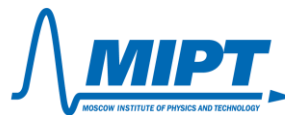


Нейросетевые методы построения мультимодальных карт и их использование для навигации и управления роботами

Дмитрий Александрович Юдин

к.т.н., заведующий лабораторией интеллектуального транспорта МФТИ - НКБ ВС,
Центр когнитивного моделирования МФТИ, старший научный сотрудник AIRI



О чем поговорим

- 01 Мотивация
- 02 Методы построения мультимодальных карт
- 03 Плотные (Dense) методы
- 04 Разреженные (Sparse) методы
- 05 Базовые модели сегментации изображений с открытым словарем запросов (open vocabulary)
- 06 Наши исследования и приложения
- 07 Направления дальнейшего развития

01



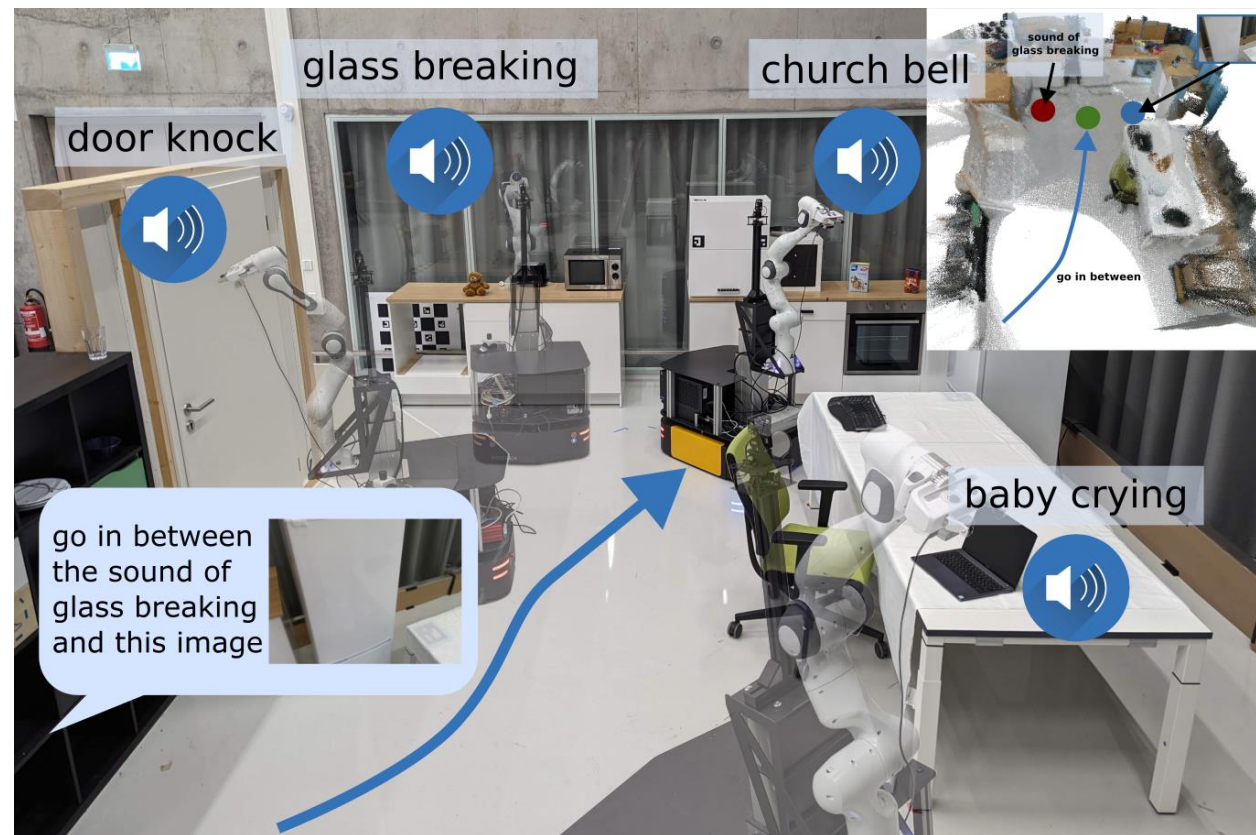
Мотивация

Будущее построения 3D-карт: пространственные запросы и ответы по трехмерным картам

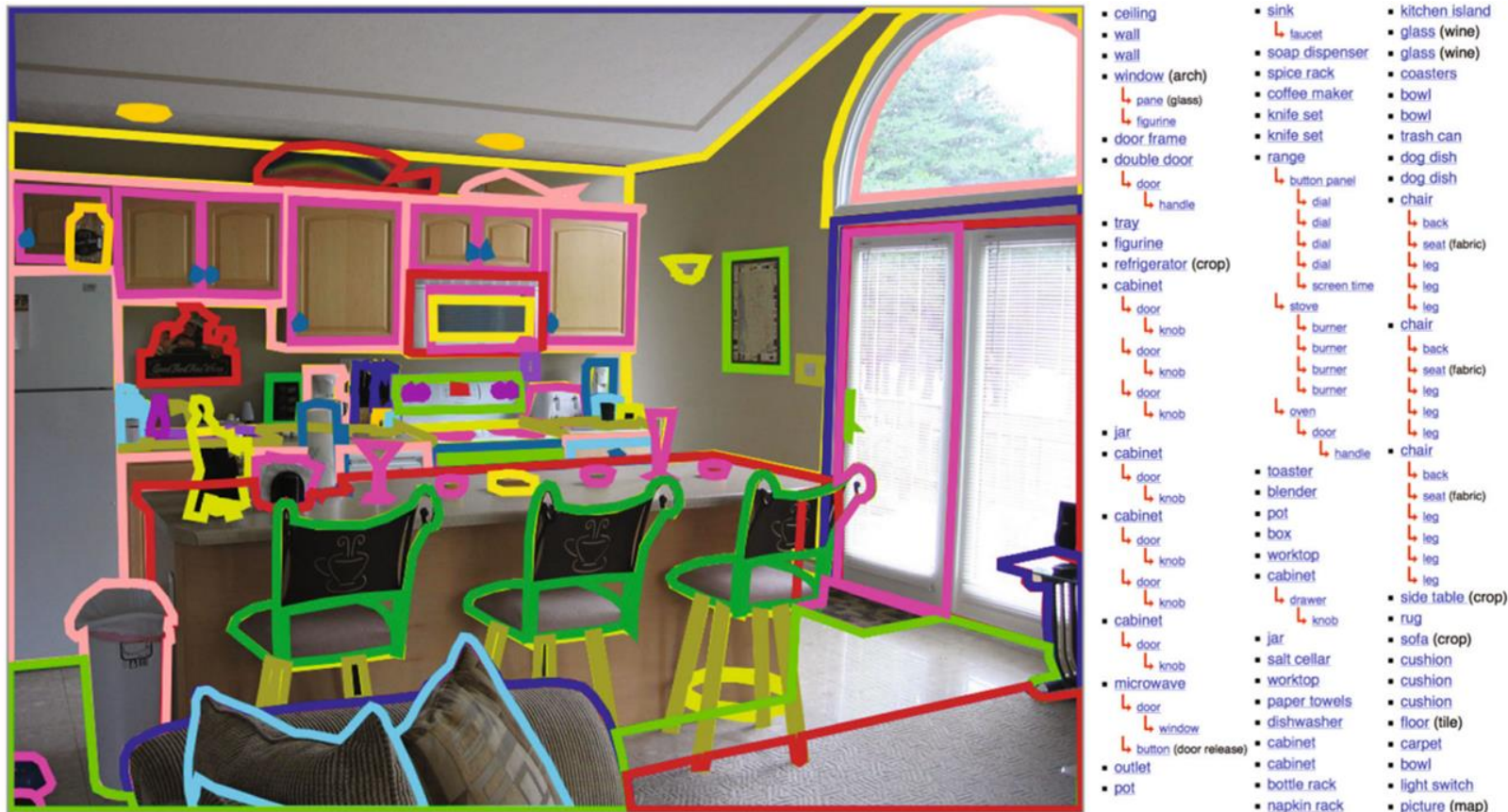
- **Текстовый запрос к 3D карте.**
Пример: “Место для сидения (Sit)”



- **Мультимодальный запрос к 3D карте.**
Пример: “Локация между местом, где был слышен звук разбитого стекла, и местом, показанном на фотографии”



Что делать с тысячами категорий объектов, которые надо сегментировать?



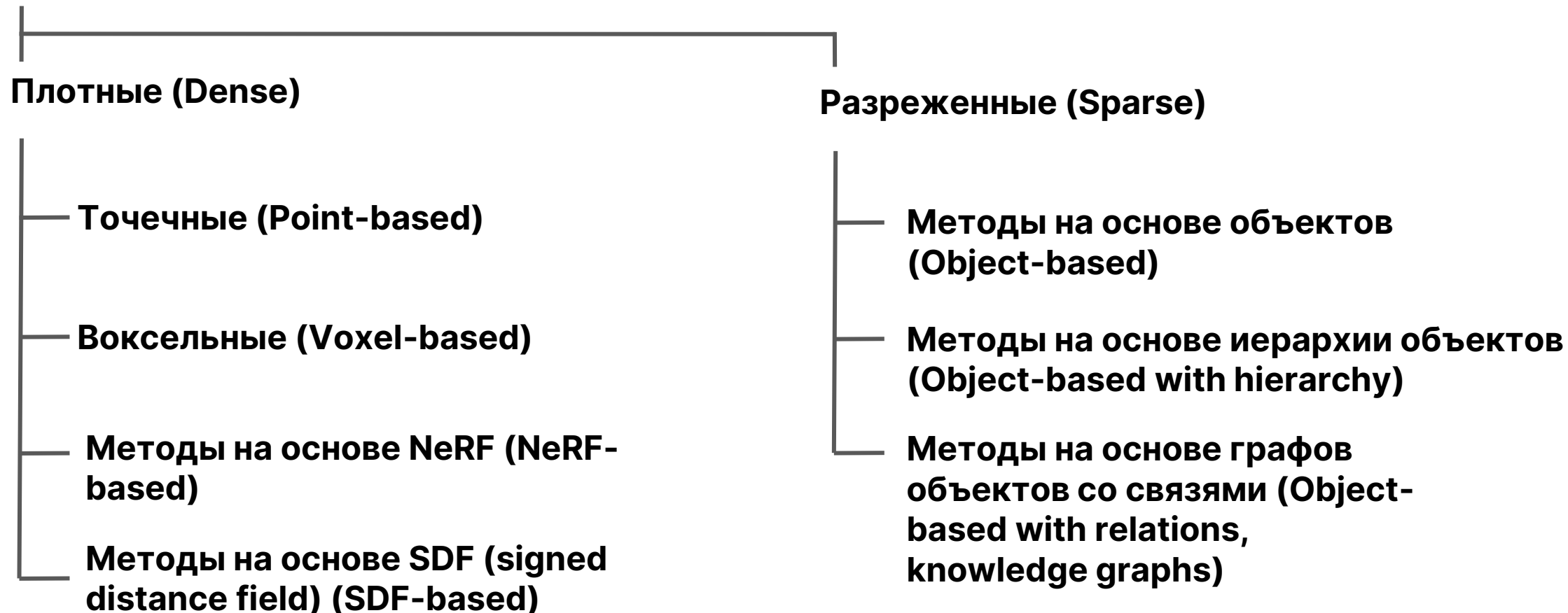
- Популярный набор данных ADE20K содержит 27,574 изображений внутри и вне помещений с 707,868 объектами, принадлежащими к 3,688 категориям(классам)
- Что делать, если мы хотим сегментировать объект, описанный на естественном языке?
- При подготовке собственных наборов данных мы хотим размечать объекты на изображениях в несколько кликов мыши

02

Методы построения мультимодальных карт

Какие методы бывают?

Методы построения мультимодальных карт



03

Плотные (Dense) методы

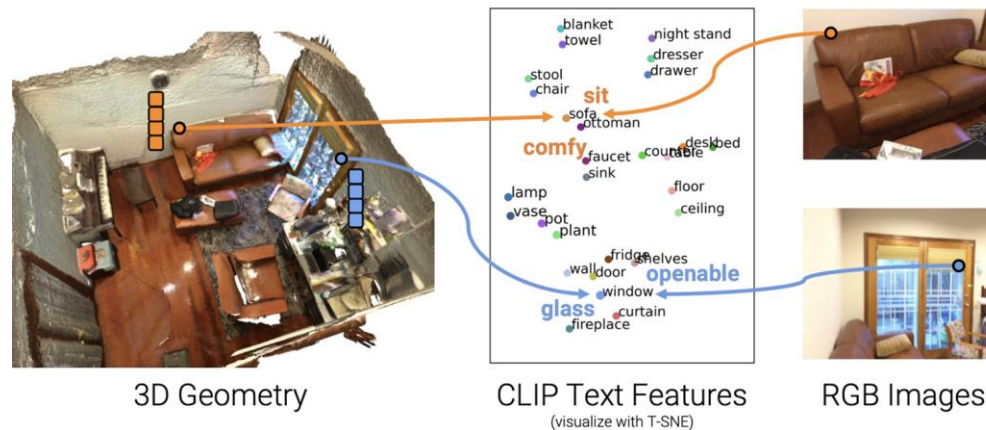
Постановка задачи

$$PCL_{Query} = \mathcal{A}((RGB_1 \dots RGB_N), [D_1 \dots D_N], (Pose_1 \dots Pose_N), Query), Query = (Text, [Image]).$$

$$I : Map = \mathcal{M}((RGB_1 \dots RGB_N), [D_1 \dots D_N], (Pose_1 \dots Pose_N)) = (Embedding_1 \dots Embedding_K),$$

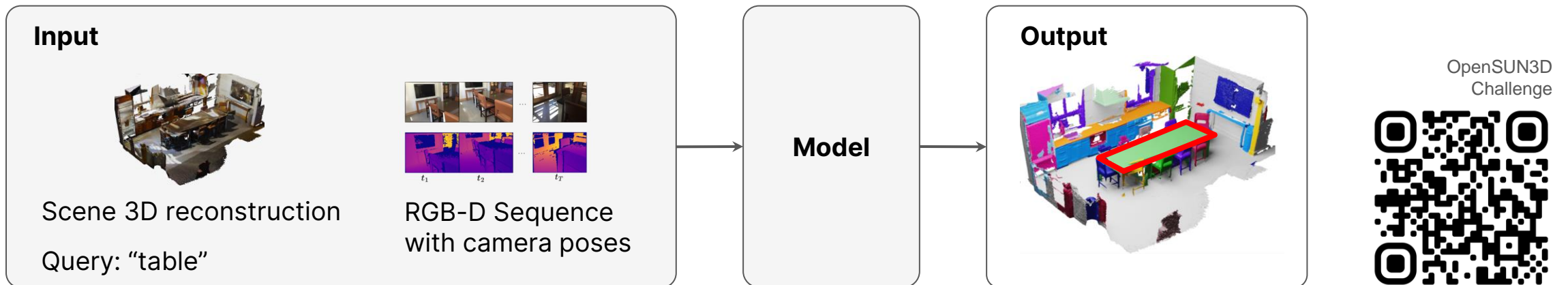
where N – image sequence length, K – number of map points.

$$II : PCL_{Query} = \mathcal{Q}(Map, Query) = (X_1 \dots X_K, Y_1 \dots Y_K, Z_1 \dots Z_k, Score_1 \dots Score_K).$$

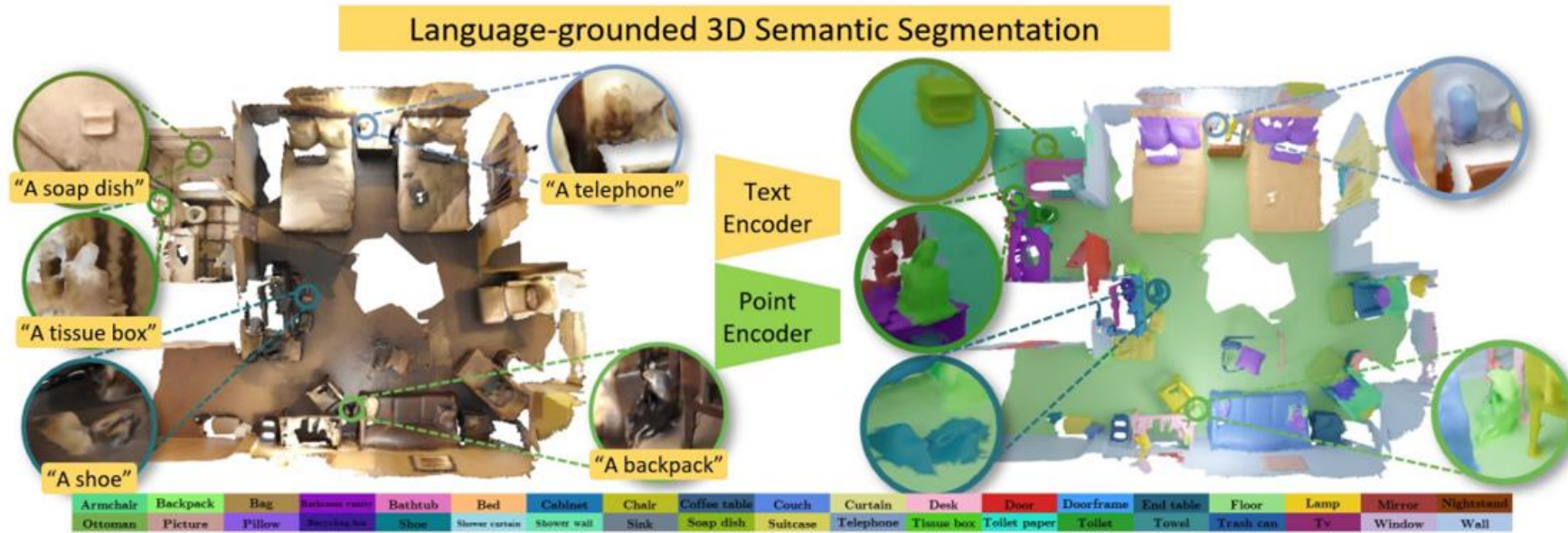


OpenSUN3D Workshop Challenge on 3D Open-Vocabulary Scene Understanding (ICCV2023)

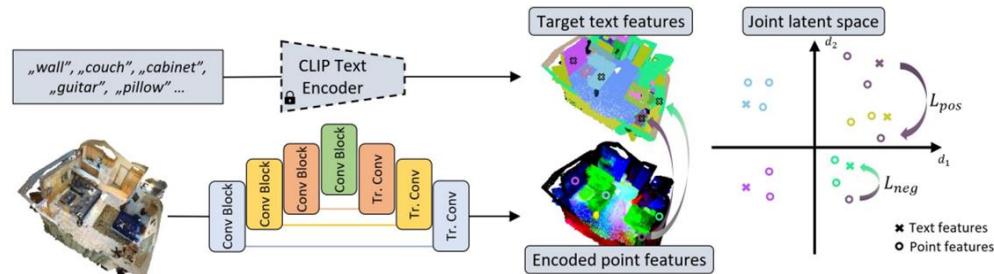
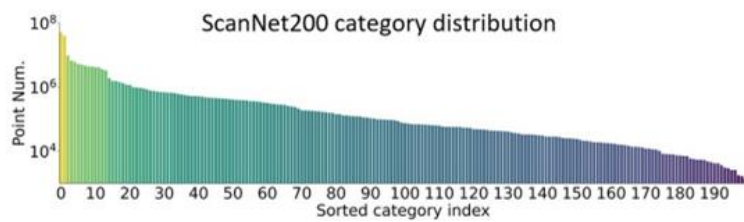
- TASK: Given an open-vocabulary, text-based query, the aim is to localize and segment the object instances that fit best with the given prompt, which might describe object properties such as semantics, material type, affordances and situational context.
- INPUT: An RGB-D sequence and the 3D reconstruction of a given scene, camera parameters, and a text-based input query.
- OUTPUT: Instance segmentation of the point cloud that corresponds to the vertices of the provided 3D mesh reconstruction, segmenting the objects that fit best with the given prompt.



