**Aalto-yliopisto**
**Aalto-universitetet**
**Aalto University**

**UNIVERSITY OF TURKU**

—

# Semantic Mapping
# From Segmentation to Scene Graphs

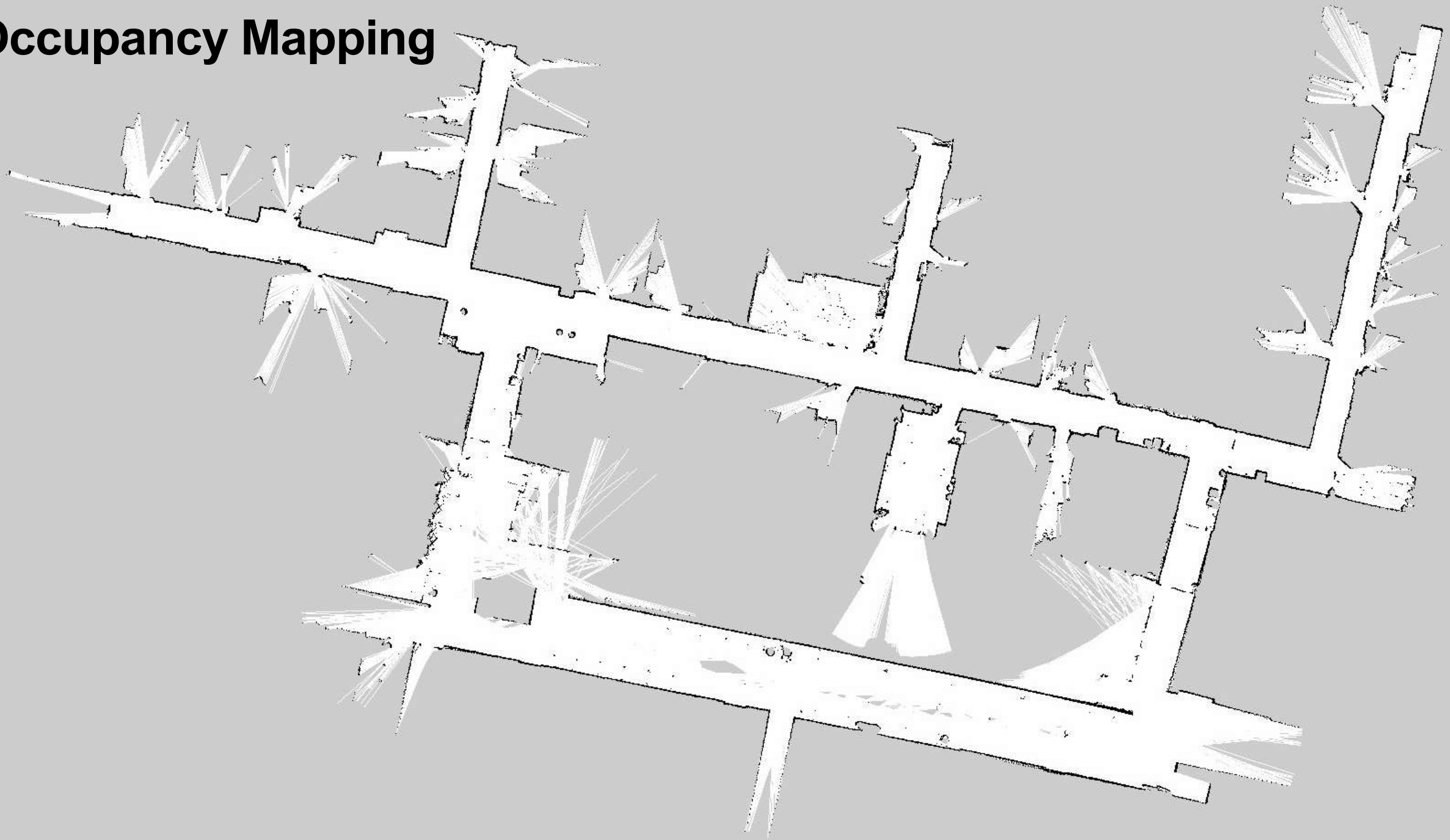**Francesco Verdoja**
Academy Research Fellow
Dept. of Electrical Engineering and Automation

9.12.2025

# Lecture outline

- Recap of occupancy mapping
- Semantics?
- Metric-semantic mapping
- Open-vocabulary semantic mapping
- Scene Graphs
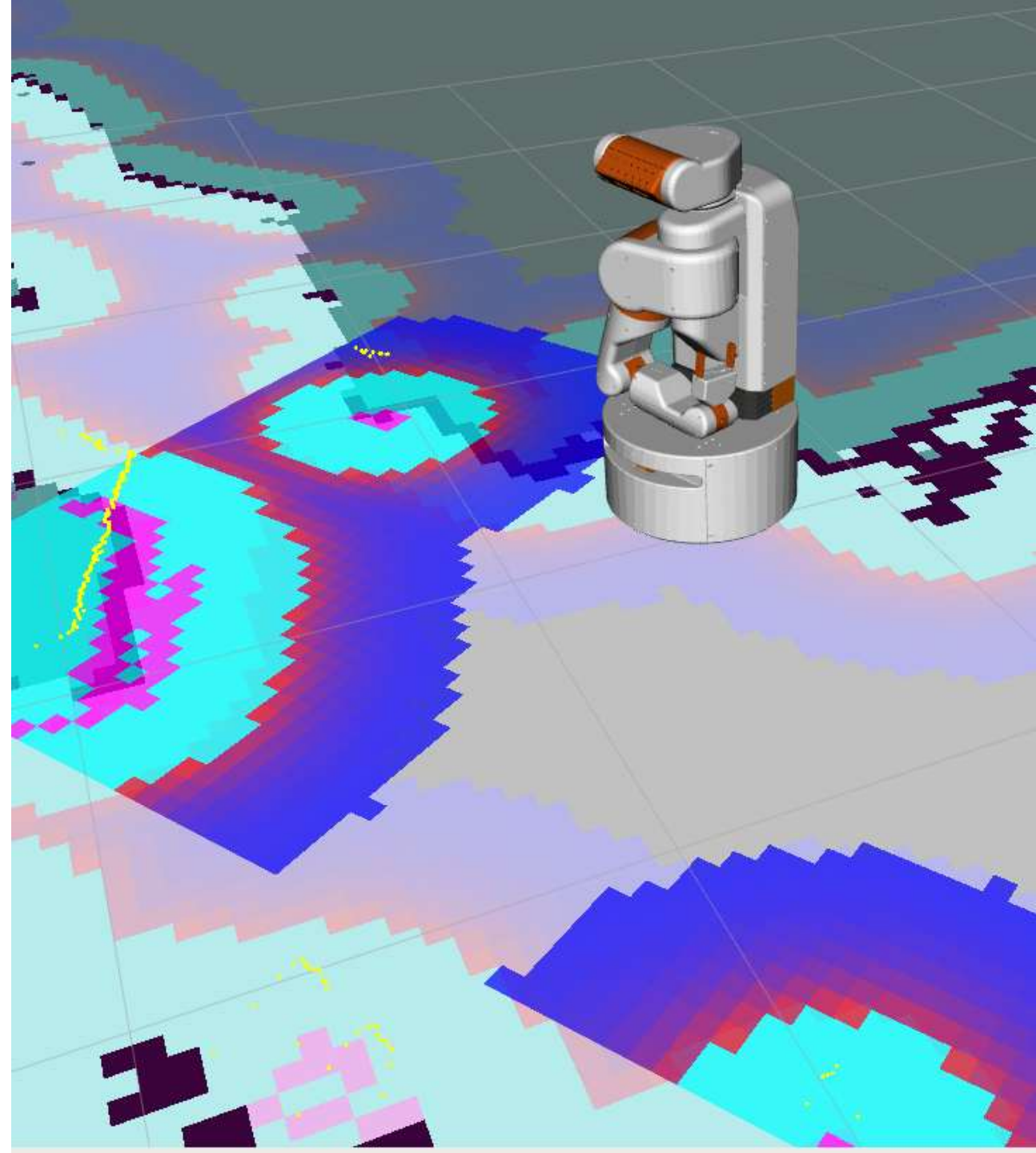- Adding time to semantic maps

