3D scene, thereby improving the quality of the output results of **fine-tuning the close-set model result**.

Our Proposed MOSS



Our Proposed MOSS - Contribution1



Input a cropped image:

- Cropped regions lack global information.

Input a purely global image:

- Fail to localize the areas that need to be understood.
- Lead to **inconsistencies in the granularity of classification**, where the level of detail may be too fine or too coarse to align with the designated regions of the mask proposal.

Input both global images and masks:

- Constrain regions requiring enhanced understanding.
- Obtain contextual information to improve inference accuracy.

- **How to obtain a high-quality 2D mask?**

Our Proposed **MOSS** - Contribution2

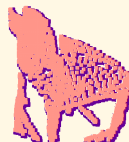
Density-based Directional Expansion Algorithm



IOU?
The number of pixels does not correspond



$$\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$



① Expansion in the equal direction

② Directional expansion based on density

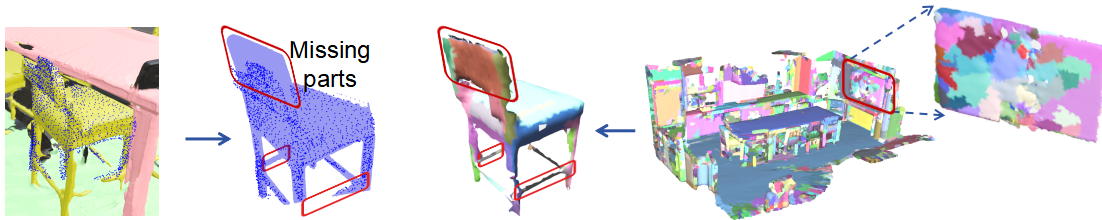
③ ...

④ Calculate IoU > the threshold

- **What else can a high-quality 2D mask offer?**

Our Proposed **MOSS** - Contribution3

Fine-Tuning of 3D Mask Proposals based on SuperPoints



- The coarse 3D mask proposal has missing areas compared to the ground truth.
- The prior knowledge from the 2D can fill these gaps.
- Pixels from the 2D mask are projected to 3D, and the matching points are added to the 3D mask.



Ground truth

- To improve efficiency, the raw point cloud is transformed into superpoint clusters.
- SuperPoints: points are grouped into geometrically homogeneous regions.
- Instead of individual points, superpoint clusters are used as the unit for merging.

Experiments & Results (Imcomplete)

Method	mAP	mAP50	mAP25	head	comm	tail
Mask3D (Closed Vocab.)	26.9	36.2	41.4	39.8	21.7	17.9
SAM3D	6.1	14.2	21.3	7	6.2	4.6
OVIR-3D	13	24.9	32.3	14.4	12.7	11.7
Open3DIS	23.7	29.4	32.8	27.8	21.2	21.8
OpenScene (2D Fusion)	11.7	15.2	17.8	13.4	11.6	9.9
OpenScene (3D Distill)	4.8	6.2	7.2	10.6	2.6	0.7
OpenScene (2D-3D Ens.)	5.3	6.7	8.1	11	3.2	1.1
OpenMask3D	15.4	19.9	23.1	17.1	14.1	14.9
OpenMask3D(Our Step2)	16.2	21.3	24.8	22.2	13.4	12.5
Open3DIS	18.6	23.1	27.3	24.7	16.9	13.3
Open-YOLO 3D	24.7	31.7	36.2	27.8	24.3	21.6
MOSS(Ours)	25.2	32.7	36.5	25.6	25.2	25.9

* The experimental results are still being updated
(hyperparameters are being finalized)

Phase Summary & Next Work

- Tasks in the 2D domain can successfully guide 3D spatial understanding tasks;
- The instance matching between 2D and 3D in the form of masks works effectively
- Conduct experiments with additional datasets;
- Attempt to replace the black-box CLIP model with an integrated VLM (Visual-Language Model) during the inference phase;
- Optimize code details to reduce the inference time per scene.



Research Statement

Summer Research Internship Ver.

My research interests:

- Currently lie in open-set tasks for 3D scenes and AIGC for 3D vision.
- **In the future**, I aim to explore the application of generative models in 3D reconstruction tasks and 3D scene re-editing.

My preferred paper from the group:

- ...