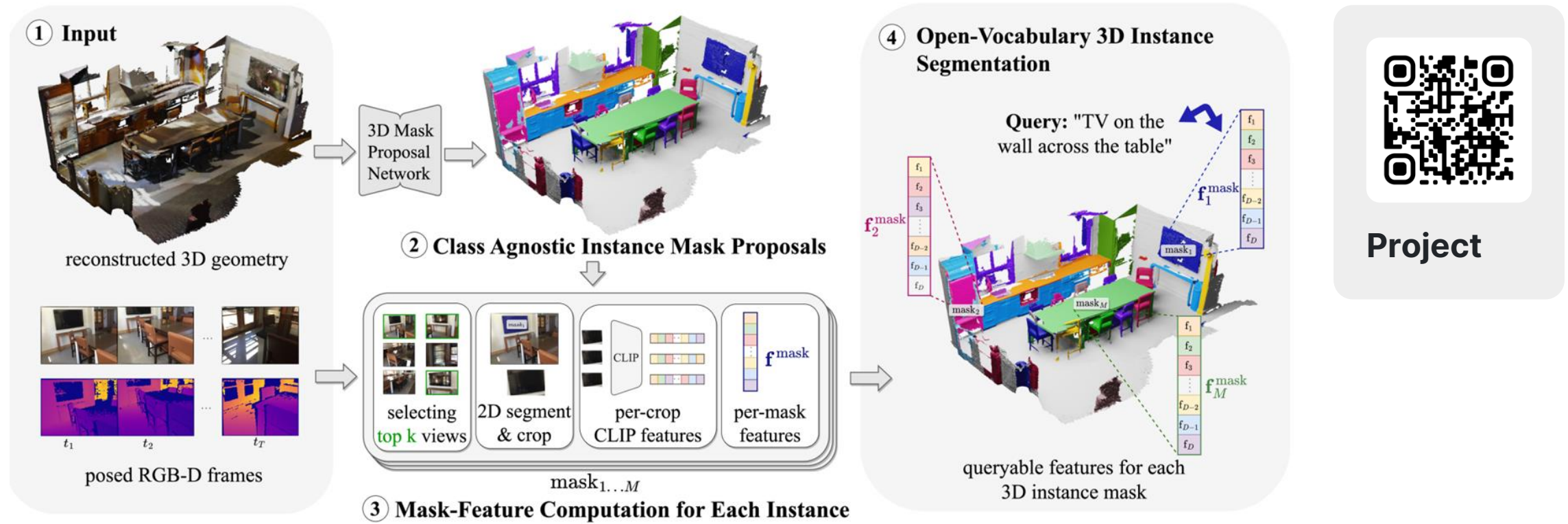
Language-Grounded Indoor 3D Semantic Segmentation in the Wild



Project
Github

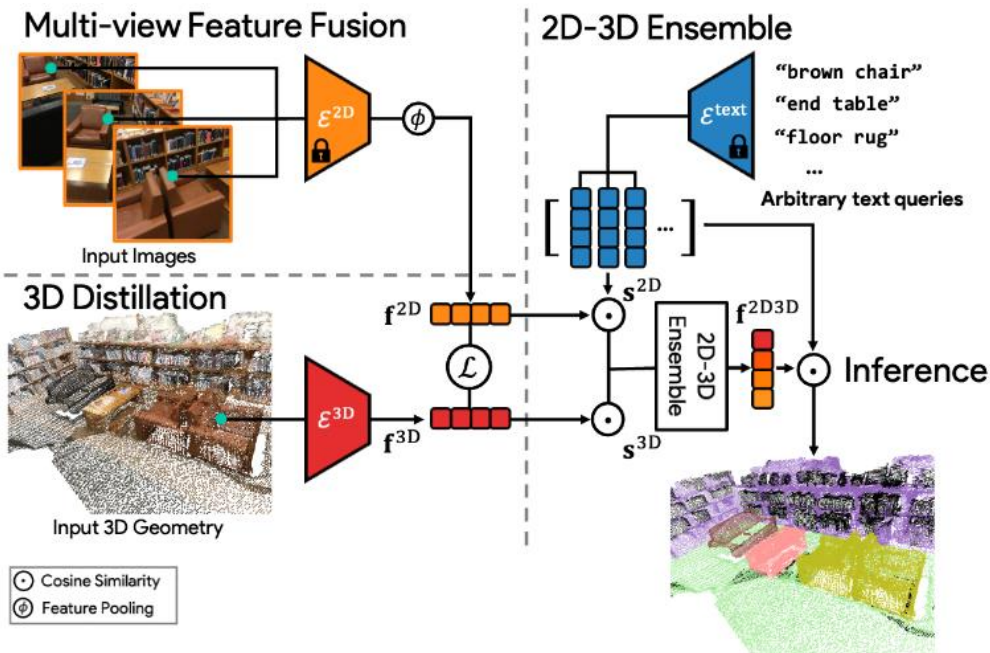


OpenMask3D: Open-Vocabulary 3D Instance Segmentation

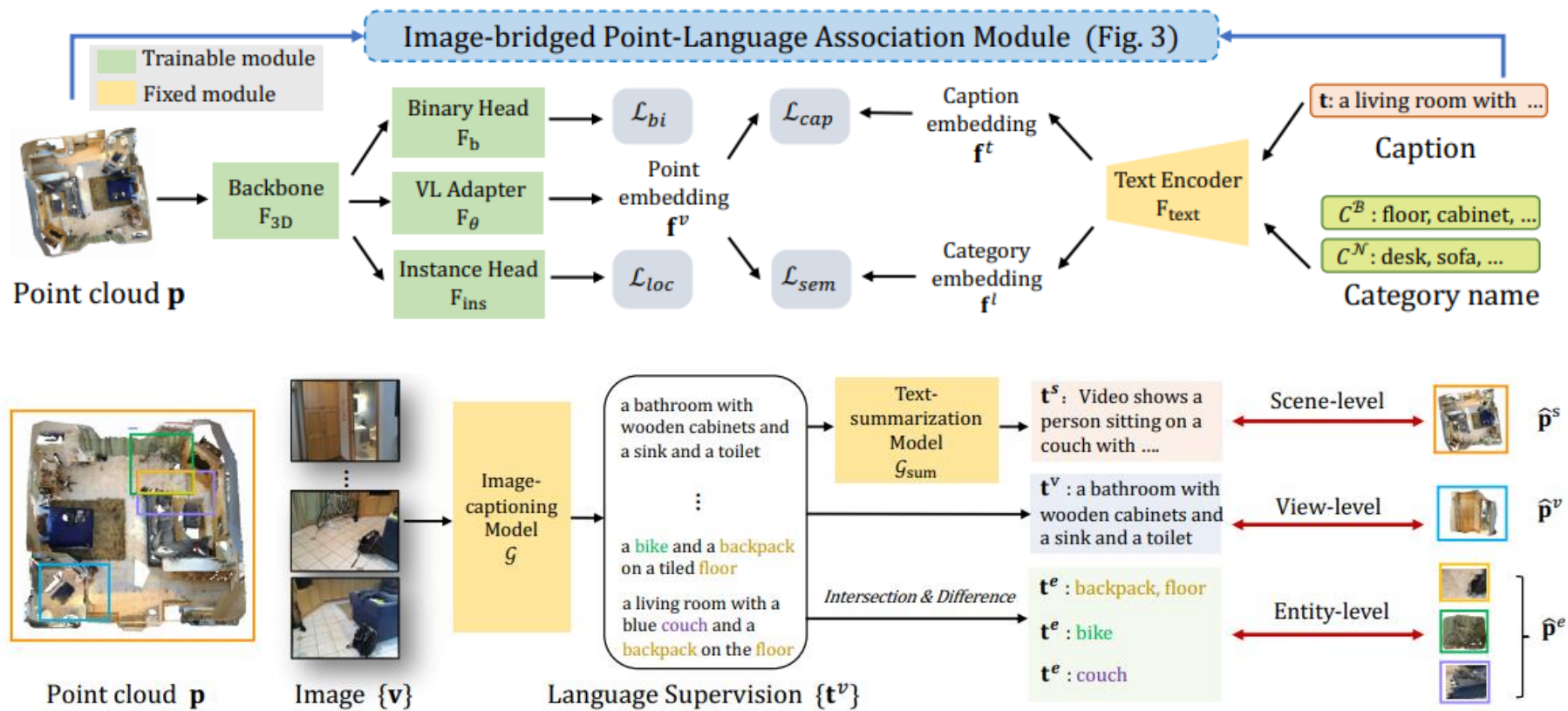


OpenScene: 3d scene understanding with open vocabularies

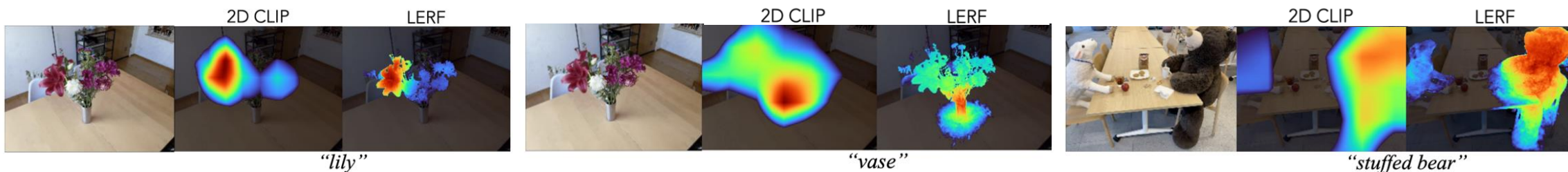
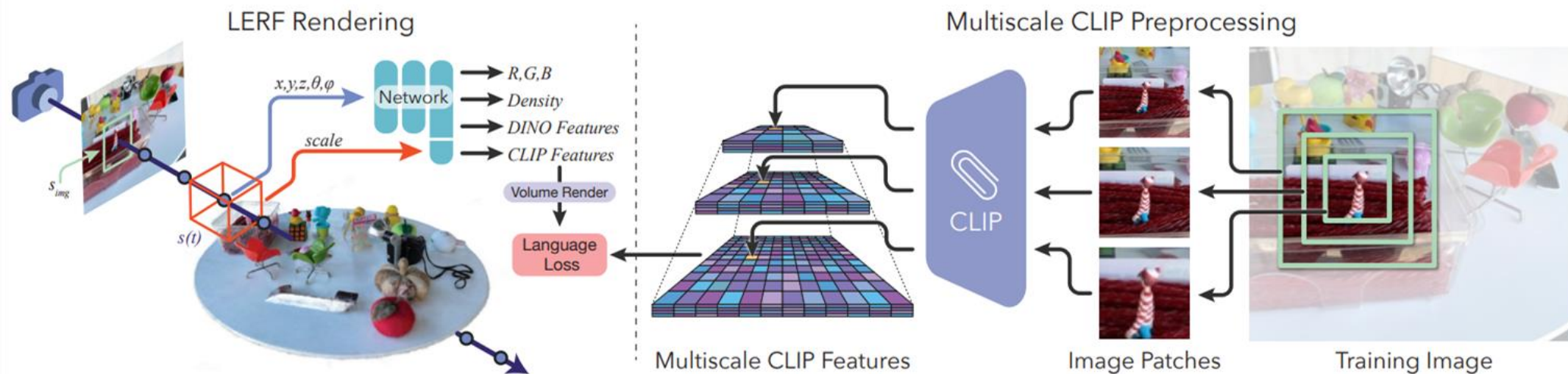
Produce text-image-3D co-embeddings



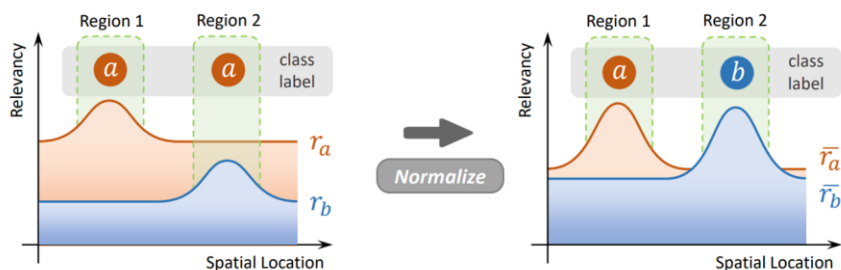
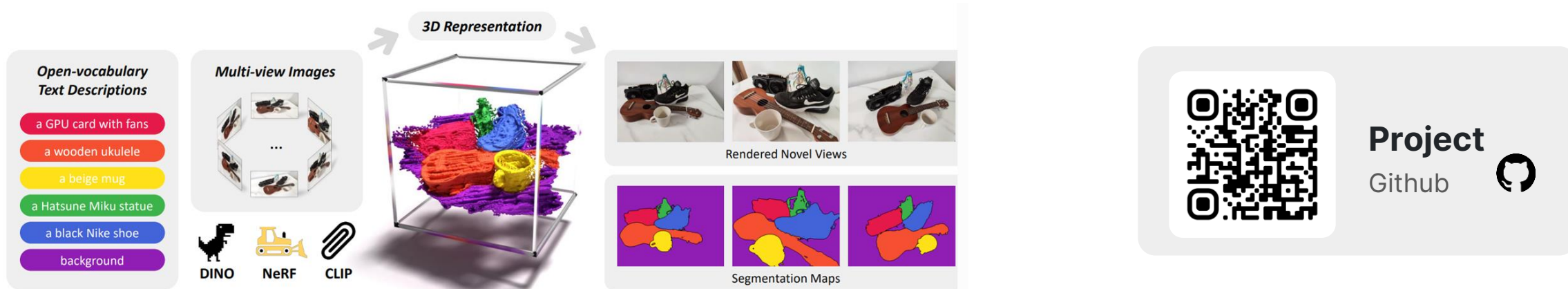
PLA: Language-Driven Open-Vocabulary 3D Scene Understanding



LERF: Language Embedded Radiance Fields



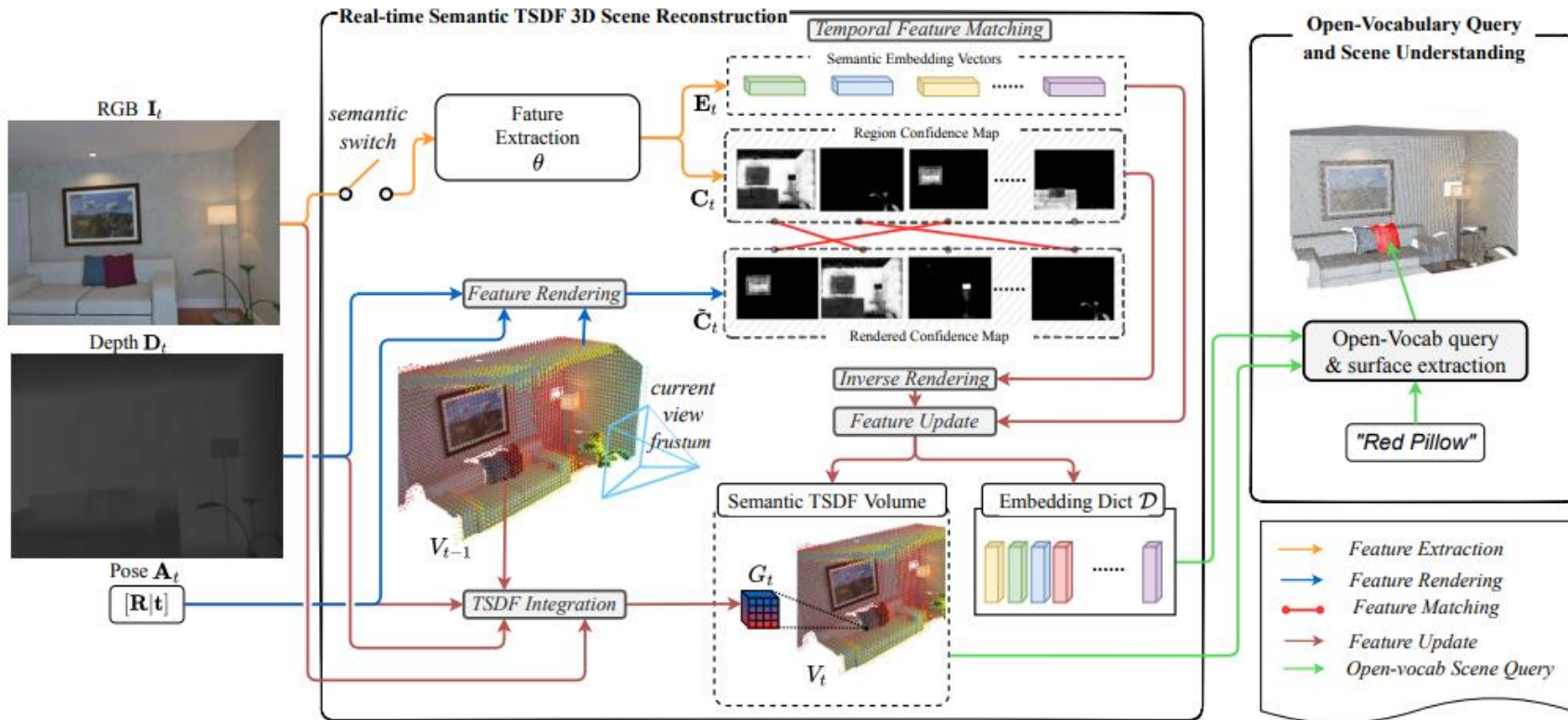
3D-OVS: Weakly Supervised 3D Open-vocabulary Segmentation



Methods	<i>bed</i>		<i>sofa</i>		<i>lawn</i>		<i>room</i>		<i>bench</i>		<i>table</i>	
	mIoU	mAP	mIoU	mAP	mIoU	mAP	mIoU	mAP	mIoU	mAP	mIoU	mAP
2D												
LSeg [8]	56.0	87.6	04.5	16.5	17.5	77.5	19.2	46.1	06.0	42.7	07.6	29.9
ODISE [16]	52.6	86.5	48.3	35.4	39.8	82.5	52.5	59.7	24.1	39.0	39.7	34.5
OV-Seg [19]	79.8	40.4	66.1	69.6	81.2	92.1	71.4	49.1	88.9	89.2	80.6	65.3
3D												
FFD [4]	56.6	86.9	03.7	09.5	42.9	82.6	25.1	51.4	06.1	42.8	07.9	30.1
Sem(ODISE) [45]	50.3	86.5	27.7	22.2	24.2	80.5	29.5	61.5	25.6	56.4	18.4	30.8
Sem(OV-Seg) [45]	89.3	96.7	66.3	89.0	87.6	95.4	53.8	81.9	94.2	98.5	83.8	94.6
LERF [2]	73.5	86.9	27.0	43.8	73.7	93.5	46.6	79.8	53.2	79.7	33.4	41.0
3D-OVS	89.5	96.7	74.0	91.6	88.2	97.3	92.8	98.9	89.3	96.3	88.8	96.5

Mitigating CLIP features' ambiguities with normalized relevancy maps.