**A!**

**Occupancy Mapping**

# Capabilities supported by occupancy maps

- Global path planning given a goal point

- Local path planning / Obstacle avoidance

- Replanning around blockages

- Localization

- Navigation

- Docking

- Frontier-based exploration
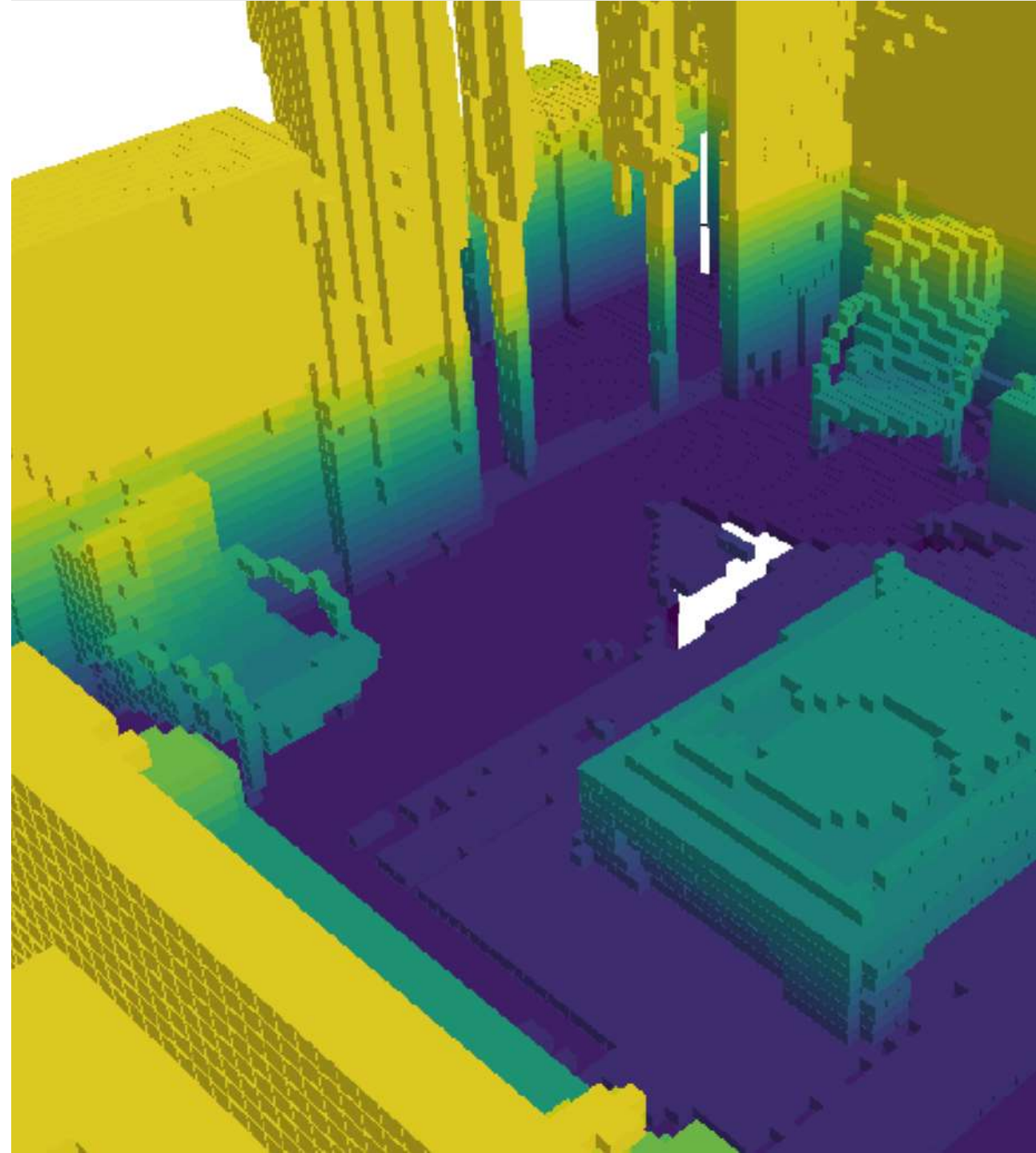
- ...

**A!**

# Applications

# Limitations

*"Move the chair closer to the table"*

- Object instances?

- Object extents?

- Affordances?

- Grid independence?

**Occupancy not enough for complex tasks (mobile manipulation, natural HRI)**

**A!**

# Semantics?
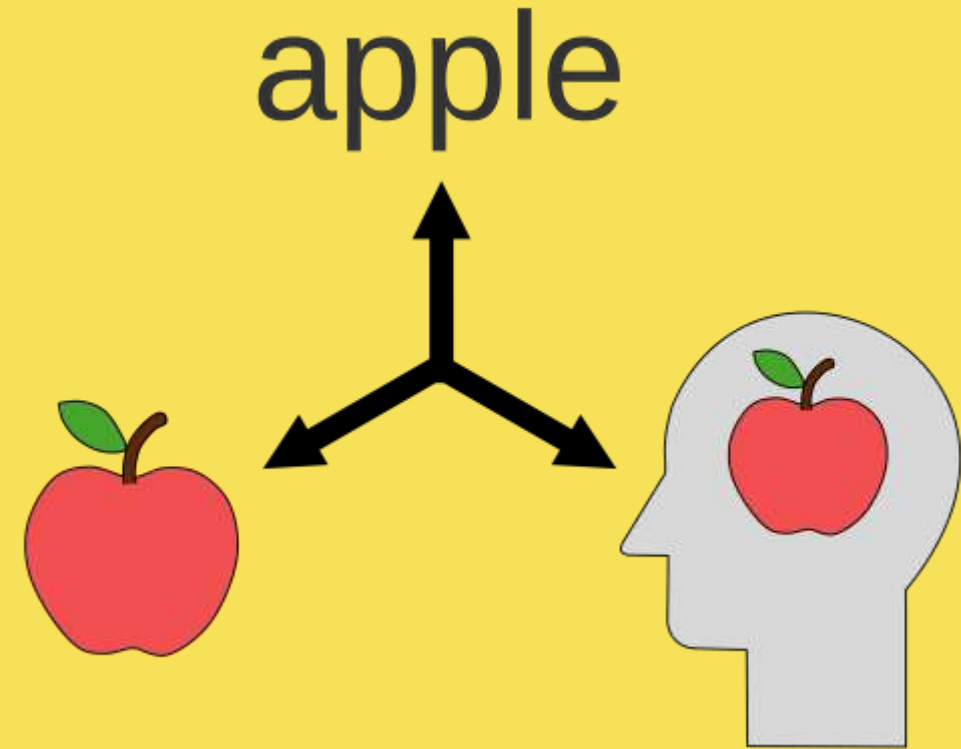
A!

# Semantics

- *Meaning* of things and words
- *Grounding* of symbols in reality

**Enable decisions that depend on:**

- Object identity (shelf vs wall)
- Function (charging stations)
- Affordances (door is *openable*)
- Human language (*"move the chair"*)

apple

A!

# Semantic Mapping

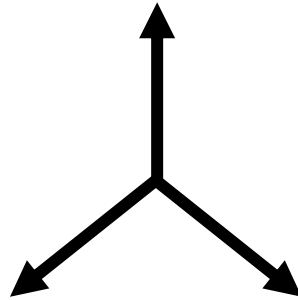***Vocabulary***

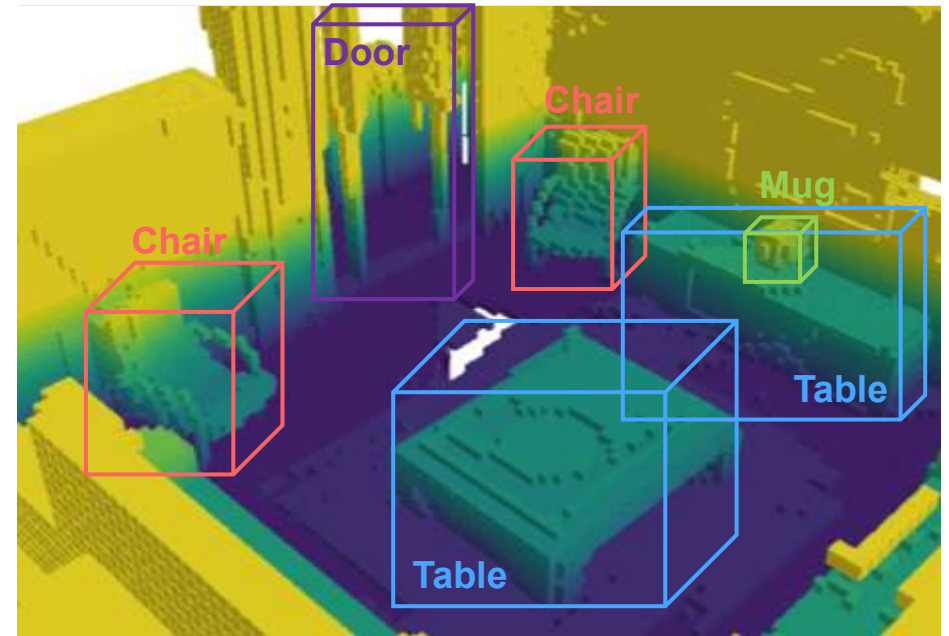Chair [seat, move, …]

Table [carrier, …]

Mug [drink, move, …]

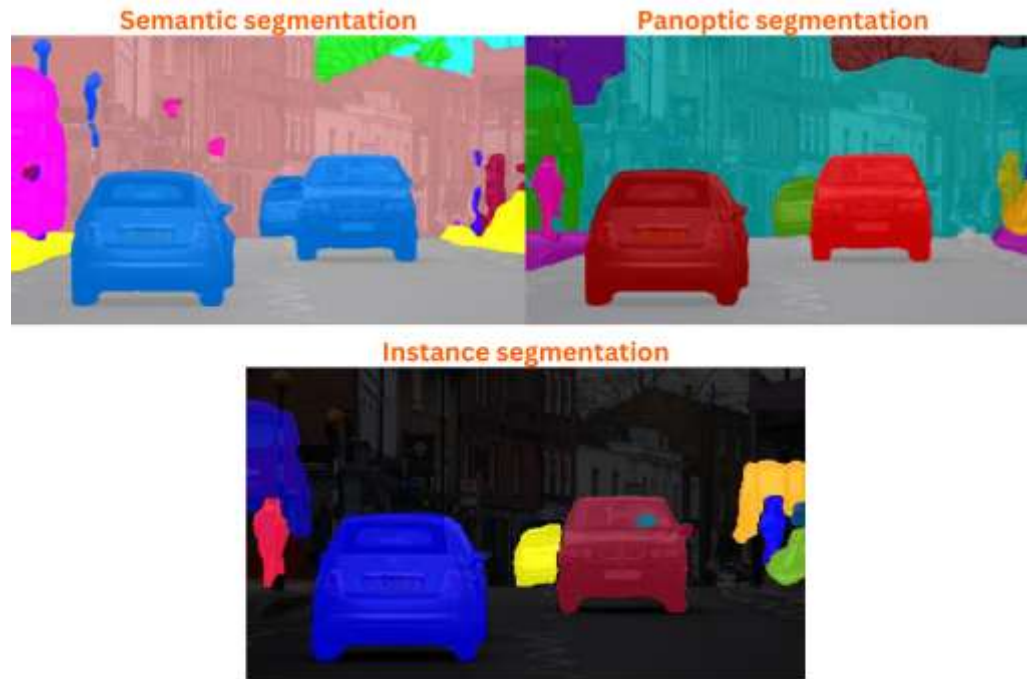Door [open, close, pass, …]

…

***Perception***
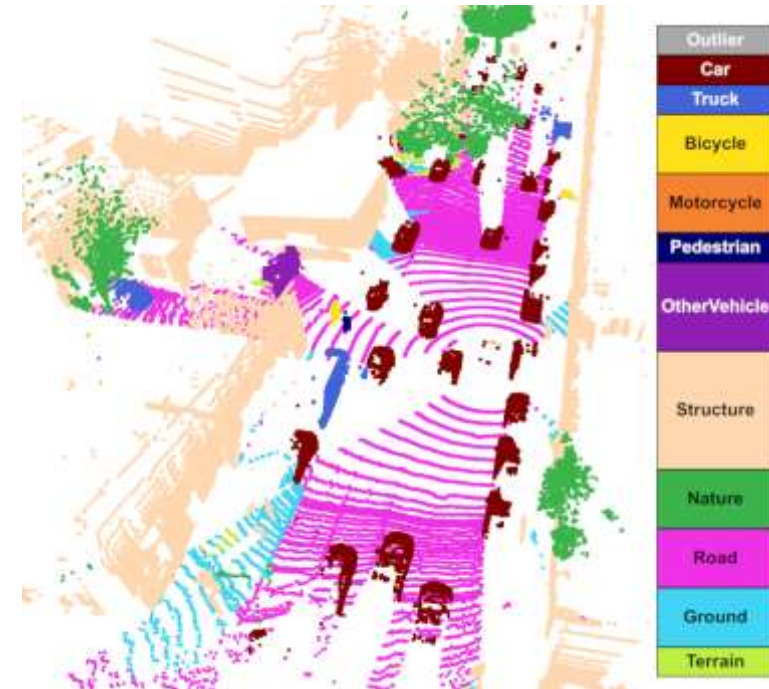
***Mapping***



A!