Open-Fusion: Real-time Open-Vocabulary 3D Mapping and Queryable Scene Representation



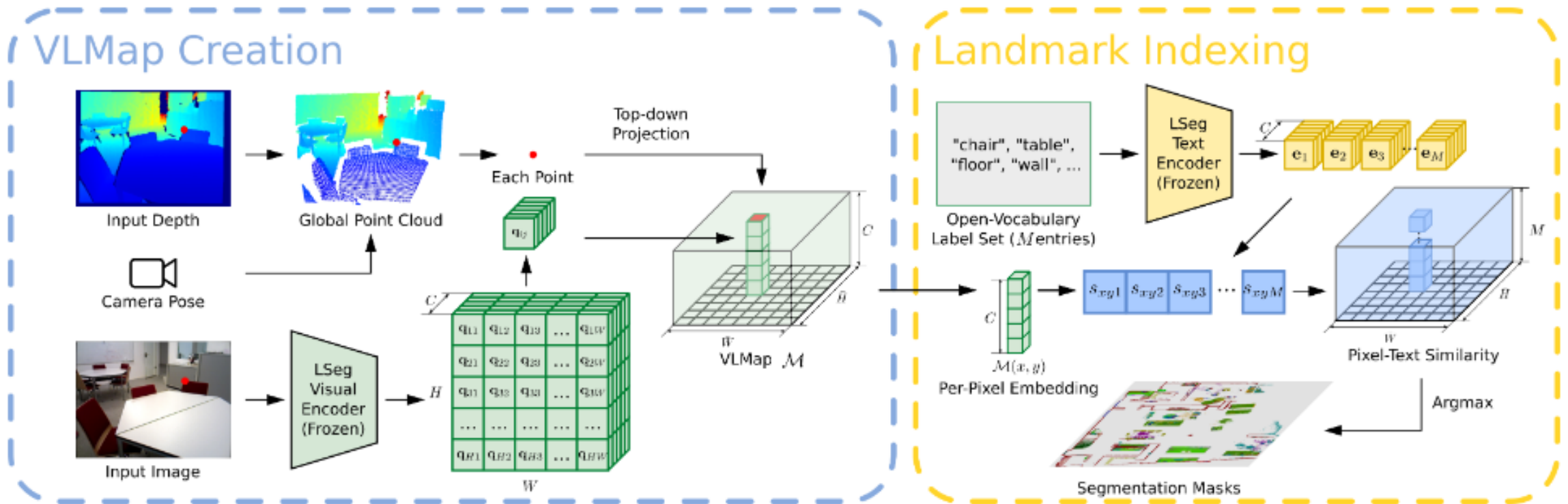
QUANTITATIVE COMPARISON OF OPEN-SET SEMANTIC SEGMENTATION AND 3D SCENE REPRESENTATION TIME BETWEEN OPEN-FUSION AND EXISTING METHODS ON THE SCANNET DATASET.

	Method	Time (FPS) \uparrow		Accuracy \uparrow	
		3D-Rec. ¹	Sem-3D-Rec ²	mAcc	f-mIoU
Priv. ₃	LSeg	-	-	0.70	0.63
	OpenSeg	-	-	0.63	0.62
	CLIPSeg (rd64-uni)	-	-	0.41	0.34
	CLIPSeg (rd16-uni)	-	-	0.41	0.36
ZS. ₄	MaskCLIP	-	-	0.24	0.28
	ConceptFusion	1.5	0.15	0.63	0.58
	Open-Fusion	50	4.5	0.62	0.59

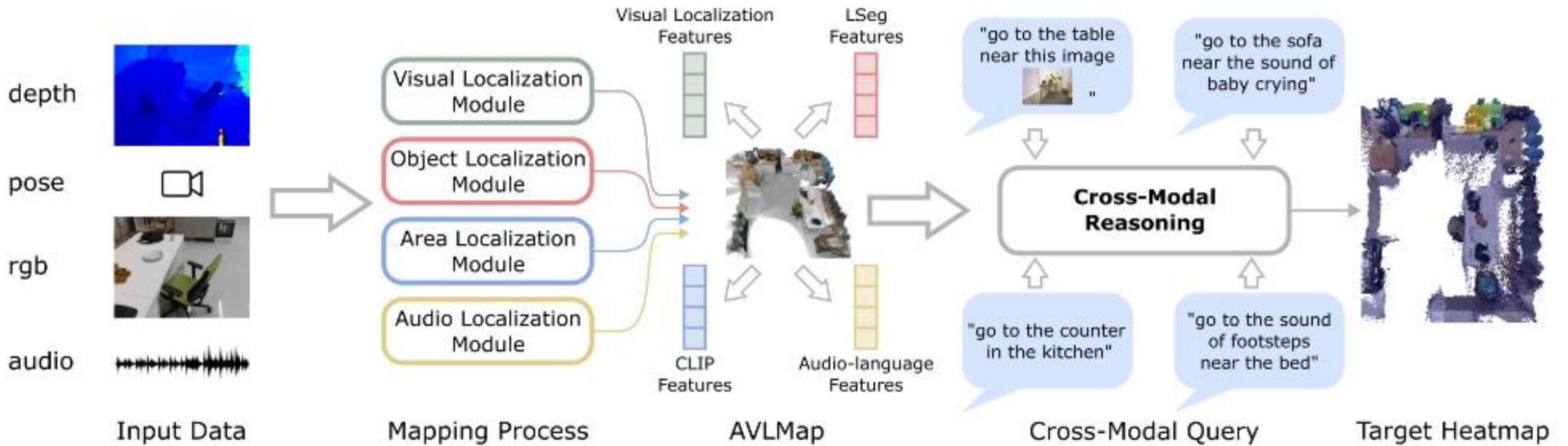
Map	Method	Representation	Foundation Model	Feature Level	Real-time ¹	Scene-specific	Sem-Query ²
2D	CoW [8]	point	CLIP [3] + GradCAM [9]	point	-	\times	$O(P)$
	NLMap [10]	point	ViLD [11] + CLIP [3]	bbox	-	\times	$O(P)$
	VLMMap [12]	point	LSeg [13]	point	\times	\times	$O(P)$
3D	CLIP-Fields [14]	NeRF	Detic [15] + CLIP [3]	bbox	\times	\checkmark	-
	LERF [16]	NeRF	CLIP [3]	image patch	\times	\checkmark	-
	SemAbs[17]	occupancy	CLIP [3] + GradCAM [9]	point	\times	\times	$O(P)$
	ConceptFusion [18]	point	SAM [19] + CLIP [3]	bbox	\times	\times	$O(P)$
	Open-Fusion	TSDF	SEEM [20]	region	\checkmark	\times	$O(M)$



VLMaps: Visual language maps for robot navigation

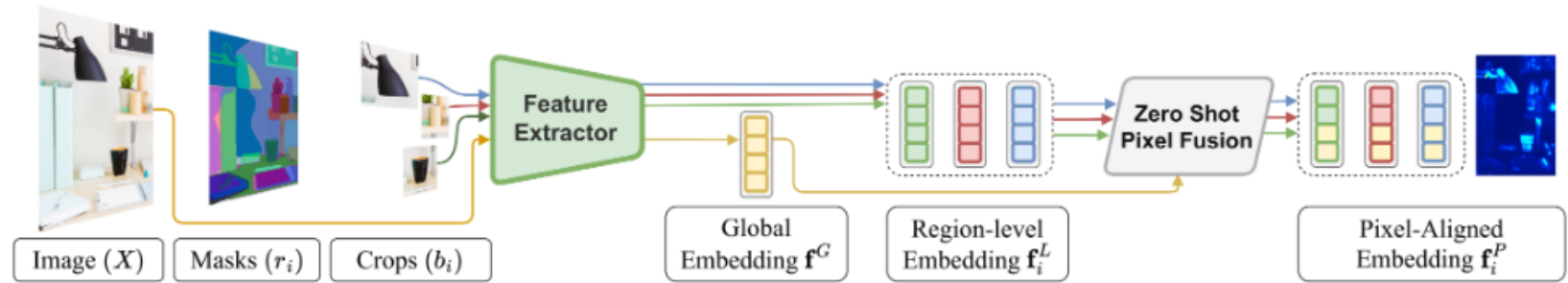


AVLMaps: Audio visual language maps for robot navigation

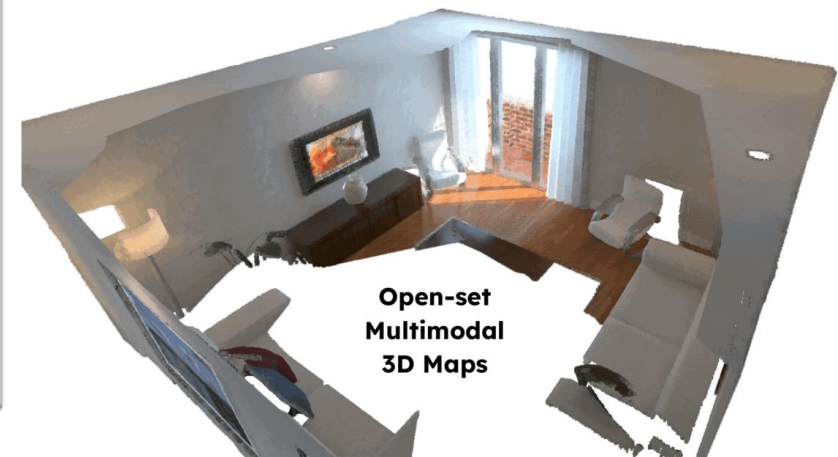
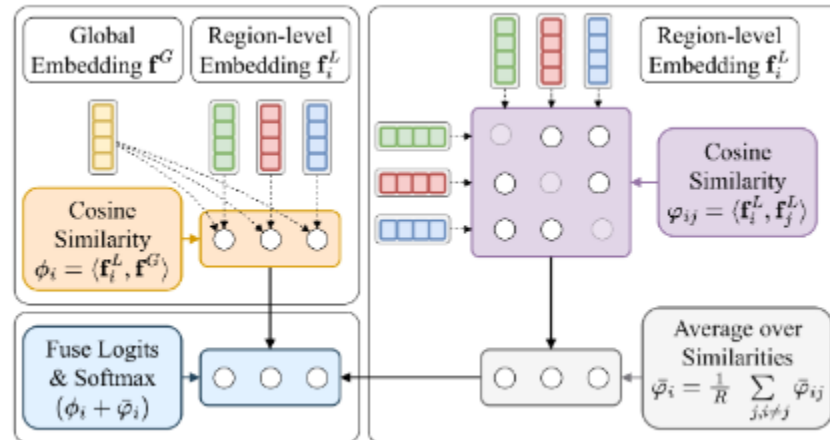


ConceptFusion: Open-set multimodal 3d mapping

Construct pixel-aligned features



Zero Shot Pixel Fusion



04

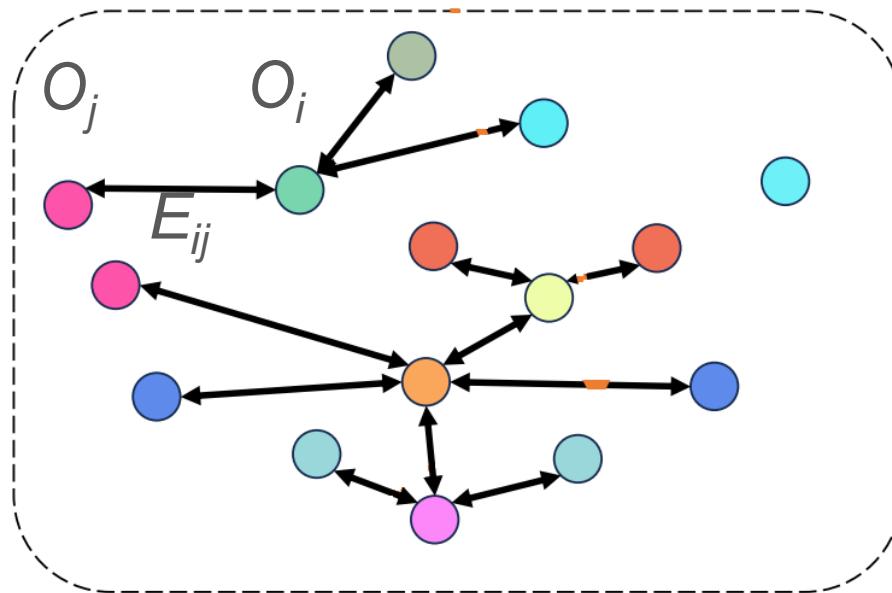
Разреженные (Sparse) методы

Постановка задачи

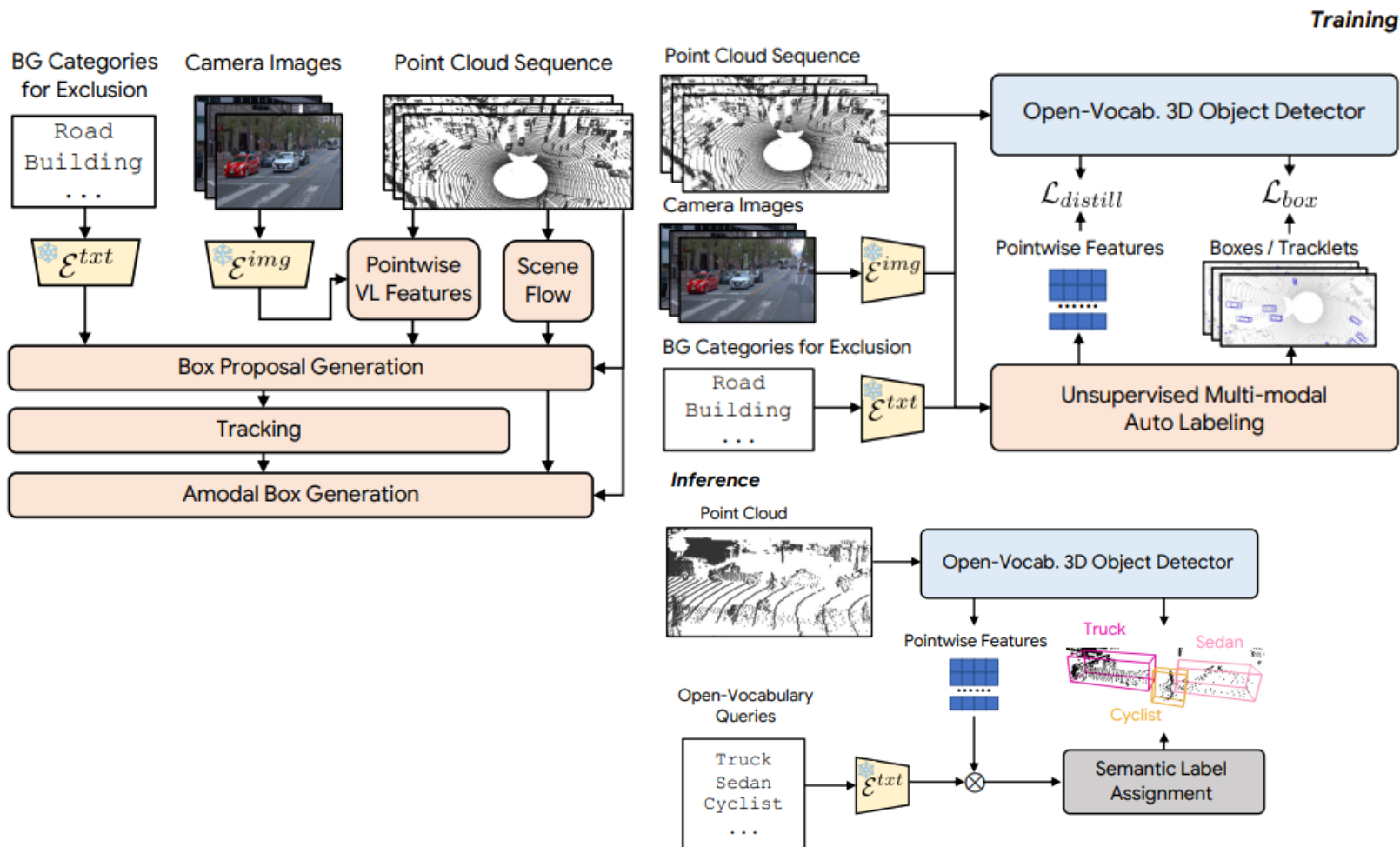
I: $Map = M((RGB_1 \dots RGB_N), (D_1 \dots D_N), (Pose_1 \dots Pose_N)) = G(O, E),$

II: $O_{Query} = Q(Map, Query).$

$O_i = (X_{ic}, Y_{ic}, Y_{ic}, Embedding_i, PCL_i, Caption_i, Tag_i).$



Unsupervised 3D Perception with 2D Vision-Language Distillation for Autonomous Driving



Algorithm 1 Unsupervised multi-modal auto labeling.

Input: A sequence of images across T frames for each of the K cameras $\{\mathbf{I}_t^k\}$; a sequence of LiDAR point locations $\{\mathbf{P}_t\}$.

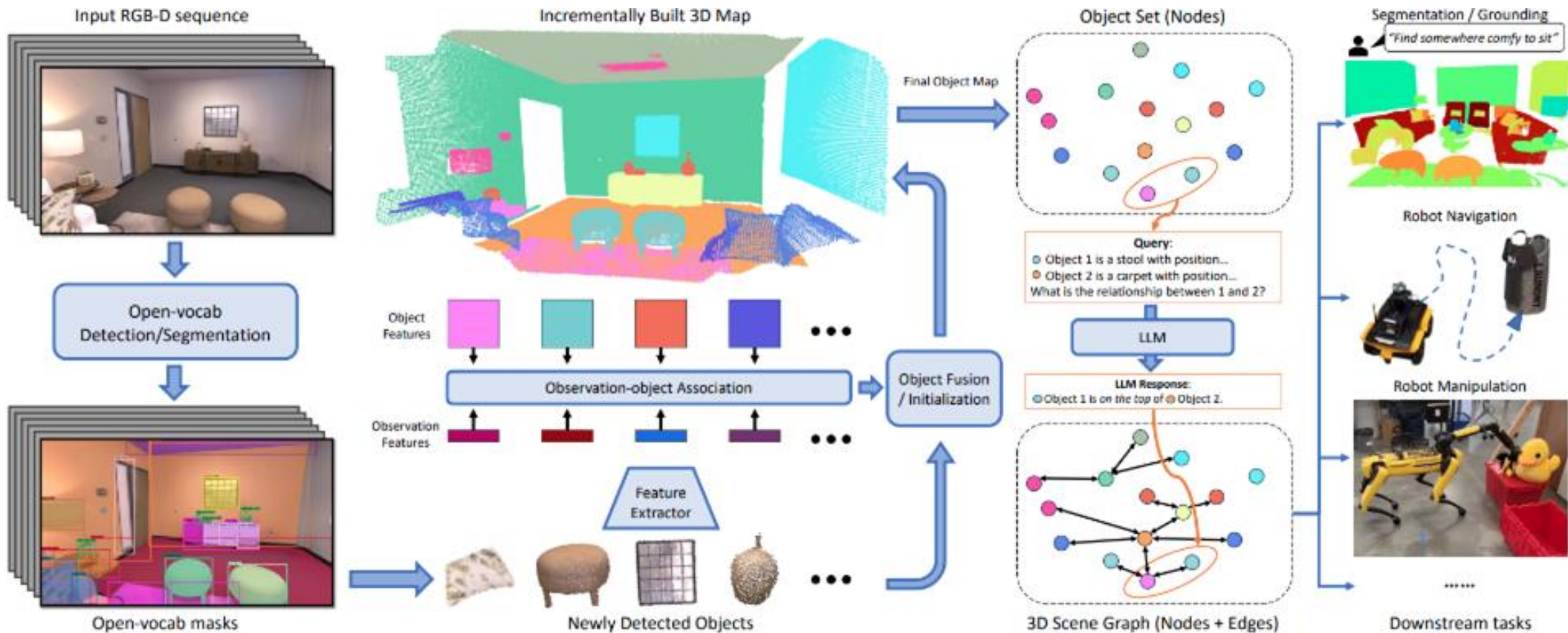
Requires: Cosine similarity threshold for background categories ϵ^{bg} ; minimum scene flow magnitude ϵ^{sf} ; maximum ratio of background points within a box r^{bg} ; a set of a prior background categories \mathbf{C}^{bg} ; a pre-trained open-vocabulary model with image encoder \mathcal{E}^{img} and text encoder \mathcal{E}^{txt} .

Output: Amodal 3D bounding boxes $\{\mathbf{B}_t\}$ and their track IDs $\{\mathbf{T}_t\}$; point-wise open-vocabulary features $\{\mathbf{F}_t^{vl}\}$.

Function:

- 1: **for** $t = 1$ to T **do**
- 2: $\{\mathbf{V}_t^k\} \leftarrow \mathcal{E}^{img}(\{\mathbf{I}_t^k\})$ \triangleright 2D VL features
- 3: $\mathbf{F}_t^{vl} \leftarrow \text{Unprojection}(\{\mathbf{V}_t^k\}, \mathbf{P}_t)$ \triangleright 3D VL features
- 4: **if** $t \neq T$ **then**
- 5: $\mathbf{F}_t^{sf} \leftarrow \text{NSFP++}(\mathbf{P}_t, \mathbf{P}_{t+1})$ \triangleright Scene flow
- 6: **else**
- 7: $\mathbf{F}_t^{sf} \leftarrow -\text{NSFP++}(\mathbf{P}_t, \mathbf{P}_{t-1})$
- 8: **for** $i = 1$ to N_t **do**
- 9: $(\mathbf{M}_t^{sf})_i \leftarrow \mathbb{1}(\|\mathbf{F}_t^{sf}\|_i \geq \epsilon^{sf})$
- 10: $(\mathbf{M}_t^{bg})_i \leftarrow \mathbb{1}(\max_{c \in \mathbf{C}^{bg}} \frac{(\mathbf{F}_t^{vl})_i \cdot \mathcal{E}^{txt}(c)}{\|\mathbf{F}_t^{vl}\|_i \|\mathcal{E}^{txt}(c)\|} \geq \epsilon^{bg})$
- 11: $\tilde{\mathbf{P}}_t, \tilde{\mathbf{F}}_t^{sf} \leftarrow \mathbf{P}_t[\mathbf{M}_t^{sf}], \mathbf{F}_t^{sf}[\mathbf{M}_t^{sf}]$
- 12: $\mathbf{B}_t^{vis} \leftarrow \text{InitialBoxProposal}(\tilde{\mathbf{P}}_t, \tilde{\mathbf{F}}_t^{sf}, \mathbf{M}_t^{bg}, r^{bg})$
- 13: $\{\mathbf{T}_t\} \leftarrow \text{Tracking}(\{\mathbf{B}_t^{vis}\})$
- 14: $\{\mathbf{B}_t\} \leftarrow \text{AmodalBoxGeneration}(\{\mathbf{B}_t^{vis}\}, \{\mathbf{T}_t\}, \{\mathbf{P}_t\})$
- 15: **return** $\{\mathbf{B}_t\}, \{\mathbf{T}_t\}, \{\mathbf{F}_t^{vl}\}$

ConceptGraphs: Open-Vocabulary 3D Scene Graphs for Perception and Planning



Project
Github

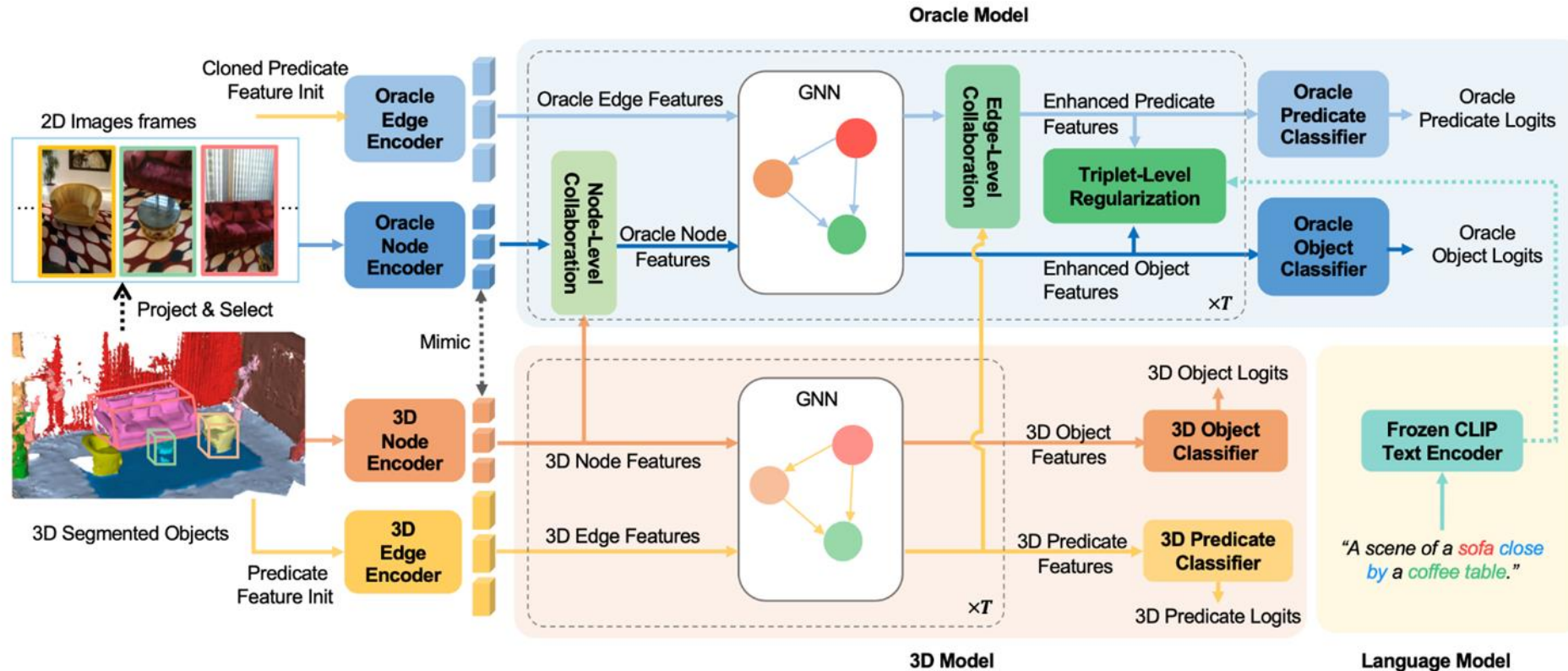
VL-SAT: Visual-Linguistic Semantics Assisted Training for 3D Semantic Scene Graph Prediction

Dataset used: 3DSSG, (scenes from 3RSCAN)

- 1553 scenes
- 160 cat. objects
- 26 cat. predicates

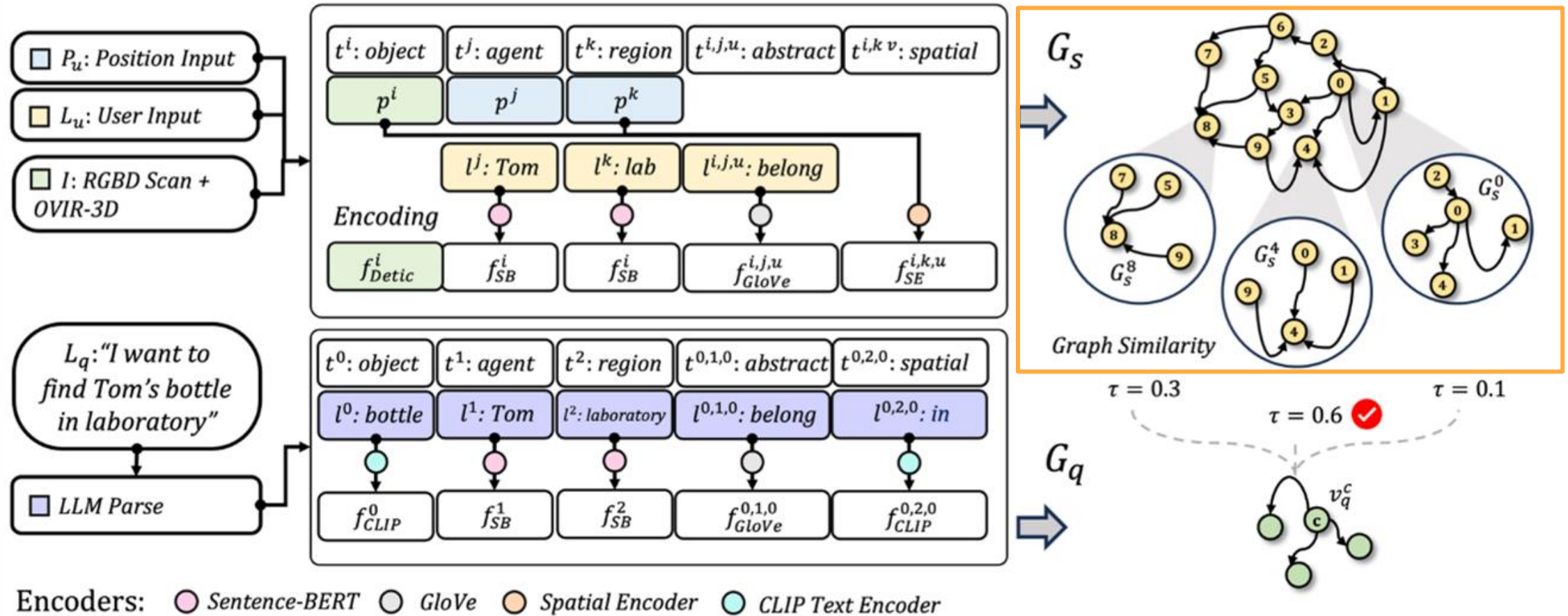
Target:

- build an accurate 3D scene graph, i.e.:
- classify objects
- classify relationships between objects



Wang, Z., Cheng, B., Zhao, L., Xu, D., Tang, Y., & Sheng, L. (2023). VL-SAT: Visual-Linguistic Semantics Assisted Training for 3D Semantic Scene Graph Prediction in Point Cloud. In Proceedings of the IEEE/CVF Conference on Computer

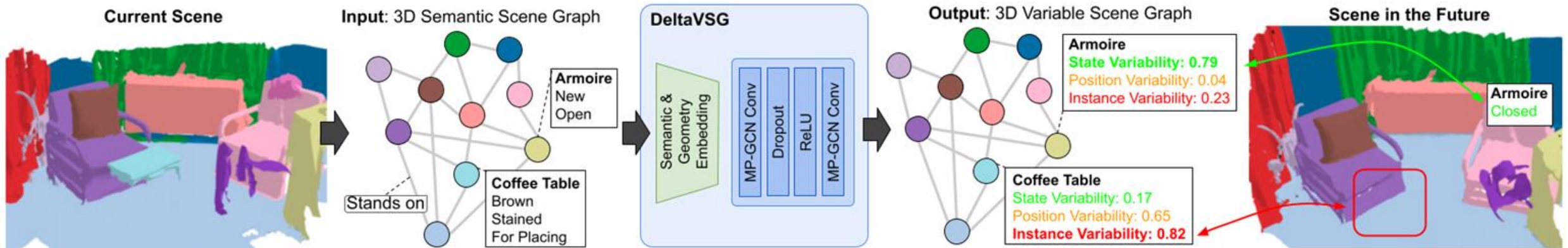
Context-Aware Entity Grounding with Open-Vocabulary 3D Scene Graphs built once



Chang, H., Boyalakuntla, K., Lu, S., Cai, S., Jing, E., Keskar, S., ... & Boularias, A. (2023). Context-aware entity grounding with open-vocabulary 3d scene graphs. *arXiv preprint arXiv:2309.15940*.

<https://github.com/changhaonan/ovsg>

3D Variable Scene Graphs



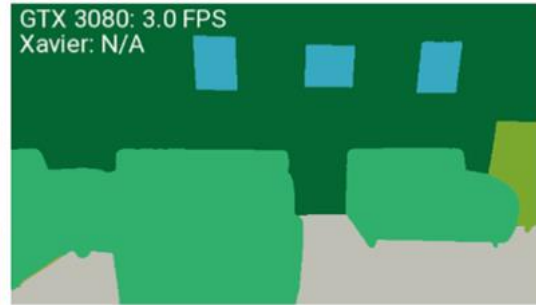
https://github.com/ethz-asl/3d_vsg

Looper, S., Rodriguez-Puigvert, J., Siegart, R., Cadena, C., & Schmid, L. (2023, May). 3d vsg: Long-term semantic scene change prediction through 3d variable scene graphs. In *2023 IEEE International Conference on Robotics and Automation (ICRA)* (pp. 8179-8186). IEEE.

Foundations of Spatial Perception for Robotics: Hierarchical Representations and Real-time Systems



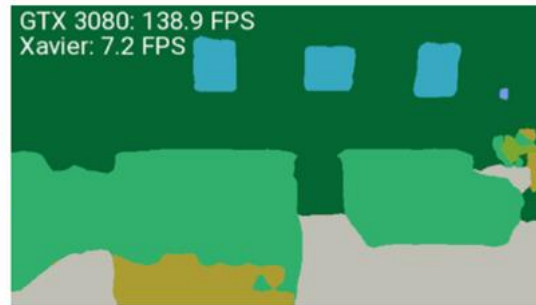
(a) Original Image



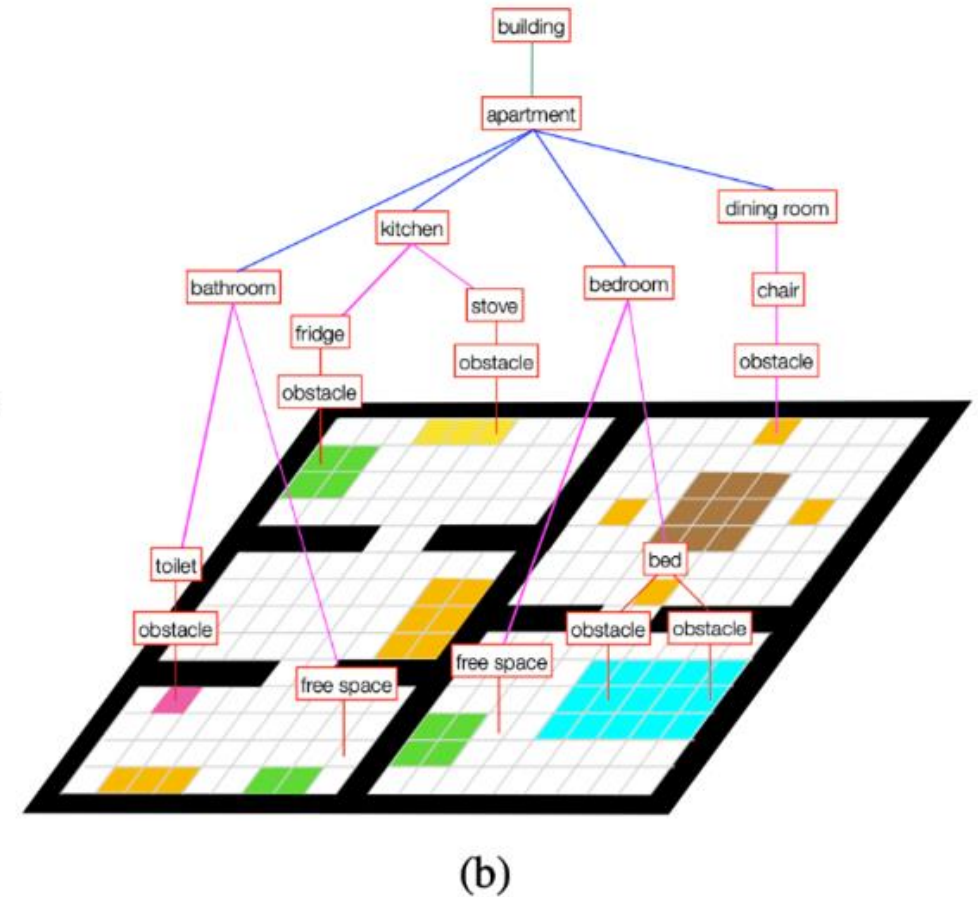
(b) OneFormer



(c) HRNet



(d) MobileNet

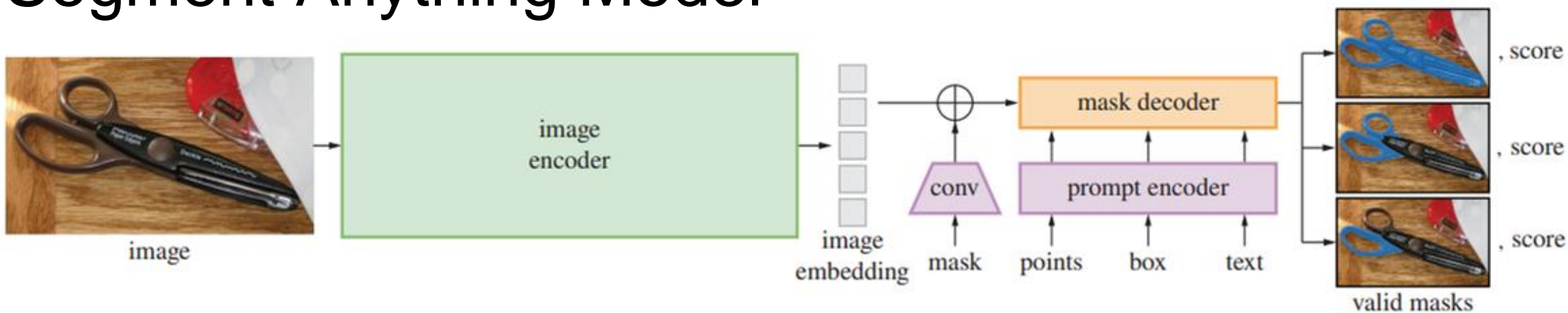


Hughes, N., Chang, Y., Hu, S., Talak, R., Abdulhai, R., Strader, J., & Carlone, L. (2023). Foundations of Spatial Perception for Robotics: Hierarchical Representations and Real-time Systems. *arXiv preprint arXiv:2305.07154*.