# Where do you get semantics?

**RGB camera**
Mask2Former, YOLO, Segment Anything…



Semantic segmentation

Panoptic segmentation

Instance segmentation

Requires raycasting!

**LiDAR**
RangeNet++, RandLA-net…



Outlier
Car
Truck
Bicycle
Motorcycle
Pedestrian
OtherVehicle
Structure
Nature
Road
Ground
Terrain

**A!**

# Metric-semantic mapping
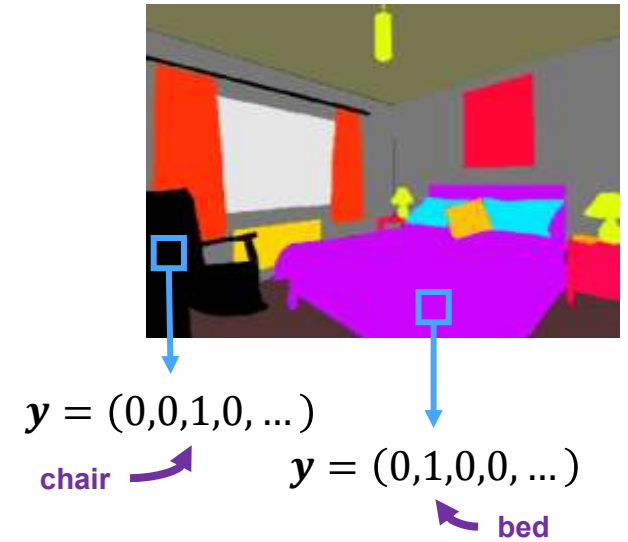
# Extending Occupancy Mapping

- Occupancy maps: O = {occupied, free}
- Semantic maps: C = {chair, table, door, mug, …}
- Both are **categorical distributions** $\mathrm{Cat}(K, \boldsymbol{p})$
  - $K > 0$ number of categories ($K = 2$ for occupancy map, *Bernoulli distribution*)
  - $\boldsymbol{p} = (p_1, p_2, \ldots, p_K)$ probabilities of individual categories ($p_i \geq 0, \sum p_i = 1$)
  - Mode (i.e., most likely category): $i \mid p_i = \max(p_1, \ldots, p_K)$

**A!**

# Semantic mapping as Bayesian inference

- ***For any map voxel*** $\mathrm{Cat}(K, \boldsymbol{v})$***:***
  $\boldsymbol{v} = (v_1, \dots, v_K)$ where $v_i \geq 0$ and $\sum v_i = 1$

- ***Measurement (one-hot):***
  $\mathbf{y} = (y_1, \dots, y_K)$, where $y_i \in \{0,1\}$ and $\sum y_i = 1$

- ***Categorical likelihood:*** $p(\boldsymbol{y}|\boldsymbol{v}) = \prod v_i^{y_i}$



$\boldsymbol{y} = (0,0,1,0,\dots)$

chair

$\boldsymbol{y} = (0,1,0,0,\dots)$

bed

But how to find posterior $p(\boldsymbol{v}|\boldsymbol{y})$?

**A!**

# Dirichlet conjugate prior for categorical distributions

For any categorical distribution $\mathrm{Cat}(K, \boldsymbol{v})$:

- Given a **concentration hyperparameter** $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)$

- **Dirichlet (conjugate) prior:** $p(\boldsymbol{v}|\boldsymbol{\alpha}) \sim \mathrm{Dir}(K, \boldsymbol{\alpha}) = \frac{1}{\mathrm{B}(\boldsymbol{\alpha})} \prod v_i^{\alpha_i - 1}$

- $\boldsymbol{c} = (c_1, \dots, c_K)$, number of observations of each category

- $\boldsymbol{\alpha'} = \boldsymbol{c} + \boldsymbol{\alpha} = (c_1 + \alpha_1, \dots, c_K + \alpha_K)$, and $S(\boldsymbol{\alpha'}) = \sum \alpha_i'$
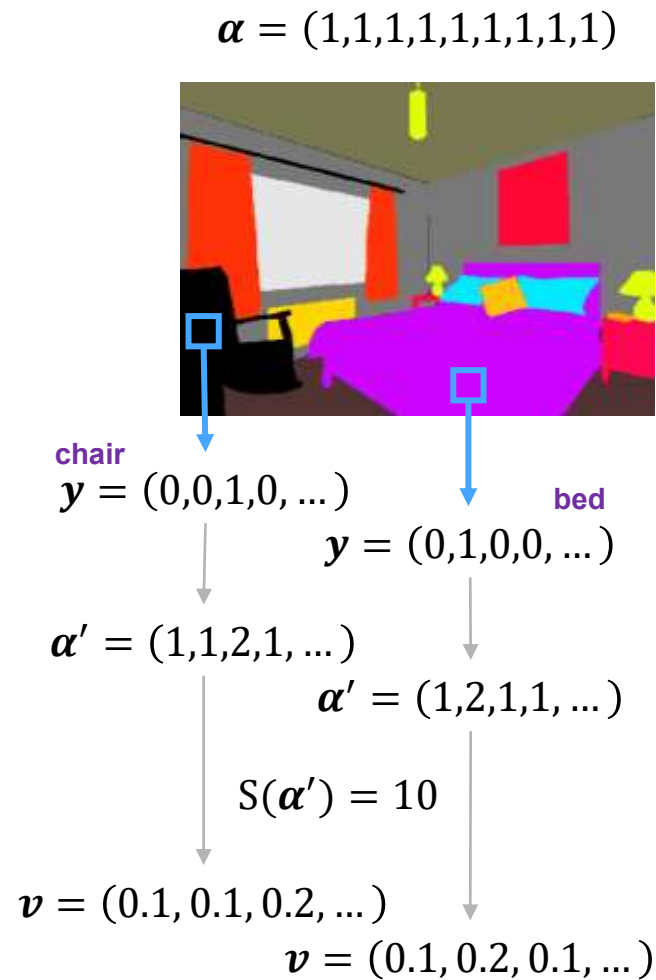
**Then:**

$$p(\boldsymbol{v}|\boldsymbol{c}, \boldsymbol{\alpha}) \sim \mathrm{Dir}(\boldsymbol{c} + \boldsymbol{\alpha}) \sim \mathrm{Dir}(\boldsymbol{\alpha'})$$
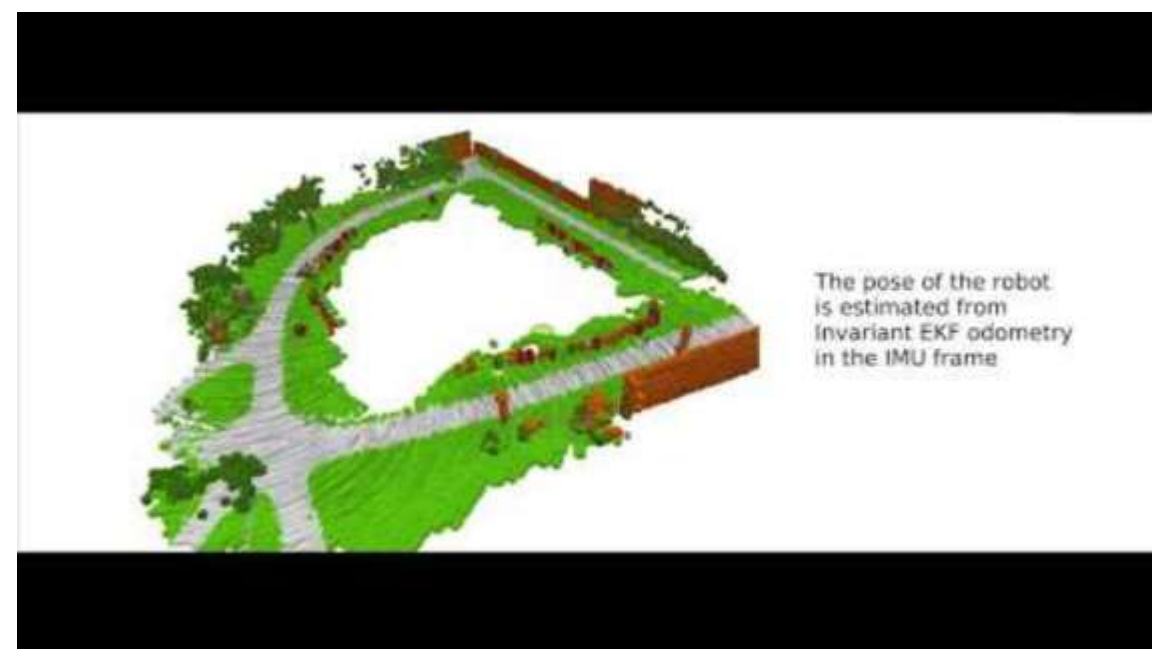
$$\mathbb{E}[v_i] = \frac{\alpha_i'}{S(\boldsymbol{\alpha'})} \qquad \mathbb{V}[v_i] = \frac{\alpha_i'(S(\boldsymbol{\alpha'}) - \alpha_i')}{S(\boldsymbol{\alpha'})^2 (S(\boldsymbol{\alpha'}) + 1)}$$
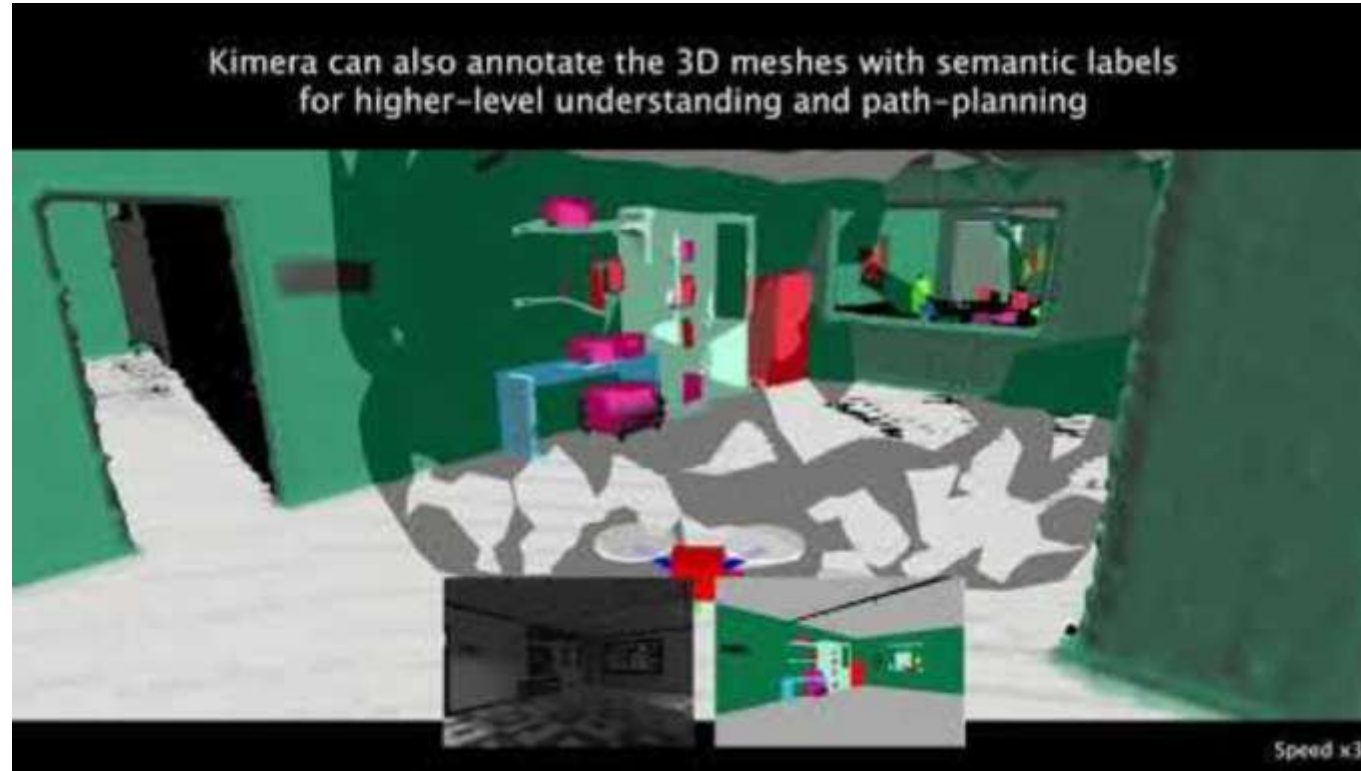
**A!**

# Back to semantic mapping

$\alpha = (1,1,1,1,1,1,1,1,1)$



**chair**
$y = (0,0,1,0,\dots)$

**bed**
$y = (0,1,0,0,\dots)$

$\alpha' = (1,1,2,1,\dots)$

$\alpha' = (1,2,1,1,\dots)$

$S(\alpha') = 10$

$v = (0.1, 0.1, 0.2, \dots)$

$v = (0.1, 0.2, 0.1, \dots)$

$$\mathbb{E}[v_i] = \frac{\alpha'_i}{S(\alpha')} \text{ and } \mathbb{E}[v] = \text{argmax}_i(v_i)$$



The pose of the robot is estimated from Invariant EKF odometry in the IMU frame

Gan, Lu, et al. "Bayesian spatial kernel smoothing for scalable dense semantic mapping."
*IEEE Robotics and Automation Letters* 5.2 (2020): 790-797.

**A!**

# 3D meshes instead of voxels (e.g., KIMERA)



A. Rosinol, et al. "Kimera: From SLAM to spatial perception with 3D dynamic scene graphs,"
*The Int. J. of Robotics Research*, vol. 40, no. 12-14, pp. 1510–1546, 2021

# A lot is borrowed from Visual SLAM
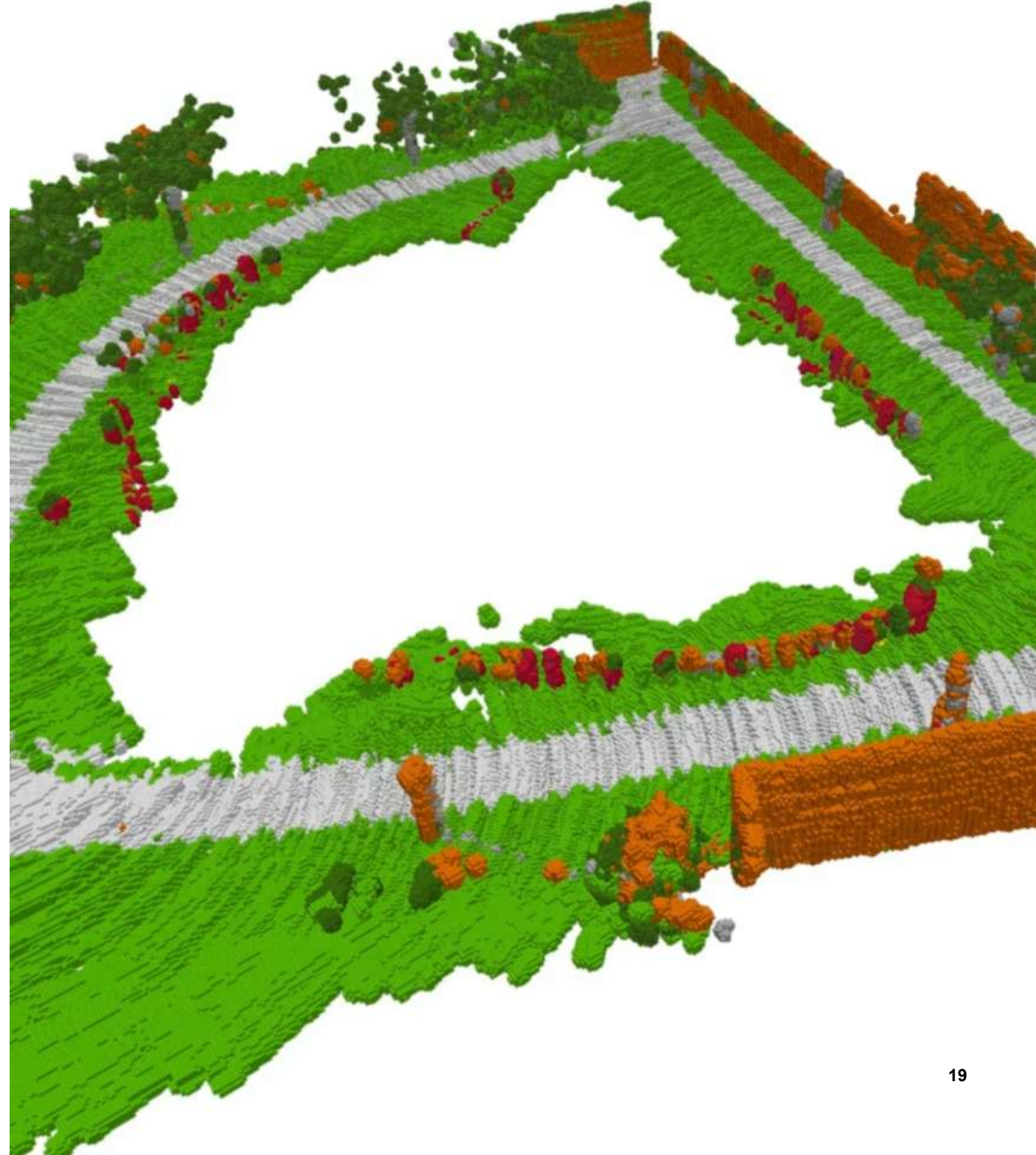
# Panoptic Maps for object instances



Narita, Gaku, et al. "Panopticfusion: Online volumetric semantic mapping at the level of stuff and things.«
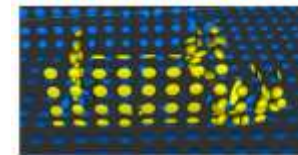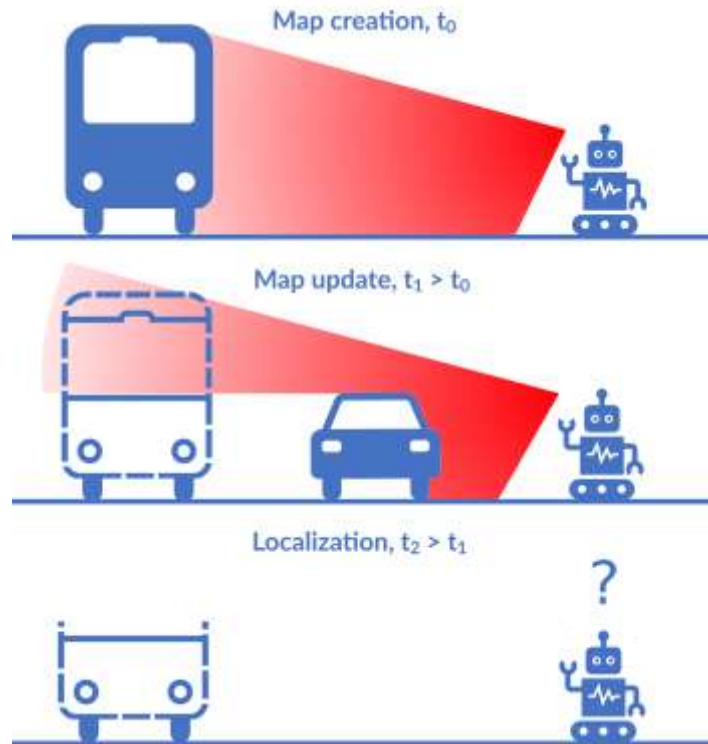*2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019.