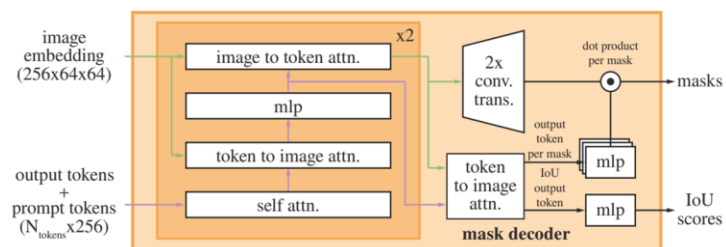
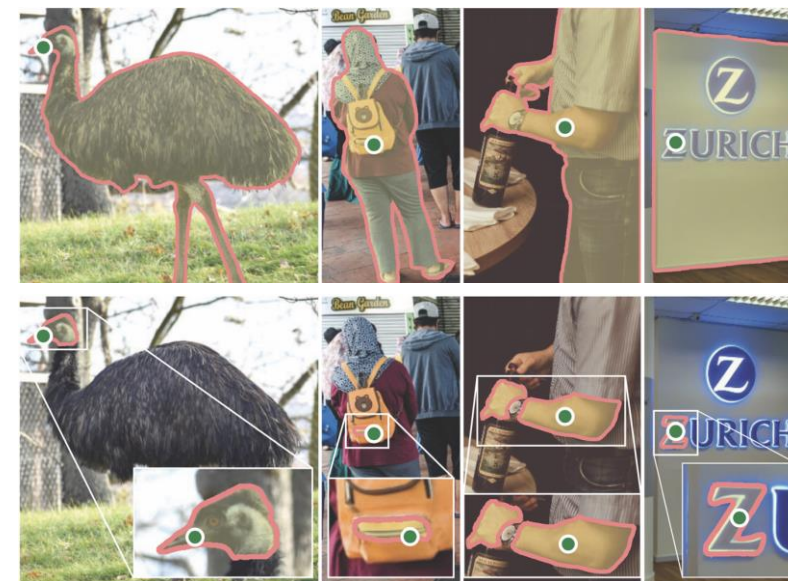
05



Базовые модели сегментации изображений с открытым словарем запросов (open vocabulary)

Segment Anything Model

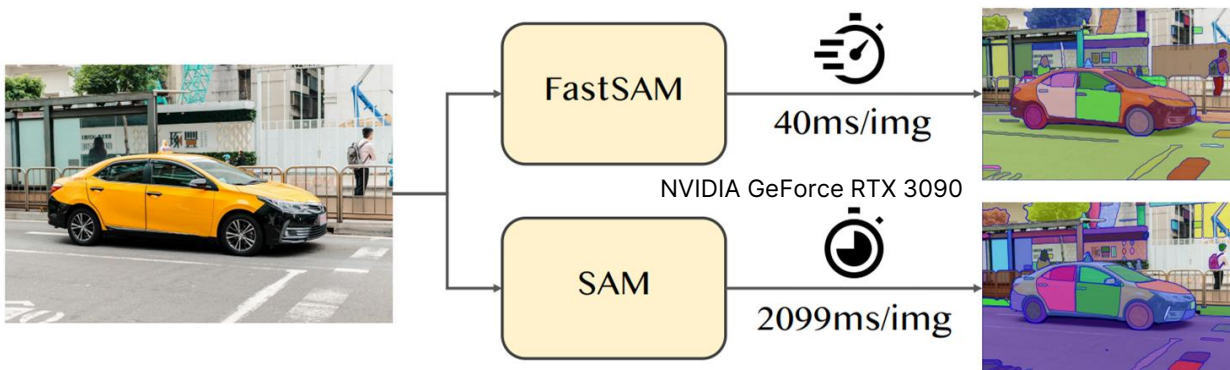


- Image encoder: Vision Transformer (ViT) ViT-H, ViT-L, ViT-B
- Prompt encoder: positional encoding for points and box corners, CLIP-based text encoder and convolutional mask encoder
- Mask decoder: lightweight attention-based model

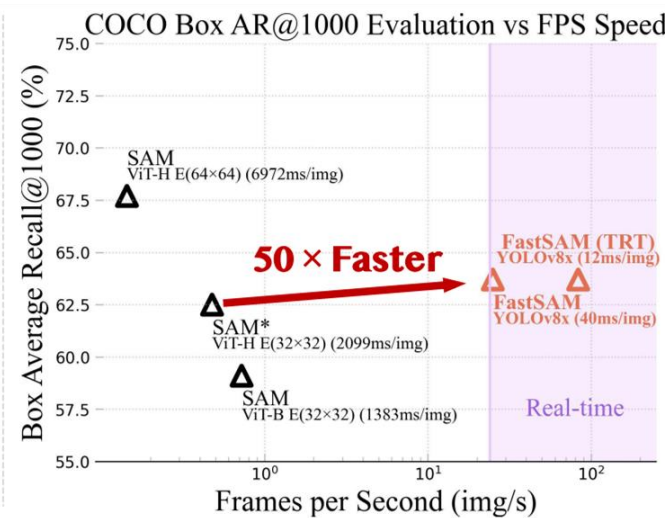
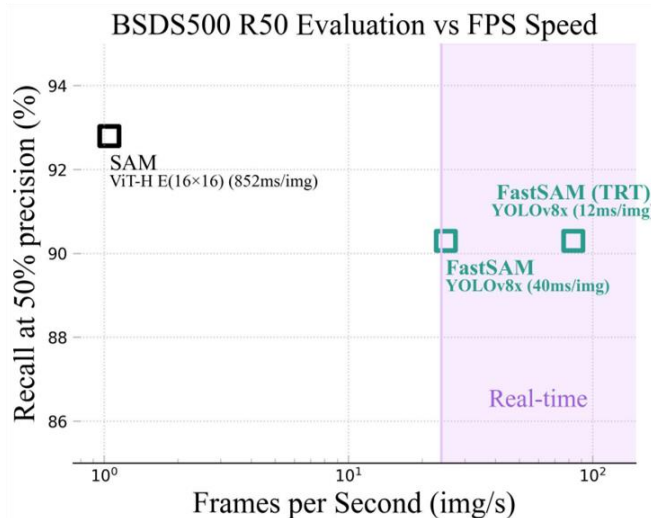


 **Project**
Github 

Fast query-based segmentation: FastSAM



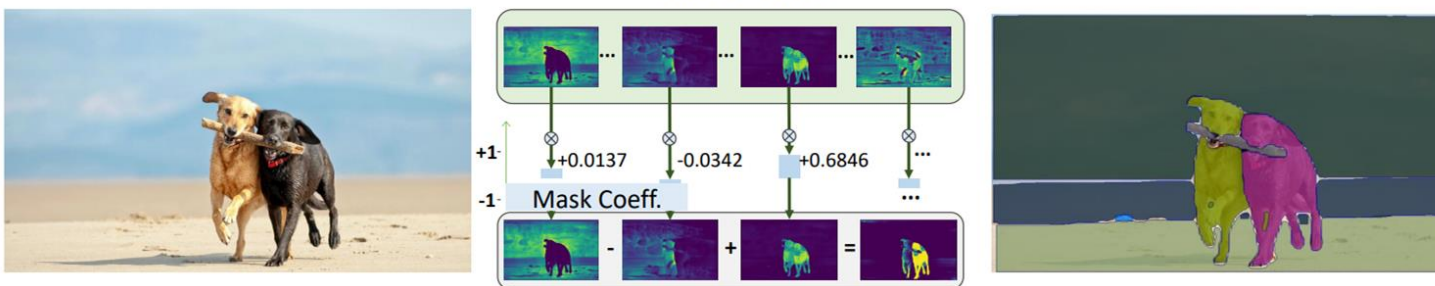
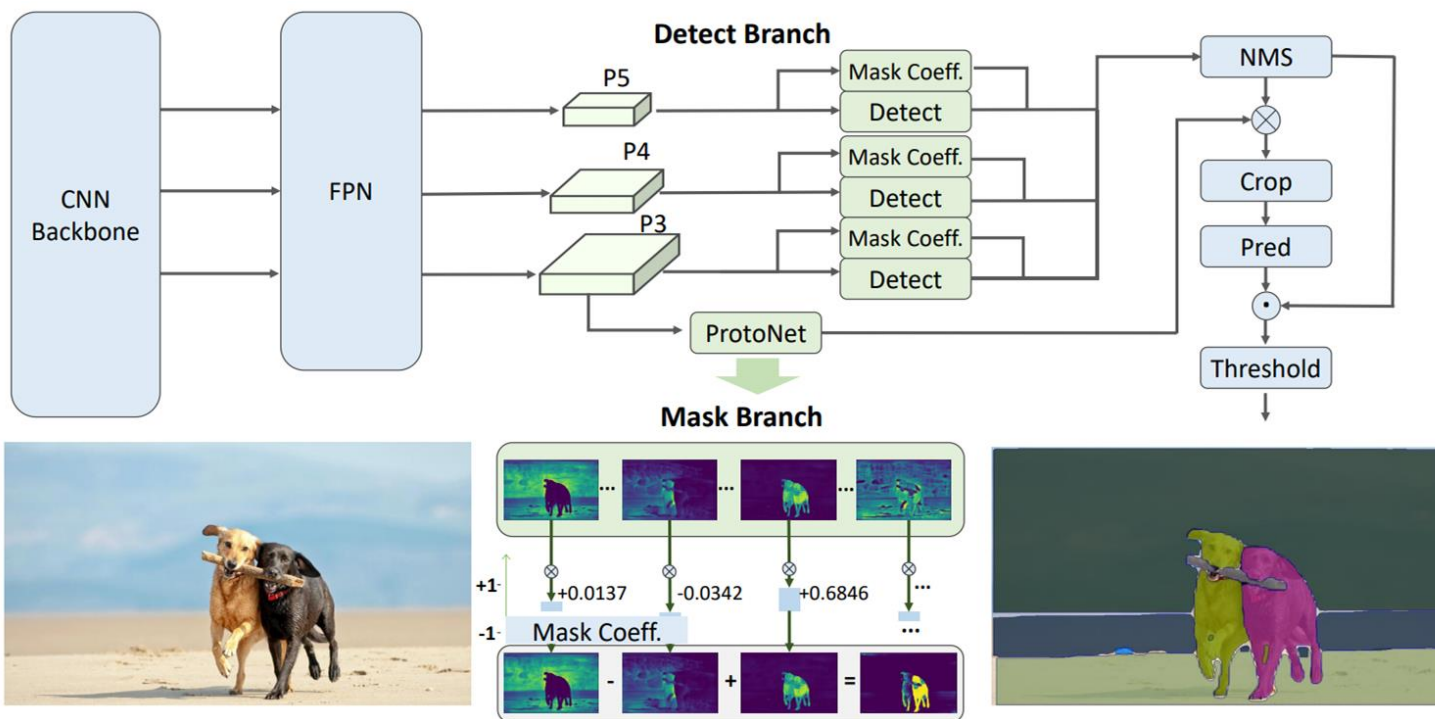
- Authors trained the existing real-time instance segmentation method YOLOv8-seg using only 1/50 of the SA-1B dataset published by SAM



Project

Github

Fast query-based segmentation: FastSAM

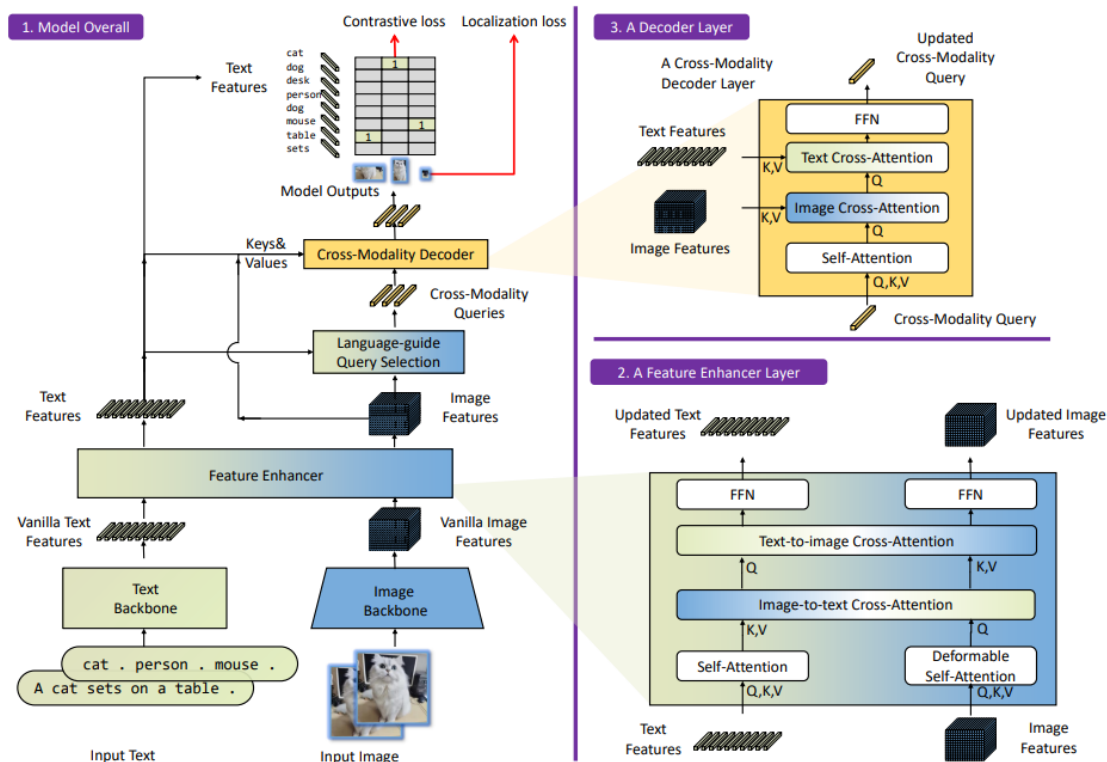


- The framework of FastSAM. It contains two stages: All-instance Segmentation (AIS) and Prompt-guided Selection (PGS).
- Authors use YOLOv8-seg to segment all objects or regions in an image.
- Then authors use various prompts to identify the specific object(s) of interest. It mainly involves the utilization of point prompts, box prompts, and text prompt. The text prompt is based on CLIP



Fast query-based segmentation: Grounded Mobile SAM

Grounding DINO [1]



Mobile SAM [2]

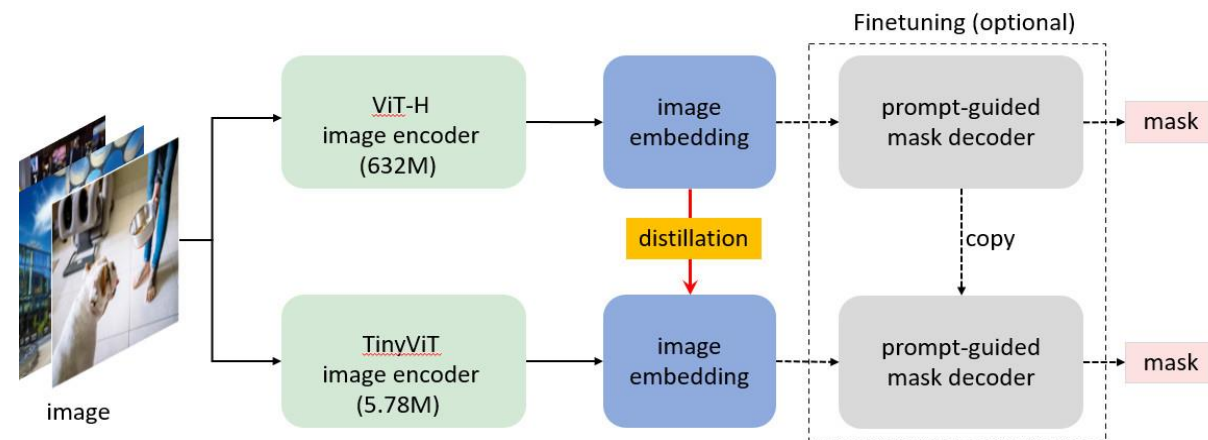


Table 6: Comparison between FastSAM and MobileSAM.