# What have we fixed?

- Object instances?  ◑
- Object extents?  ✗
- Affordances?  ◑
- Grid Independence?  ✗

**A!**

# Some efforts on addressing grid independence



Map creation, $t_0$

Map update, $t_1 > t_0$

Localization, $t_2 > t_1$

?



(a) Before the update     (b) NDT-OM     (c) Our method

Object-Oriented Grid Mapping in Dynamic Environments
Matti Pekkanen, Francesco Verdoja, and Ville Kyrki
*2024 IEEE Int. Conf. on Multisensor Fusion and Integration for Intelligent Systems (MFI)*, 2024

**A!**

# What can we do?

- Semantic-aware navigation

  - "stay on the road"

  - "stop at pedestrian crossings"

  - "never go closer than 2m to a tree"
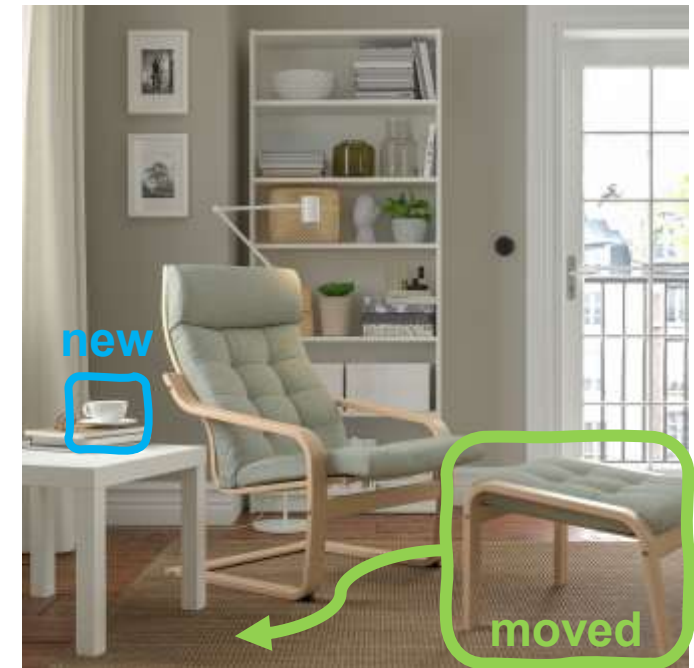
  - "rest close to a wall"

**A!**

# New challenges

Class bleeding at boundaries
(RGB+D calibration)

High memory footprint
(submapping)
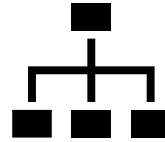
Complex map update and
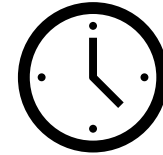vocabolary extension

# Beyond metric-semantic mapping

### Open-vocabulary

Not limited to a closed set of predefined semantic labels

### From voxels to concepts

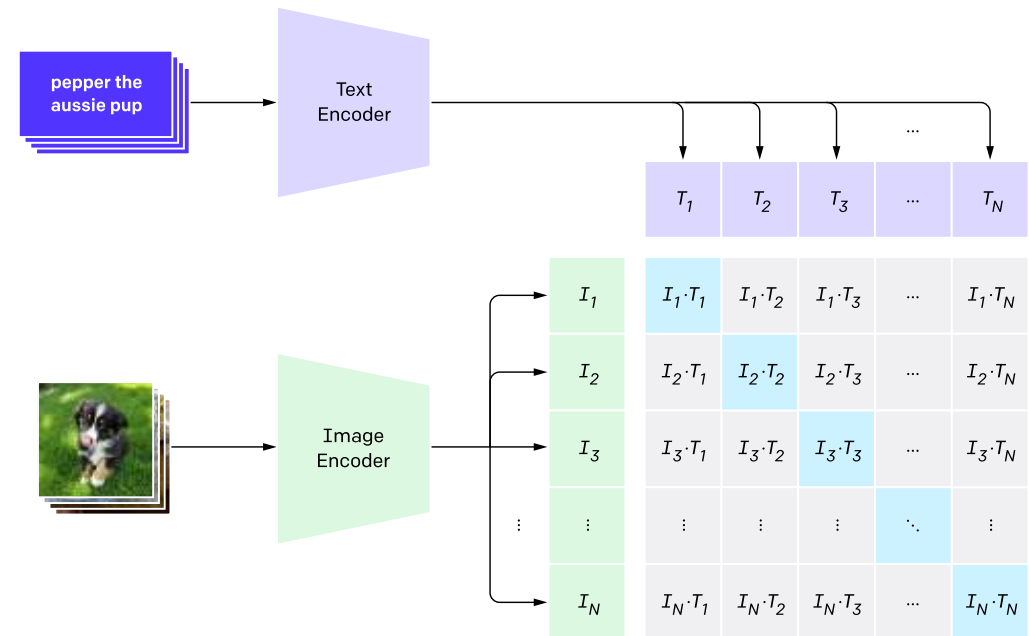Voxels are part of objects, rooms, and other semantic entities

### From 3D to 4D+

Reasoning and handling dynamic environments over time

**A!**

# Open-vocabulary semantic maps

**A!**

# Visual-Language Models

- Coupled Transformer Neural Networks
  - Text: to $N$-dim embedding $T$
  - Images: to $N$-dim embeddings $I$
- Trained on (image, text caption) dataset
- Minimize distance between $T$ and $I$ for (image, caption) pair
- Maximize distance between $T$ and $I$ for non-pairs
- **CLIP** from OpenAI: 512-dim embedding



Radford, Alec, et al. "Learning transferable visual models from natural language supervision." *International conference on machine learning*. PmLR, 2021.
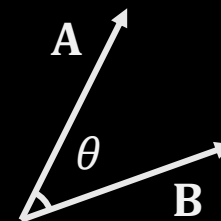
# Embeddings and Cosine Similarity

## *Embeddings*

- $N$-dim vectors $\mathbf{A} = (A_1, \ldots, A_N)$

- Often for VLMs they are unit-size, i.e., $\|\mathbf{A}\| = 1$
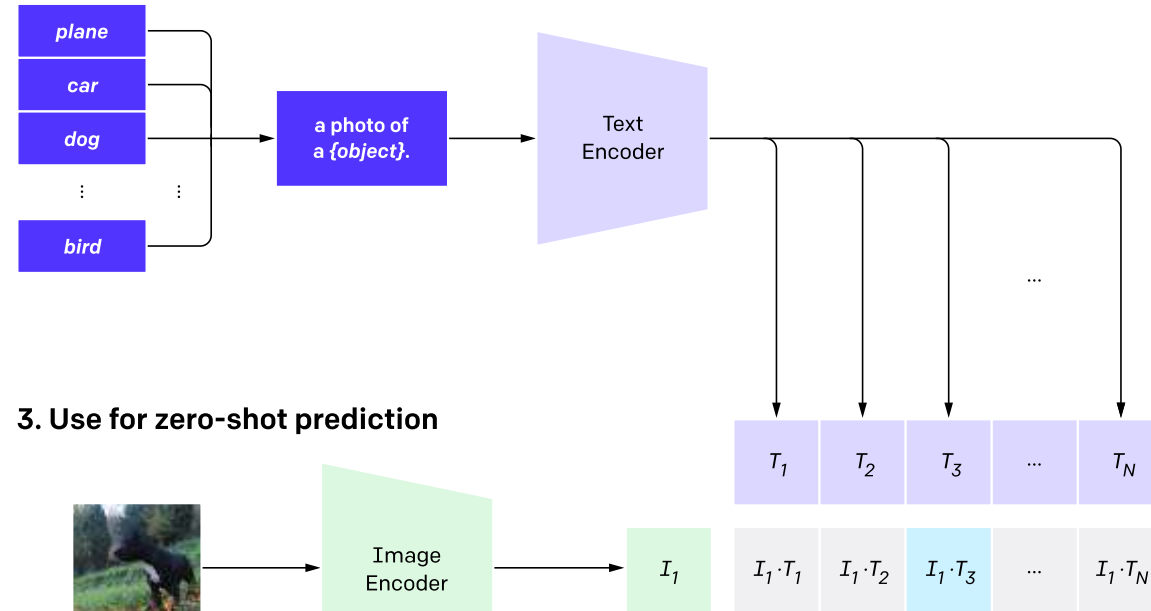
## *Cosine similarity*

- $S_C(A, B) := \cos(\theta) = \dfrac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|\|\mathbf{B}\|} = \dfrac{\sum_{i=1}^{N} A_i B_i}{\sqrt{\sum_{i=1}^{N} A_i^2} \cdot \sqrt{\sum_{i=1}^{N} B_i^2}}$