	FastSAM	MobileSAM	Ratio
Size	68M	9.66M	7
Speed	40ms	10ms	4

Table 7: mIoU comparison. With the assumption that the predicted mask from the original SAM is ground-truth, a higher mIoU indicates a better performance.


	100	200	300	400	500
FastSAM	0.27	0.33	0.37	0.41	0.41
MobileSAM	0.73	0.71	0.74	0.73	0.73

[1] Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., ... & Zhang, L. (2023). Grounding dino: Marrying dino with grounded pre-training for open-set object detection. arXiv preprint arXiv:2303.05499. (Tsinghua University, IDEA, The Hong Kong University of Science and Technology, The Chinese University of Hong Kong (Shenzhen), Microsoft Research, Redmond)

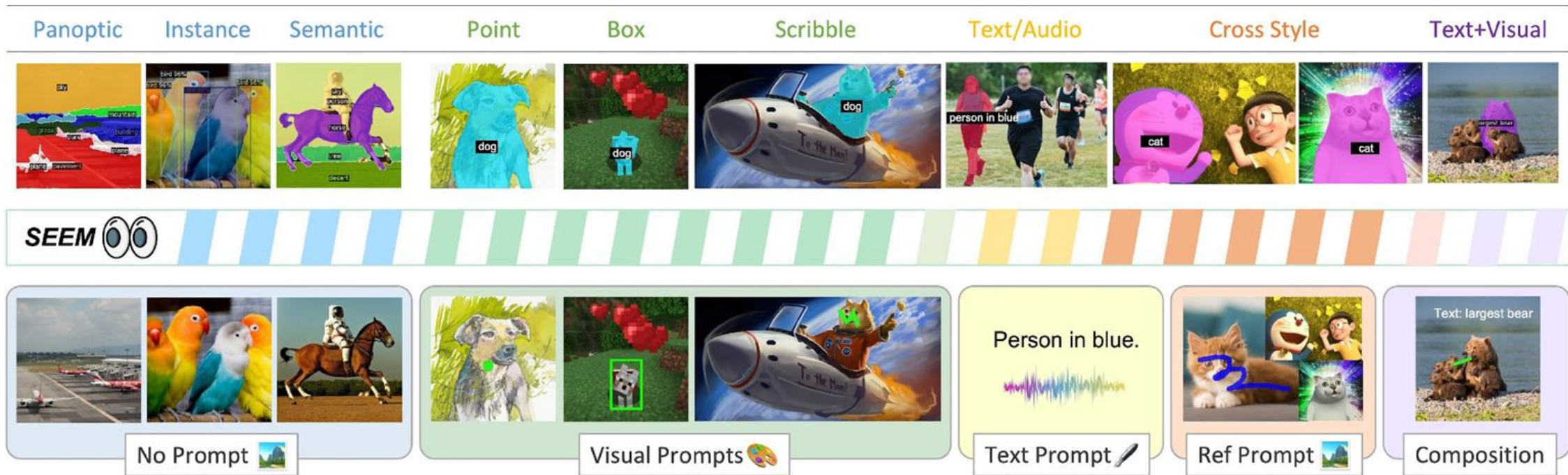
[2] Zhang, C., Han, D., Qiao, Y., Kim, J. U., Bae, S. H., Lee, S., & Hong, C. S. (2023). Faster Segment Anything: Towards Lightweight SAM for Mobile Applications. arXiv preprint arXiv:2306.14289. (Kyung Hee University)



Project

Github 


SEEM: Segment Everything Everywhere All at Once



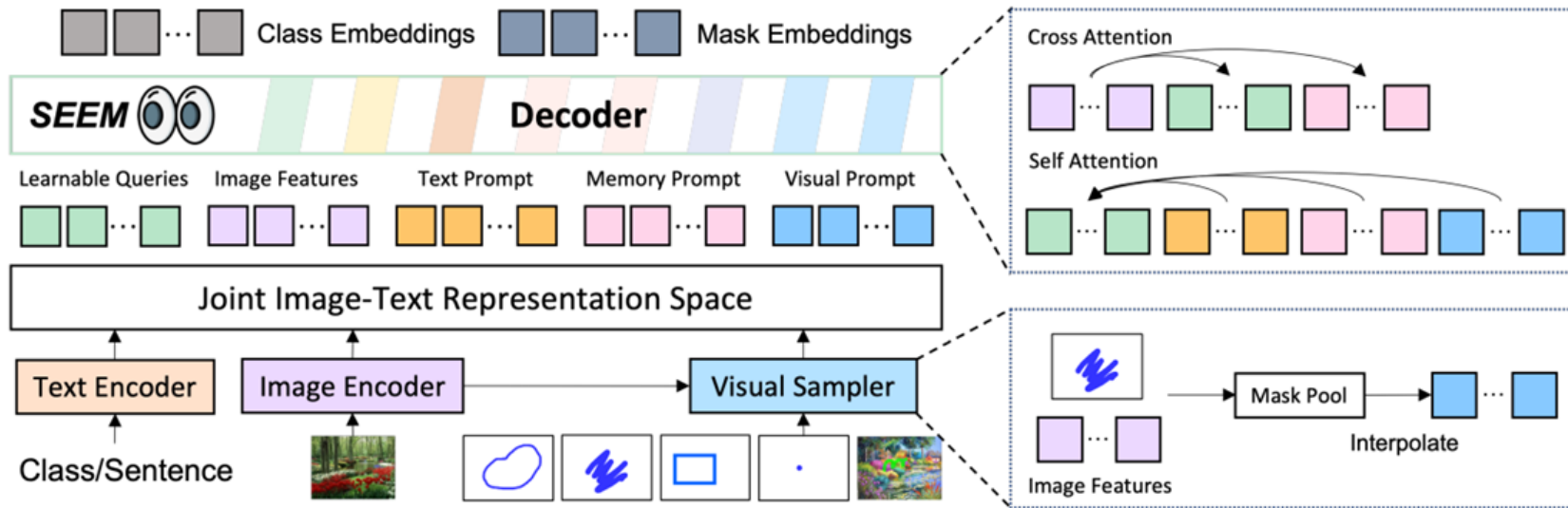
- Authors design a new prompting scheme that can encode various user intents into prompts in a joint visual-semantic space, enabling strong flexibility for various segmentation tasks and generalization capability to unseen prompts or their combinations.
- Authors build SEEM, a universal and interactive segmentation interface that integrates the newly designed prompting mechanism into a lightweight decoder for all segmentation tasks, leading to a model possessing properties of versatility, compositionality, interactivity, and semantic awareness



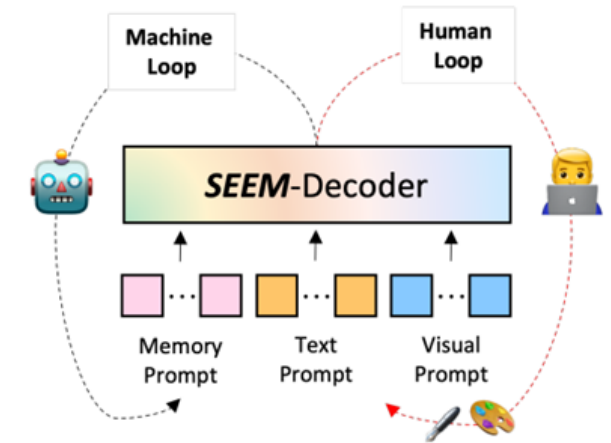
Project

Github 

SEEM-Decoder



(a) Model Architecture




(b) Human-Model Interaction

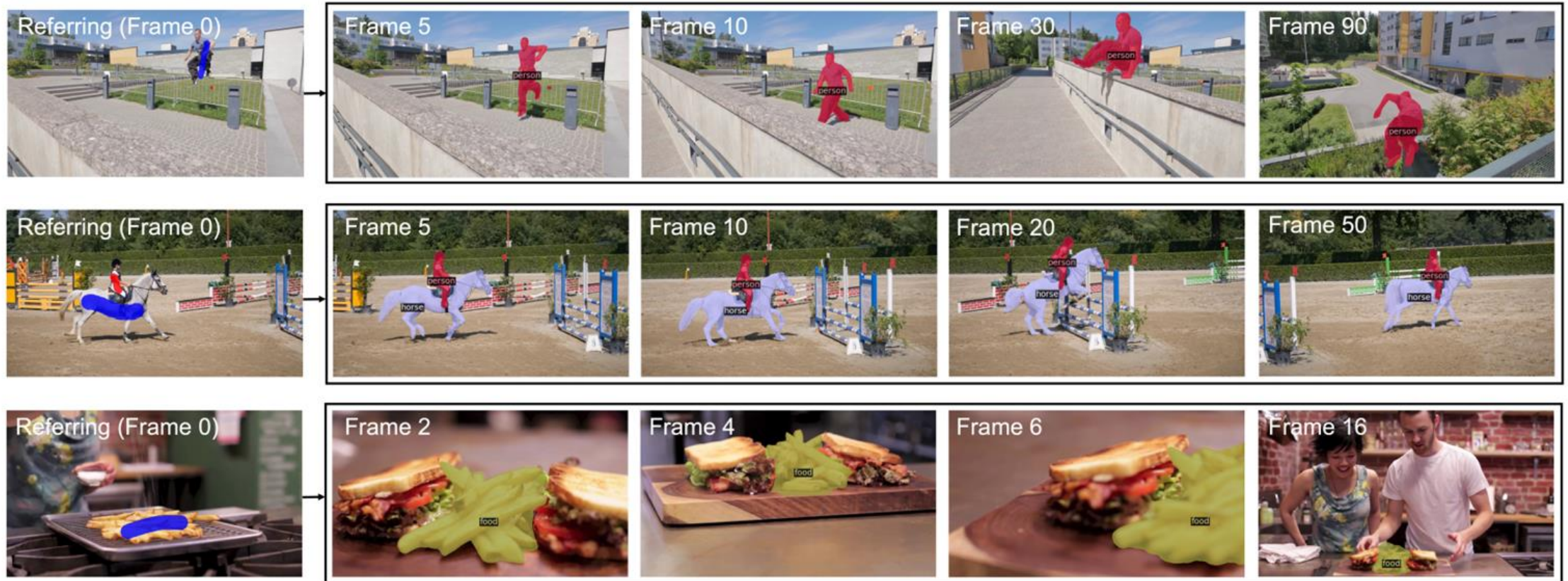
- SEEM encodes image, text, and human inputs into joint visual-semantic space as queries, features, and prompts, and then decodes queries to class and mask embeddings.
- With the benefit of SEEM decoder, the machine loop enables memorizing history mask information, and the human loop provides new corrections to the next round.
- For the vision backbone, authors use FocalT, DaViT-d3 (B), and DaViT-d5 (L). For the language encoder, they adopt a UniCL or Florence text encoder. For interactive segmentation, after one click on the image to generate the predicted mask, the next click is placed at the center of the area with the largest segmentation error



Project

Github 

SEEM: Segment Everything Everywhere All at

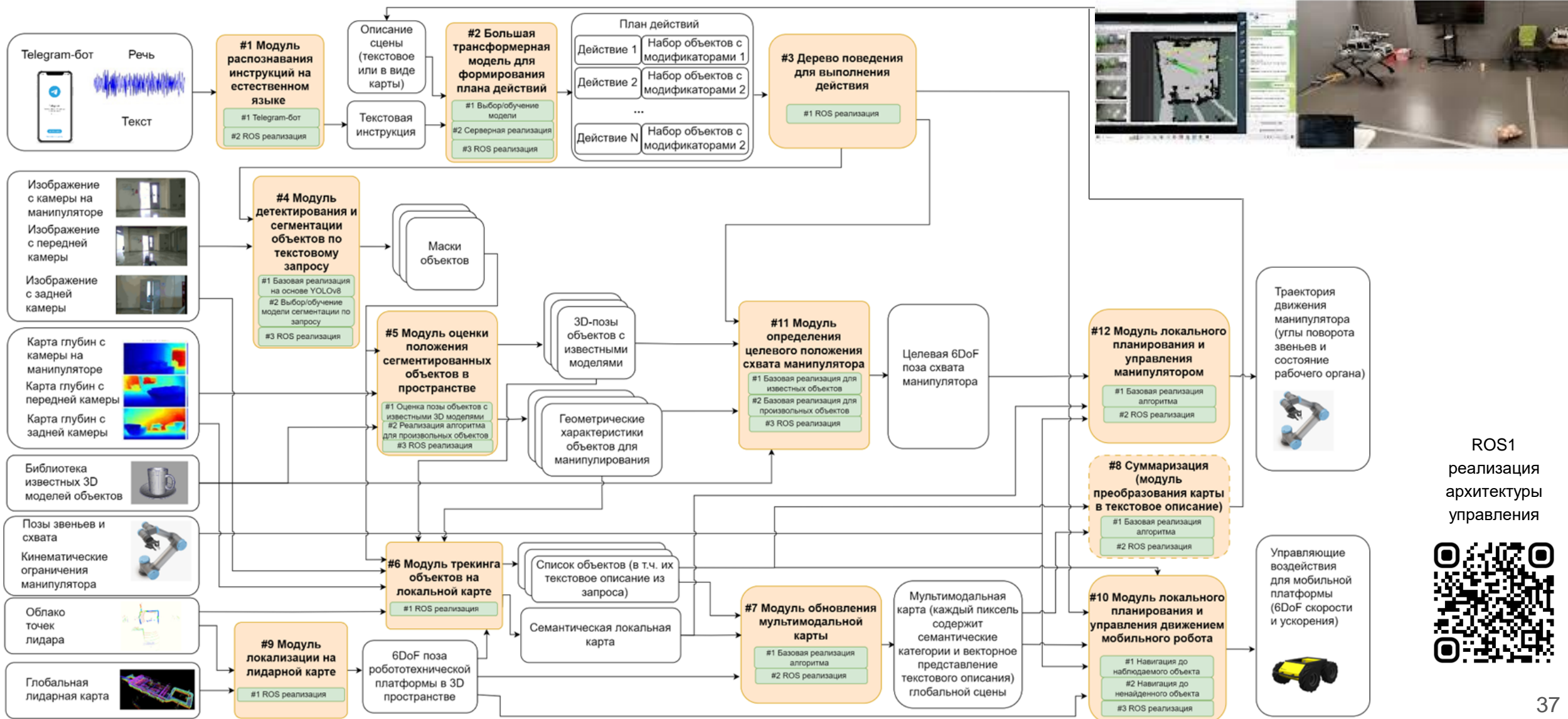


- Zero-shot video object segmentation using the first frame plus one stroke. From top to bottom, the videos are “parkour” and “horsejump-low” from DAVIS, and video 101 from YouCook2.
- SEEM precisely segments referred objects even with significant appearance changes caused by blurring or intensive deformations.

06

Наши исследования
и приложения

STRL-Robotics-LLM - Модульная система управления мобильным манипулятором

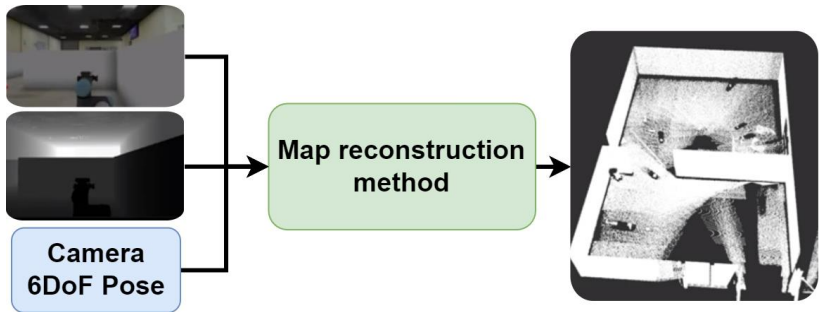


ROS1 реализация архитектуры управления

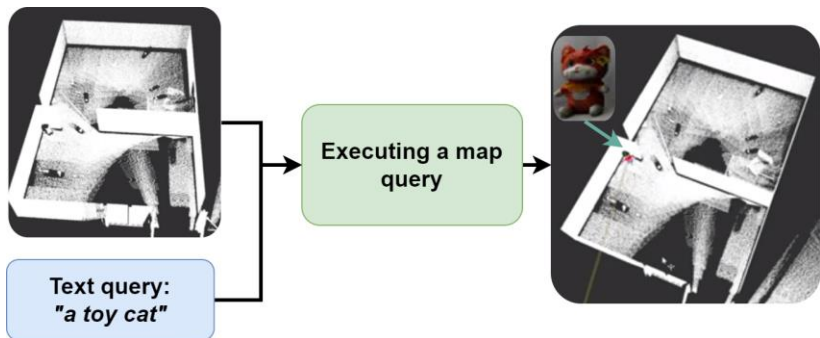


Integration with Robotic Operation System (ROS)

1st stage: map reconstruction



2nd stage: text query on the map



The proposed method has 2 stages:

- Create a map;
- Create a query on the map.