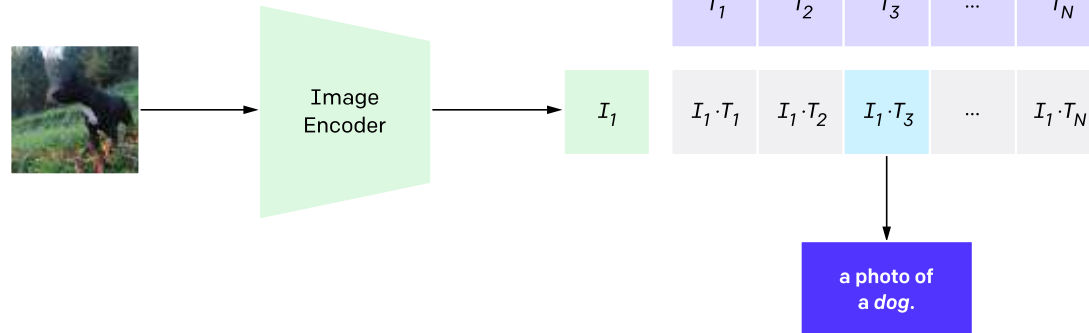- $S_c \in [-1, 1]$, with -1 opposite, +1 same, 0 orthogonal
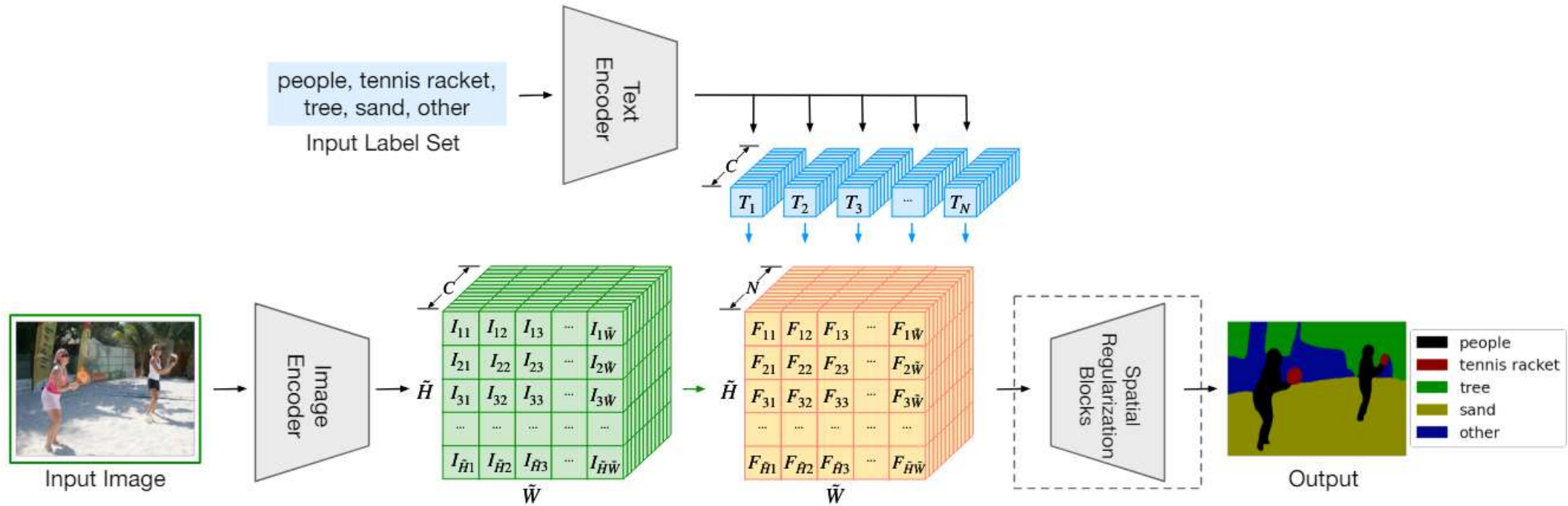
**A!**

# Querying the model

**2. Create dataset classifier from label text**

plane
car
dog
⋮    ⋮
bird

a photo of a {object}.

Text Encoder

$T_1$ | $T_2$ | $T_3$ | ... | $T_N$

**3. Use for zero-shot prediction**

Image Encoder

$I_1$

$I_1 \cdot T_1$ | $I_1 \cdot T_2$ | $I_1 \cdot T_3$ | ... | $I_1 \cdot T_N$
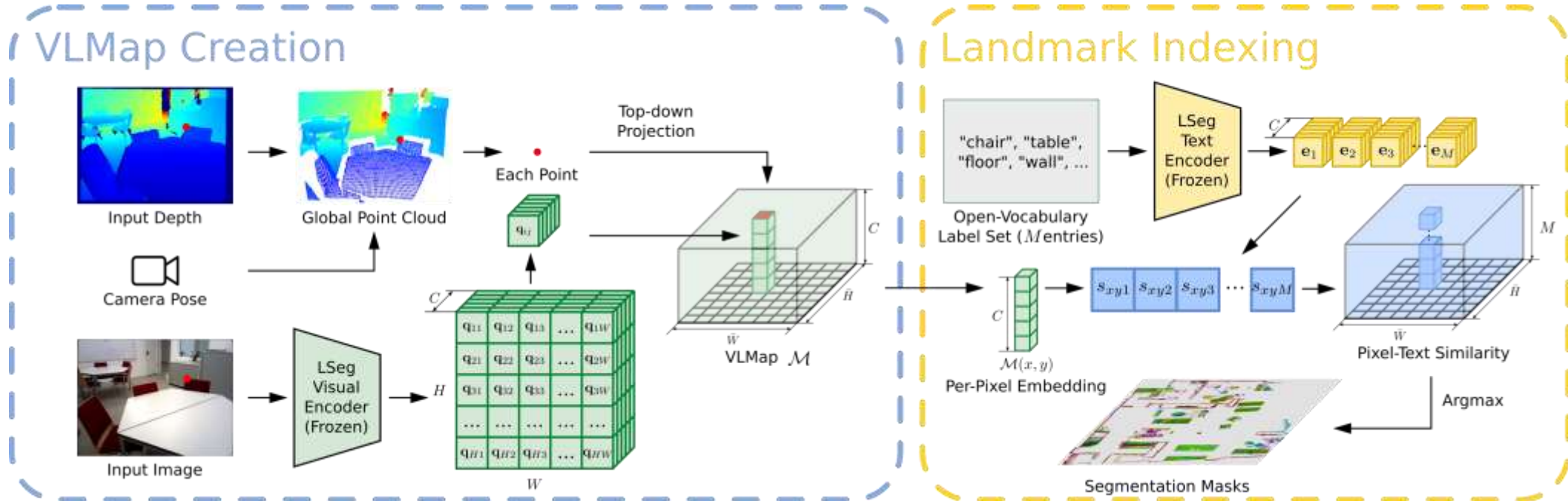
a photo of a *dog*.

**A!**

# Pixel-level CLIP embeddings (e.g., LSeg)



Li, Boyi, et al. "Language-driven Semantic Segmentation."
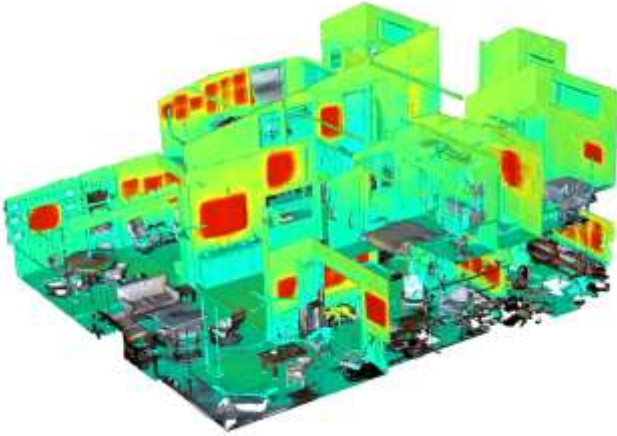*2022 International Conference on Learning Representations (ICLR)*, 2022.
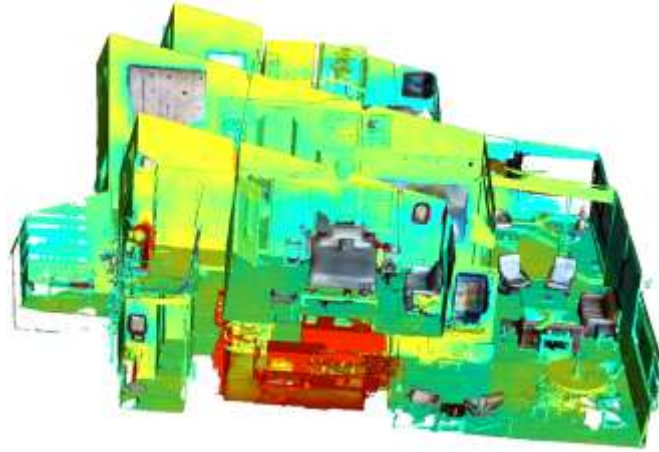
# Maps of Embeddings (e.g., VLMap)



Huang, Chenguang, et al. "Visual Language Maps for Robot Navigation."
*2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023.
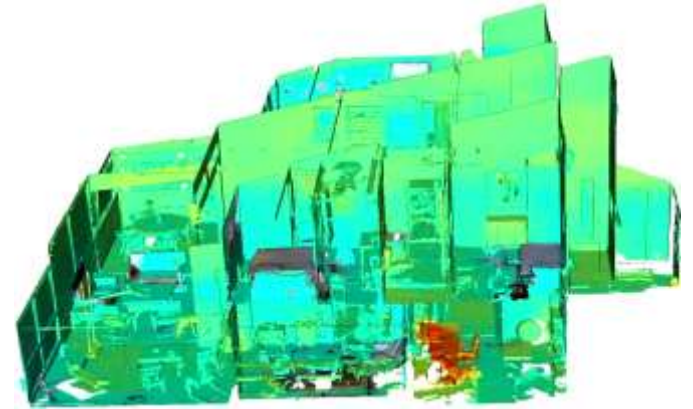
# Open-vocabulary Querying

Painting

Kitchen

Work



Matti Pekkanen, Tsvetomila Mihaylova, Francesco Verdoja, and Ville Kyrki, "Do Visual-Language Grid Maps Capture Latent Semantics?" *2025 IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS)*, IEEE, 2025.

**A!**

# Robot interaction and planning (e.g., NLMap + SayCan)



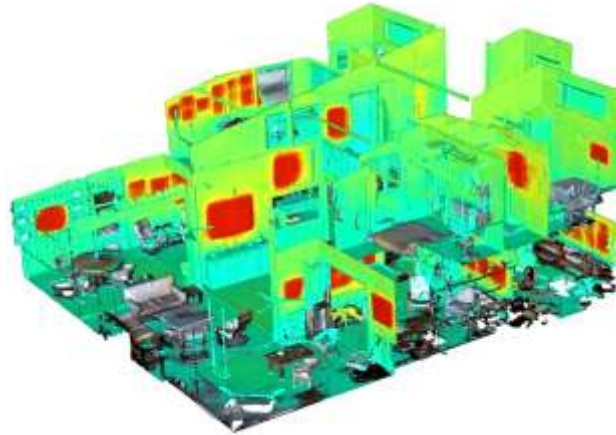We can also run frontier exploration for any novel environment.

Chen, Boyuan, et al. "Open-vocabulary Queryable Scene Representations for Real World Planning."
*2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023.

A!

# Challenges

Quality is highly dependant on VLM performance

Maps are even larger (WxLxHx512)

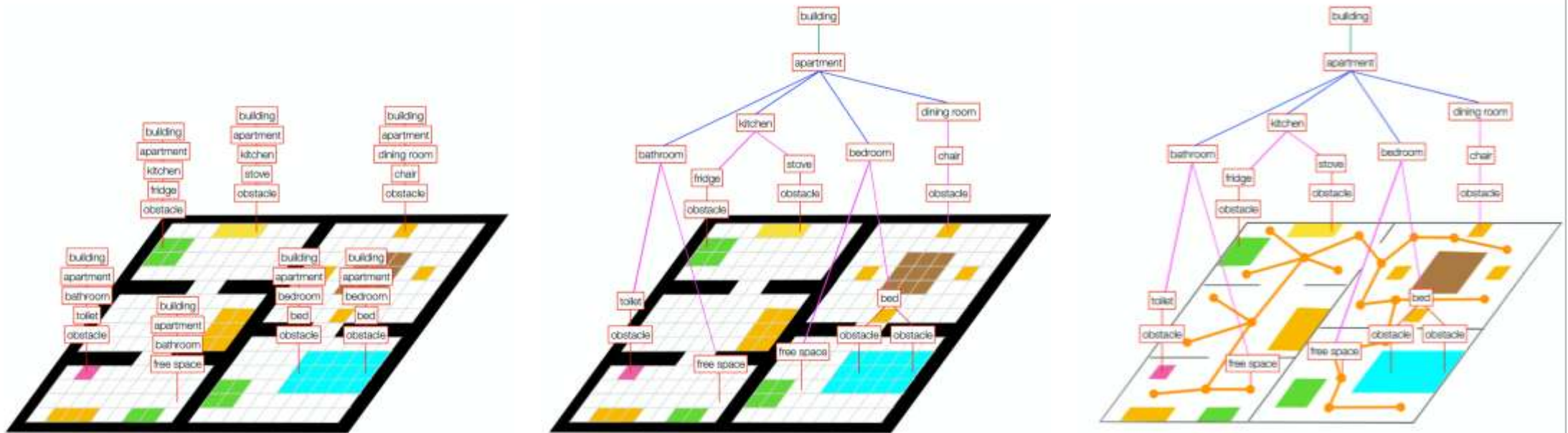What happens for queries with missing target?

# From voxels to concepts

A!

# The environment is hierarchical

# 3D Scene Graphs (3DSGs)

- **Hierarchical graph representation**
- **Objects, places, and rooms as nodes**
  - attributes (pose, shape, affordances)
  - connected to 3D mesh
  - Belong to semantic layers
- **Edges describe relations**
  - spatial (adjacency, inclusion, support)
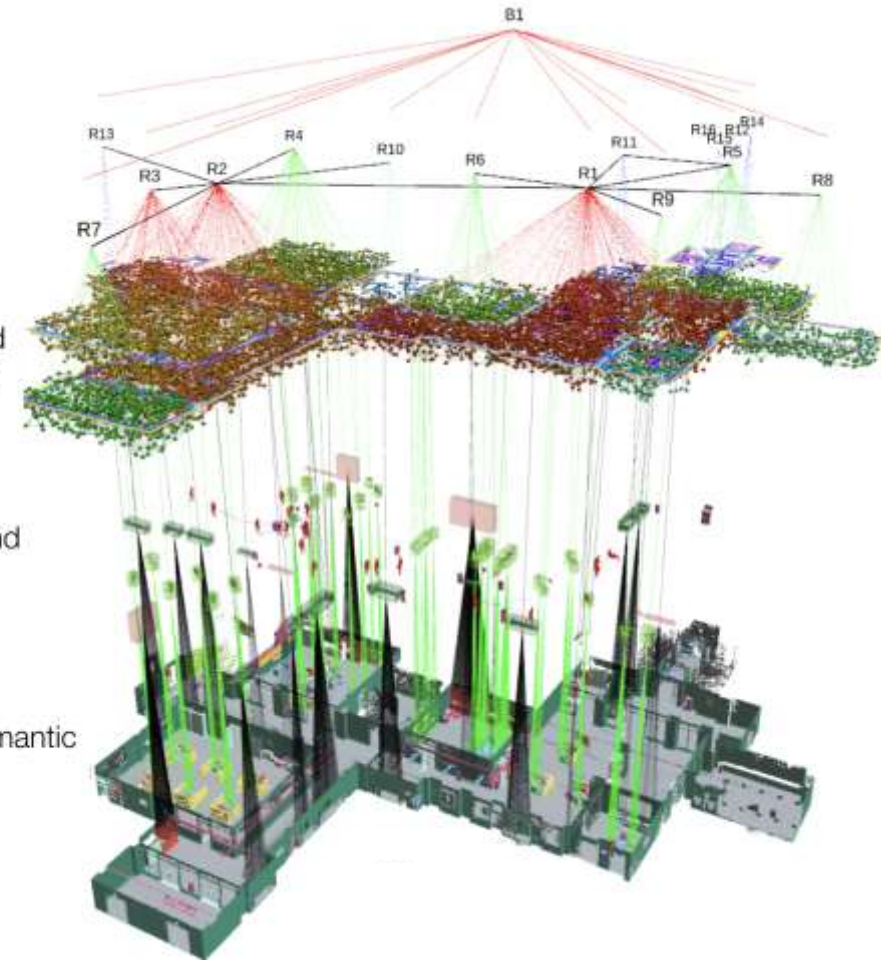  - functional (used-for, part-of)
  - …



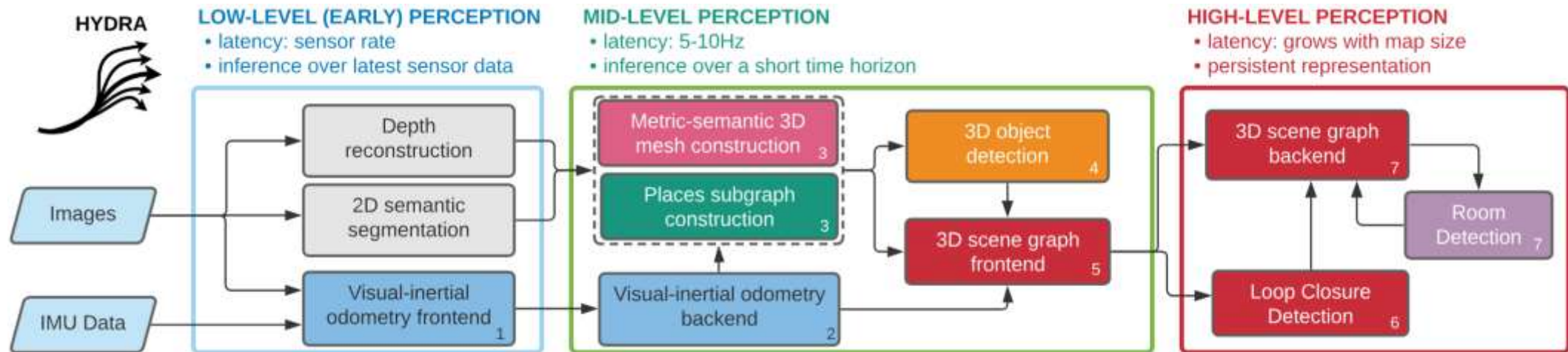Layer 5: Buildings

Layer 4: Rooms

Layer 3: Places and Structures

Layer 2: Objects and Agents

Layer 1: Metric-Semantic Mesh

A. Rosinol, et al. "Kimera: From SLAM to spatial perception with 3D dynamic scene graphs," *The Int. J. of Robotics Research*, vol. 40, no. 12-14, pp. 1510–1546, 2021

A!

# Building a 3D Scene Graph



Hughes, Nathan, Yun Chang, and Luca Carlone. "Hydra: A Real-time Spatial Perception System for 3D Scene Graph Construction and Optimization." *Robotics: Science and Systems*. 2022.

# Hydra in action



Goal: (and (ObjectAtPlace O105 P909)
            (VisitedPlace P2700)
            (Safe O130) ✓
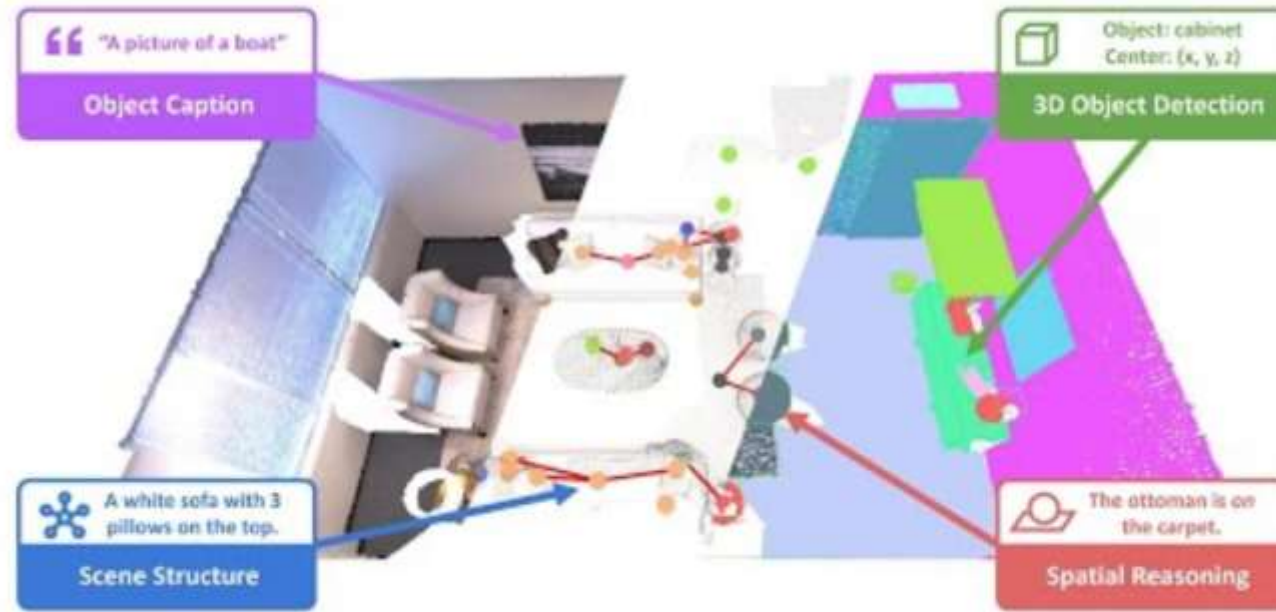            (not (VisitedPlace P1153)))

(Pick Object$_{110}$ Pose$_4$ Pose$_5$)

Aaron Ray, et al. "Task and Motion Planning in Hierarchical 3D Scene Graphs,"
*International Symposium of Robotics Research (ISRR)*, 2024

A!