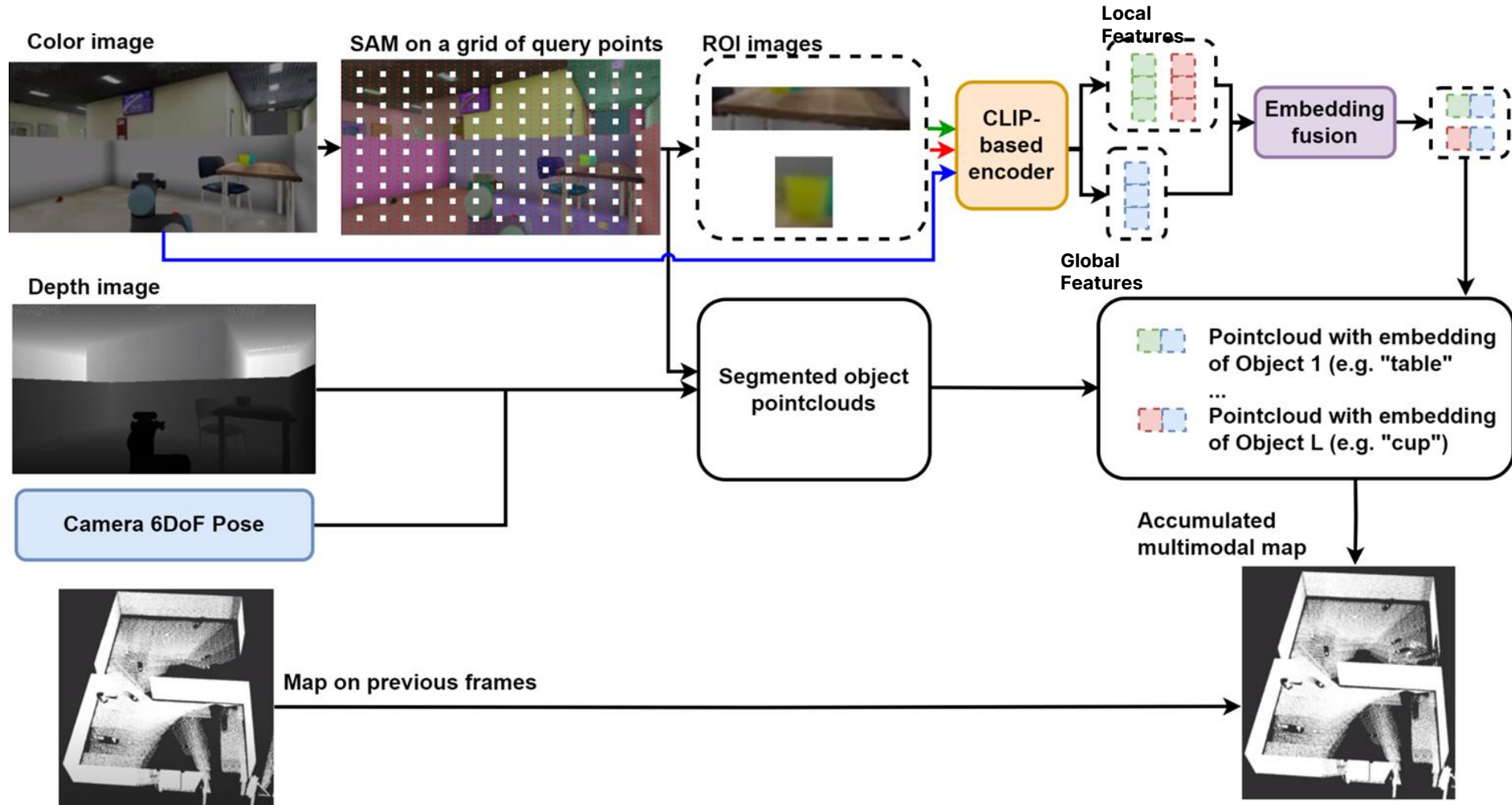


Nvidia Isaac Sim Scene

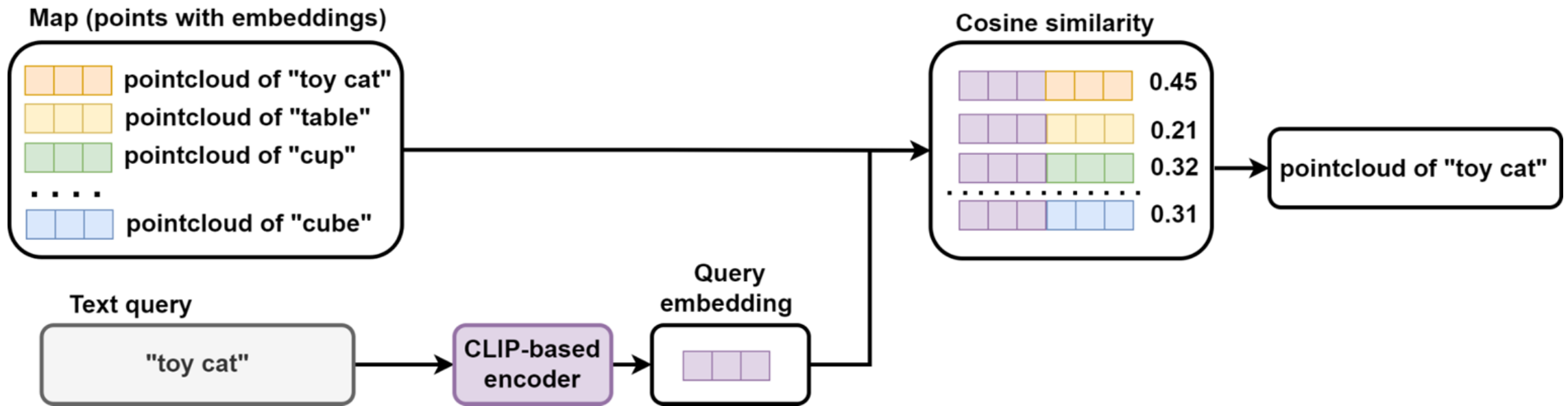


Real scene with Husky+UR5

Multimodal map reconstruction



Executing a map query

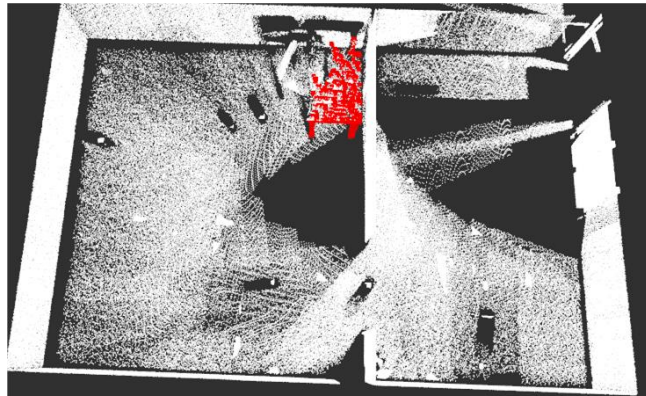


Validation of pre-trained models

Query: a table



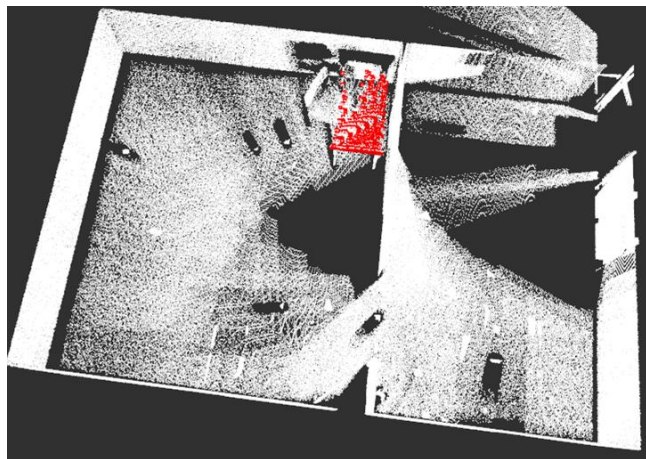
SAM
CLIP(ViT-H)



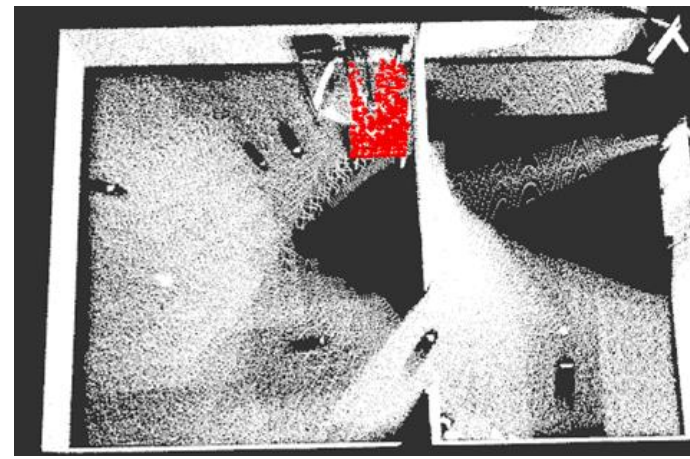
MobileSAM
CLIP(ViT-H)



SAM
CLIP(ConvNext)



MobileSAM
CLIP(ConvNext)

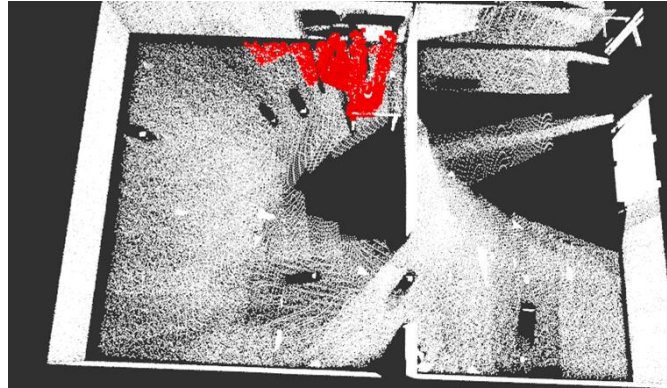


Validation of pre-trained models

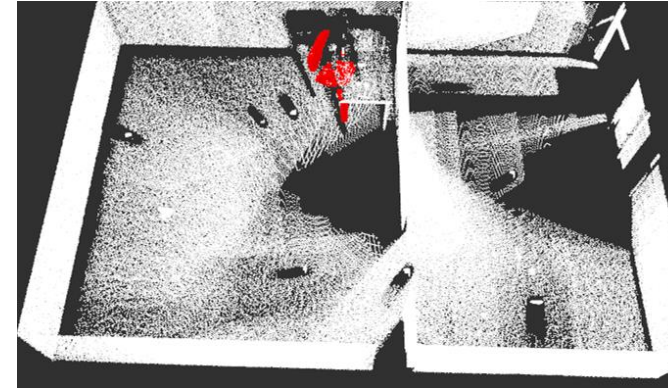
Query: a chair



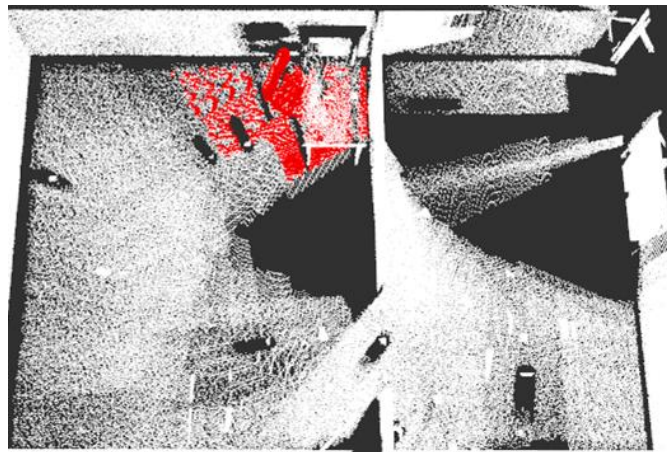
SAM
CLIP(ViT-H)



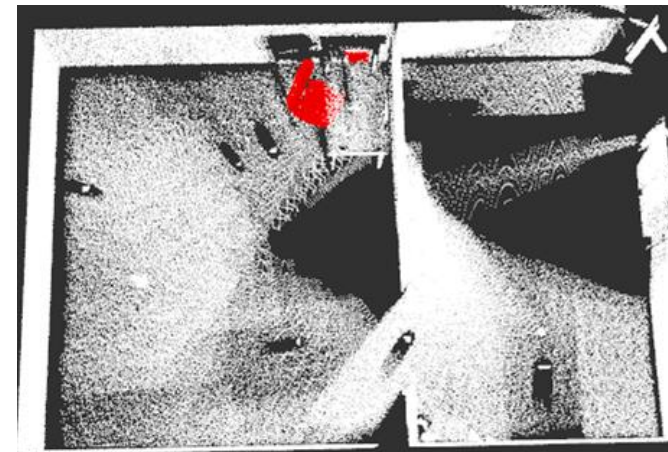
MobileSAM
CLIP(ViT-H)



SAM
CLIP(ConvNext)

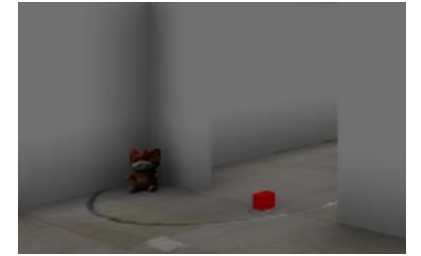


MobileSAM
CLIP(ConvNext)

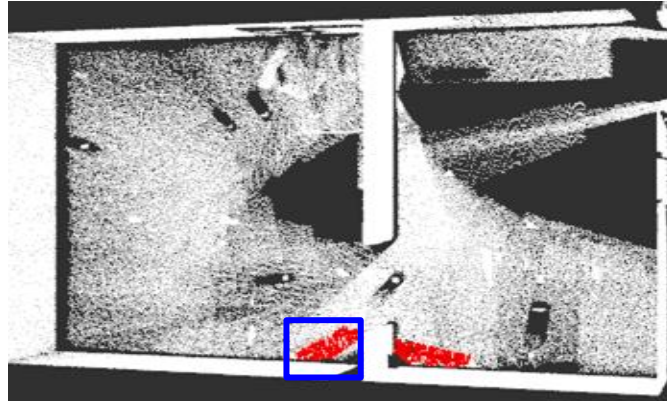


Validation of pre-trained models

Query: a toy cat



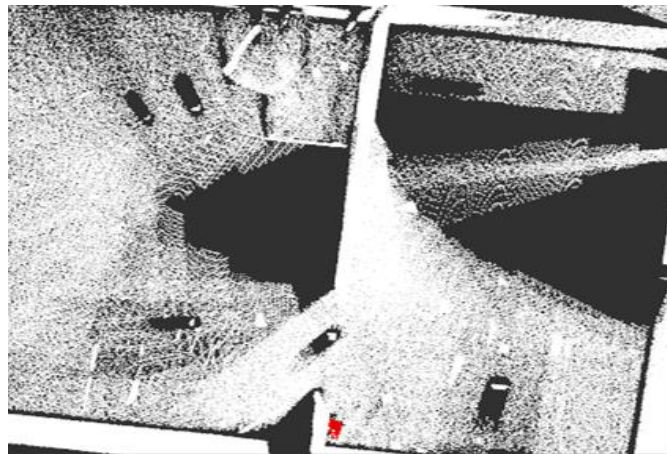
SAM
CLIP(ViT-H)



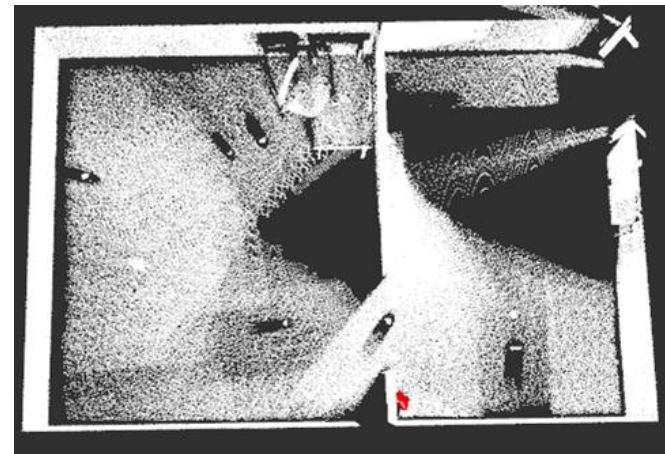
MobileSAM
CLIP(ViT-H)



SAM
CLIP(ConvNext)



MobileSAM
CLIP(ConvNext)

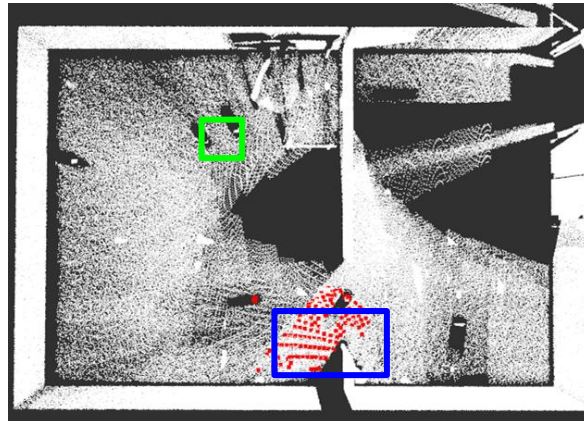


Validation of pre-trained models

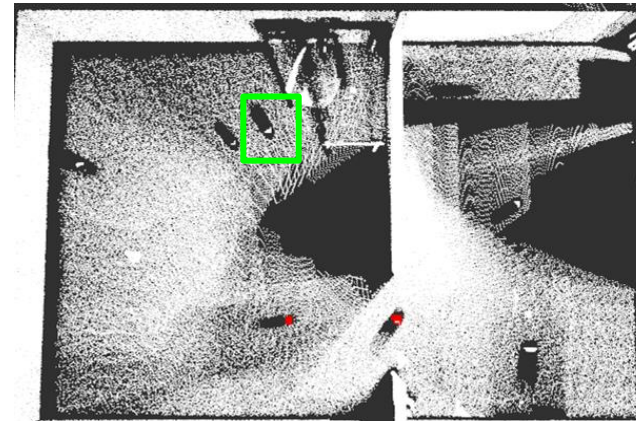
Query: a red cube



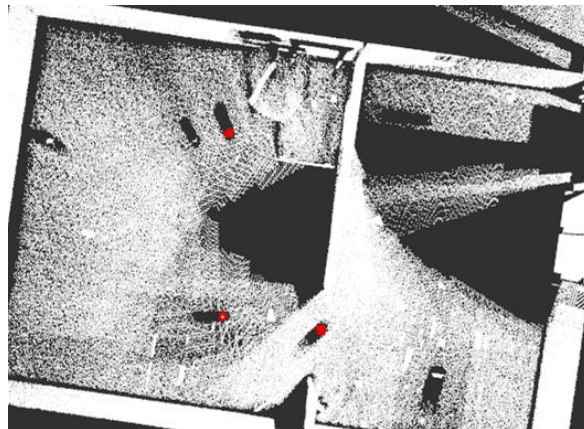
SAM
CLIP(ViT-H)



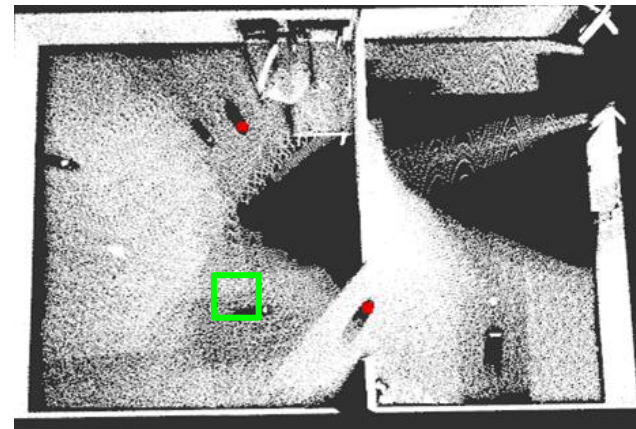
MobileSAM
CLIP(ViT-H)



SAM
CLIP(ConvNext)



MobileSAM
CLIP(ConvNext)

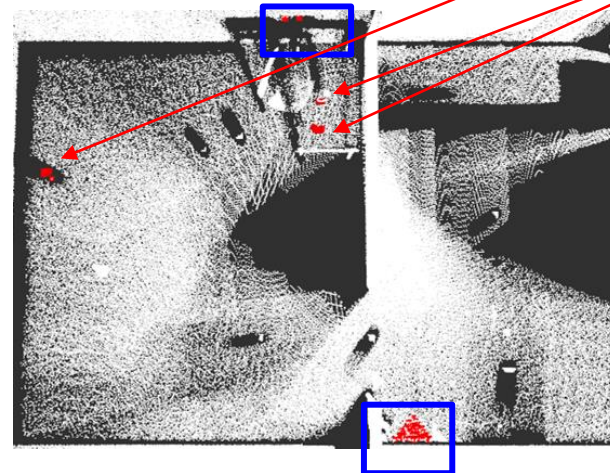
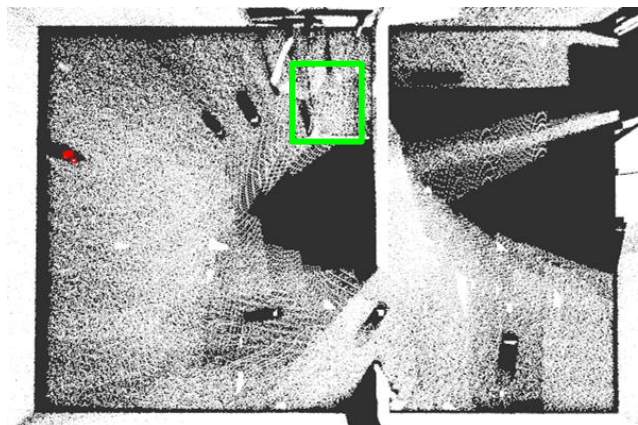


Validation of pre-trained models

Query: a cup

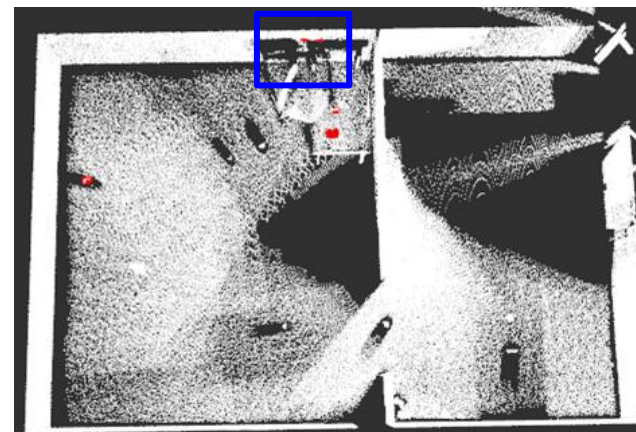
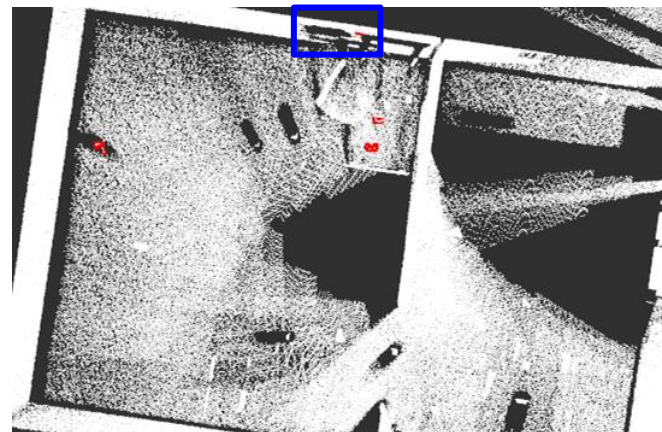


SAM
CLIP(ViT-H)



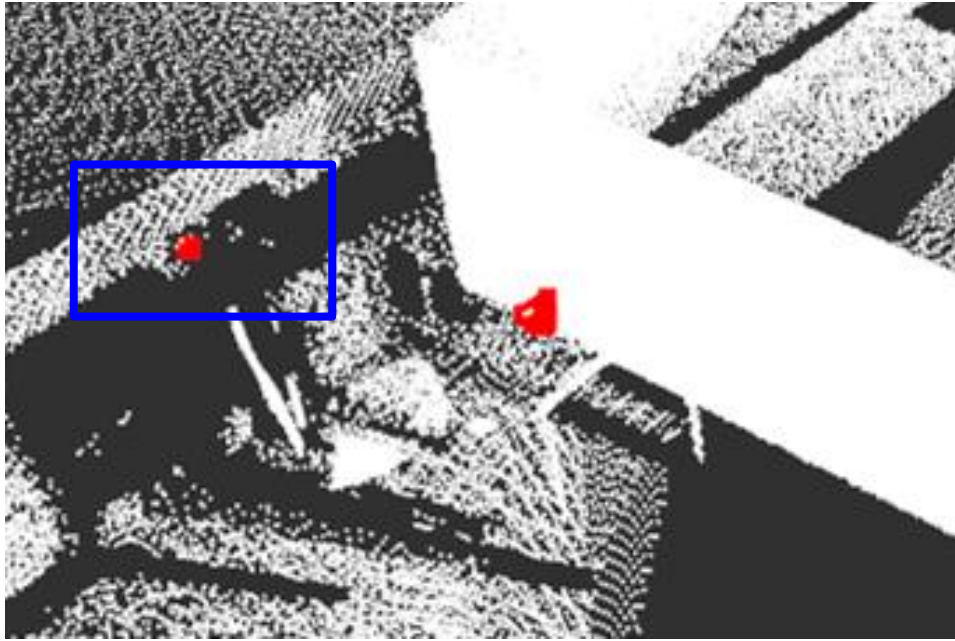
MobileSAM
CLIP(ViT-H)

SAM
CLIP(ConvNext)

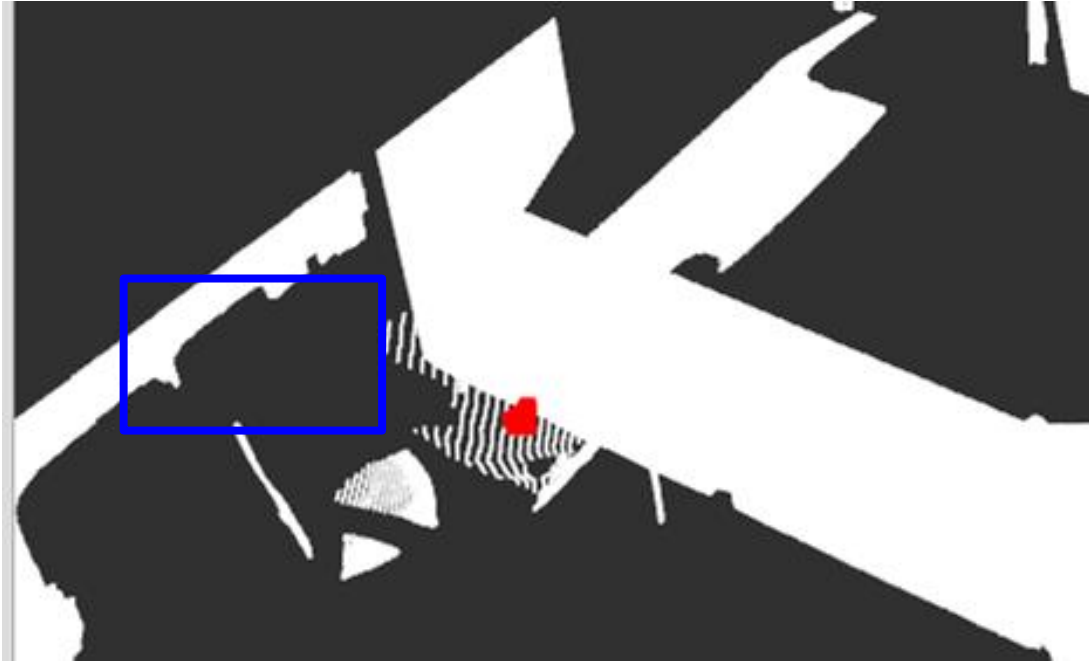


MobileSAM
CLIP(ConvNext)

Filtering the background on an object's point cloud



MobileSAM
CLIP(ConvNext) w/o filtering



MobileSAM
CLIP(ConvNext) + filtering

Performance of Multimodal Mapping Implementation for ROS

SAM Model***	CLIP Model	Time per frame in SAM, s **	Time per frame in CLIP on creating map, s **	Time in CLIP on query, s **	Memory (Reconstruct), MiB	Memory (Query), MiB
SAM ViT-H*	ViT-H-14 LAION 2B*	4.18	2.45	0.022	10'391	4'916
MobileSAM	ViT-H-14 LAION 2B	2.52 (-39.7%)	1.92(-21.6%)	0.022(-0%)	8'611(-17.1%)	4'916(-0%)
SAM ViT-H	ConvNext Base LAION 2B	4.18 (-0%)	0.73(-70.2%)	0.013(-40.9%)	7'226(-30.4%)	1'694(-65.5%)
MobileSAM	ConvNext Base LAION 2B	2.51 (-39.9%)	0.59(-75.9%)	0.013(-40.9%)	5'461(-47.4%)	1'694(-65.5%)

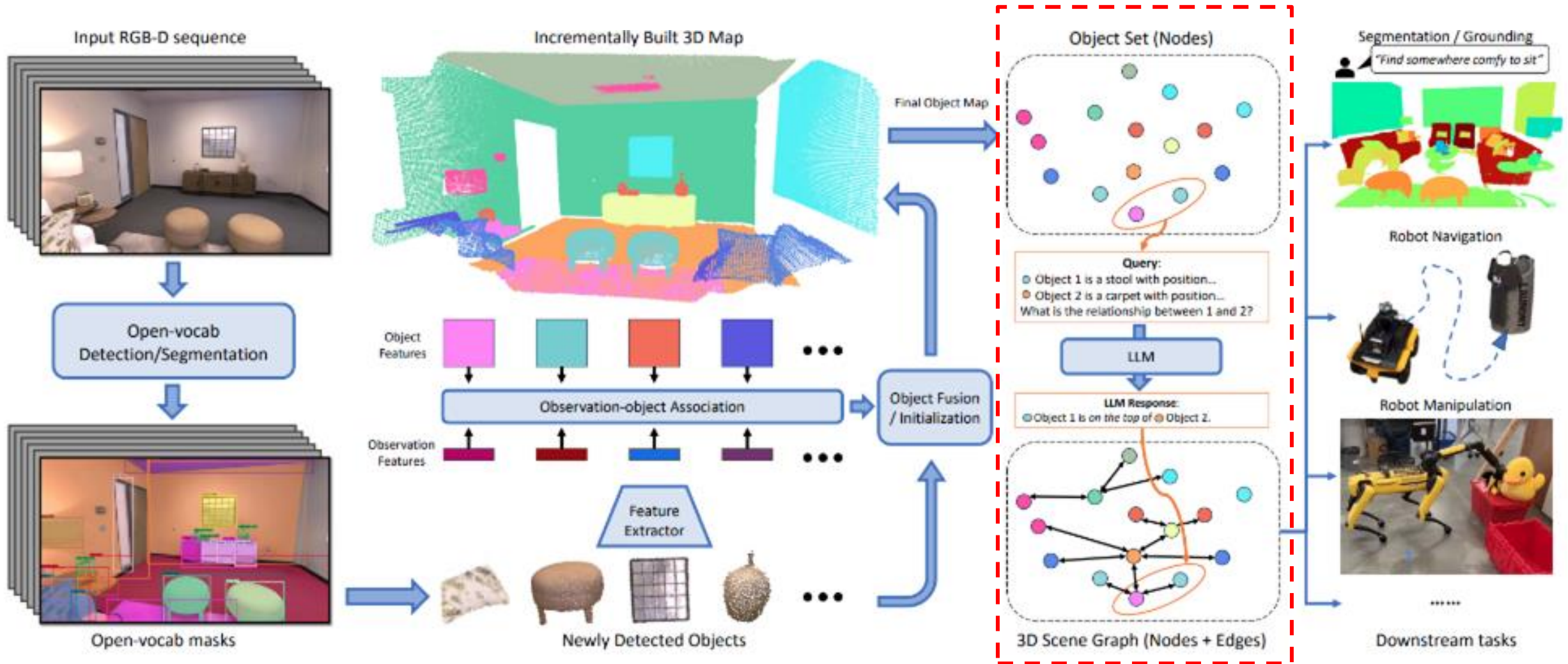
* This configuration is used in Conceptfusion

** On NVIDIA GeForce RTX 2080 Ti

*** We used SAM and MobileSAM with 32x32 seeding points

Research of ConceptGraphs

GPT4 -> Mistral (with two stage query)



Quantitative assessment of LLM performance on a subsample of the sr3d dataset (253 requests)

LLM	Method for constructing graph nodes	Method for constructing graph edges	Recall @ 1	Recall @ 2	Recall @ 3	Number of valid answers (First sentence of LLM answer as top1 prediction)
Mistral-7B-Instruct-v0.1 16-bit quant	GT 3D bbox + GT object tag	-	0.19	0.21	0.21	190
Mistral-7B-Instruct-v0.1 16-bit quant	GT 3D bbox + GT object tag	list of GT objects relations	0.34	0.35	0.36	149
Mistral-7B-Instruct-v0.1 16-bit quant	GT 3D bbox + GT object tag	Objects caption without object ids	0.31	0.33	0.34	233

Object-based with relations

Quantitative assessment of LLM performance on a new sr3d subsample (different ways to construct graph with an LLM, 526 requests)

LLM	Method for constructing graph nodes	Method for constructing graph edges	Recall @ 1 from all answers	Recall @ 1 from valid answers	Invalid answers rate	Number of valid answers (from which id can be selected)
Mixtral 8x7b	two stages: first select semantic objects, second - use their GT bboxes	-	0.60	0.64	0.07	488
Mixtral 8x7b	two stages: first select semantic objects, second - use their GT bboxes	All relations	0.66	0.69	0.05	502
Mistral-7B-Instruct-v0.1 16-bit quant	two stages: first select semantic objects, second - use their GT bboxes	-	0.45	0.46	0.03	513

07

Направления
дальнейшего развития

Что можно улучшить?

Перспективные направления:

- Построение мультимодальных карт в реальном времени
- Использование многоуровневой семантики для представления эмбедингов точек (объектов) карты
- Создание большего количества датасетов и методов для мультимодальных запросов по мультимодальным картам.
- Учет шума в 6DoF позе камеры и картах глубин входной RGB-D последовательности;
- Эффективное распознавание ребер отношений в графе объектов карты;
- Использование одного или нескольких поясняющих примеров или дополнительной информации во время запросов к карте (one/few shot);
- Учет движения (динамики) объектов на карте;
- Временные запросы к динамическим мультимодальным картам (“Место, где люди часто ходят”, “Стул, который передвинули”, “Когда проезжали стол, какие объекты были на нем?”)

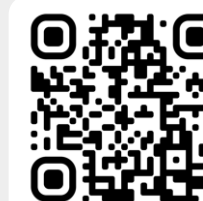
Distilled SAM: nanoSAM

Model	Jetson Orin Nano (ms)		Jetson AGX Orin (ms)		Accuracy (mIoU)			
	Image Encoder	Full Pipeline	Image Encoder	Full Pipeline	All	Small	Medium	Large
MobileSAM	TBD	146	35	39	0.728	0.658	0.759	0.804
NanoSAM (ResNet18)	TBD	27	4.2	8.1	0.706	0.624	0.738	0.7


- NanoSAM is a Segment Anything (SAM) model variant that is capable of running in real-time on NVIDIA Jetson Orin Platforms with NVIDIA TensorRT.
- NanoSAM is trained by distilling the MobileSAM image encoder on unlabeled images. For an introduction to knowledge distillation.
- Resnet18 image encoder, MobileSAM mask decoder.
- Training on COCO 2017 train.



Segment with keypoints (online using TRTPose detections)



Project

Github 

Что можно улучшить?

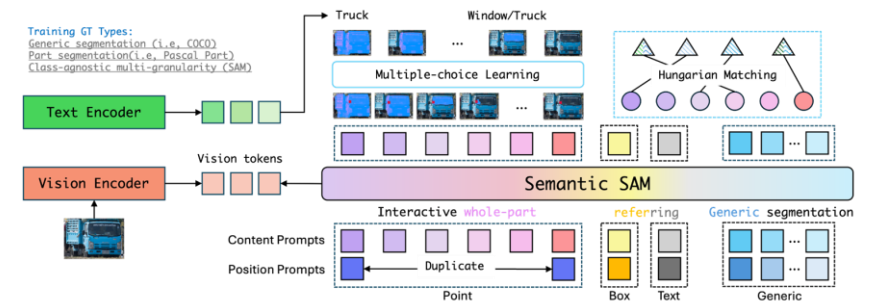
Перспективные направления:

- Построение мультимодальных карт в реальном времени
- Использование многоуровневой семантики для представления эмбедингов точек (объектов) карты
- Создание большего количества датасетов и методов для мультимодальных запросов по мультимодальным картам.
- Учет шума в 6DoF позе камеры и картах глубин входной RGB-D последовательности;
- Эффективное распознавание ребер отношений в графе объектов карты;
- Использование одного или нескольких поясняющих примеров или дополнительной информации во время запросов к карте (one/few shot);
- Учет движения (динамики) объектов на карте;
- Временные запросы к динамическим мультимодальным картам (“Место, где люди часто ходят”, “Стул, который передвинули”, “Когда проезжали стол, какие объекты были на нем?”)

Query-based Segmentation with multi-level semantics: Semantic-SAM



- Semantic-SAM model supports Generic Segmentation, Part Segmentation, Multi-level Semantic Segmentation, Multi-Level Image Editing



Project

Github

Что можно улучшить?

Перспективные направления:

- Построение мультимодальных карт в реальном времени
- Использование многоуровневой семантики для представления эмбедингов точек (объектов) карты
- Создание большего количества датасетов и методов для мультимодальных запросов по мультимодальным картам.
- Учет шума в 6DoF позе камеры и картах глубин входной RGB-D последовательности;
- Эффективное распознавание ребер отношений в графе объектов карты;
- Использование одного или нескольких поясняющих примеров или дополнительной информации во время запросов к карте (one/few shot);
- Учет движения (динамики) объектов на карте;
- Временные запросы к динамическим мультимодальным картам (“Место, где люди часто ходят”, “Стул, который передвинули”, “Когда проезжали стол, какие объекты были на нем?”)



Контакты



Дмитрий Александрович Юдин

к.т.н., заведующий лабораторией интеллектуального транспорта МФТИ - НКБ ВС, Центр когнитивного моделирования МФТИ, старший научный сотрудник AIRI

✉ yudin@airi.net

📧 @yuddim

cogmodel.mipt.ru