# Scene graphs + embeddings (e.g., ConceptGraph)



Gu, Qiao, et al. "Conceptgraphs: Open-vocabulary 3d scene graphs for perception and planning."
*2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024.

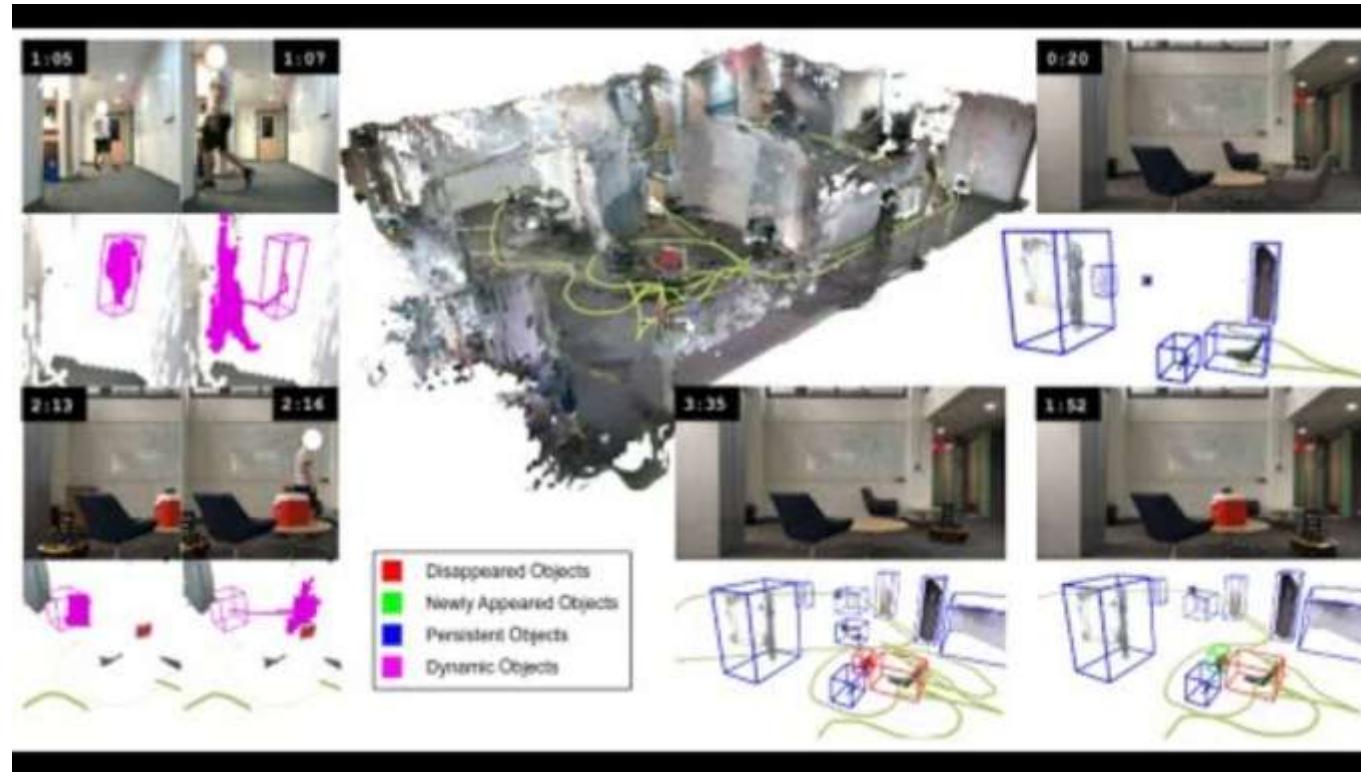# Using scene graphs in planning (e.g., SayPlan)



Rana, Krishan, et al. "SayPlan: Grounding Large Language Models using 3D Scene Graphs for Scalable Robot Task Planning." *Conference on Robot Learning*. PMLR, 2023.
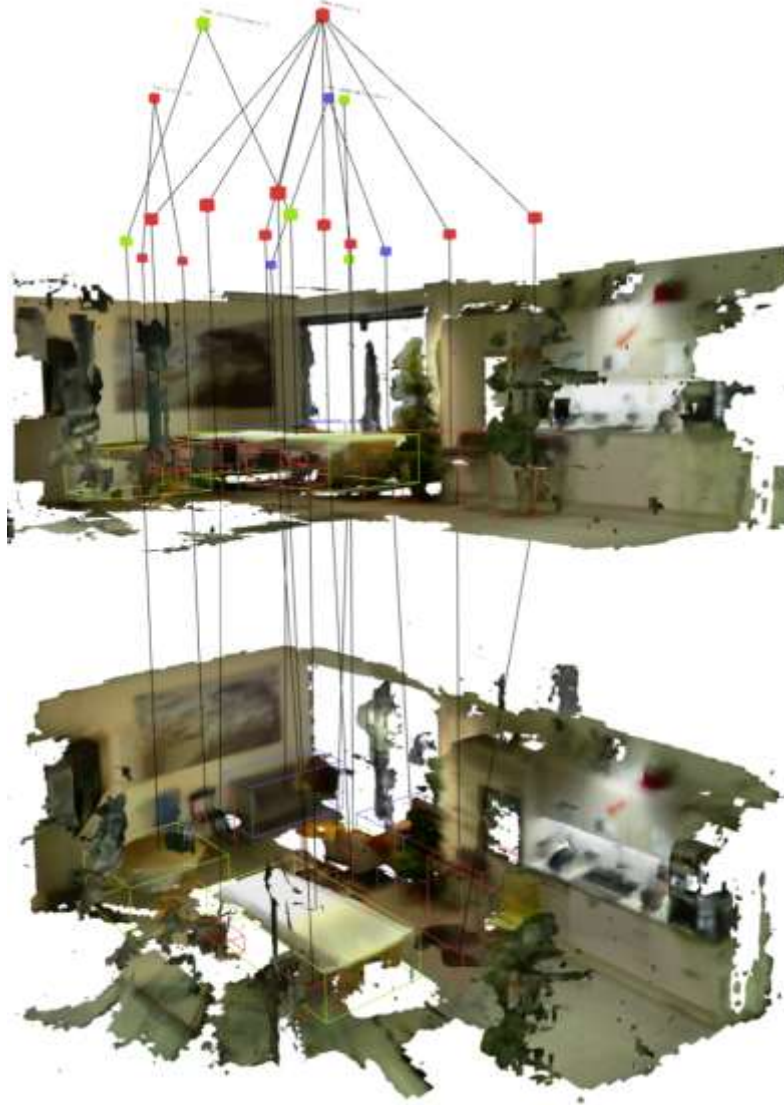
A!

# From 3D to 4D+
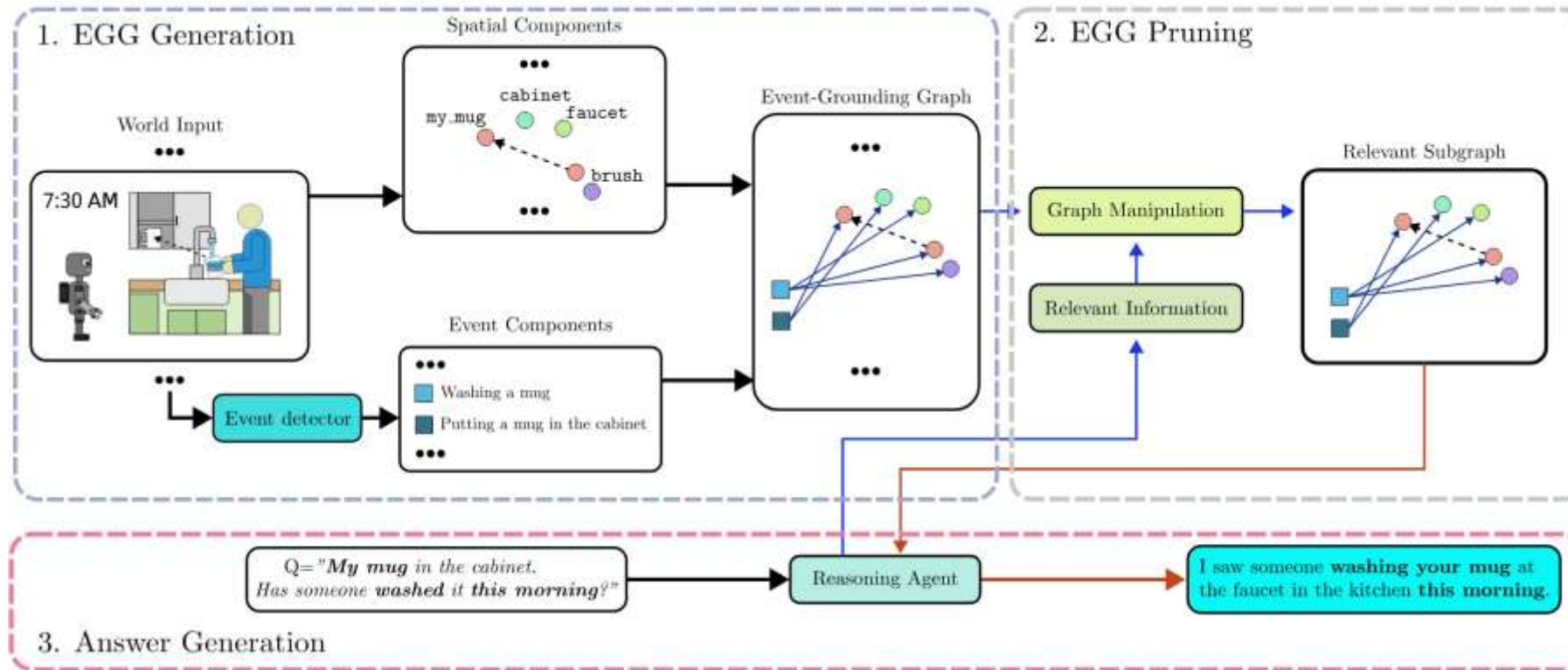
A!

# Dynamic Scene Graphs (e.g., Khronos)



Schmid, Lukas, et al. "Khronos: A Unified Approach for Spatio-Temporal Metric-Semantic SLAM in Dynamic Environments." *Robotics: Science and Systems*. 2024.

**A!**

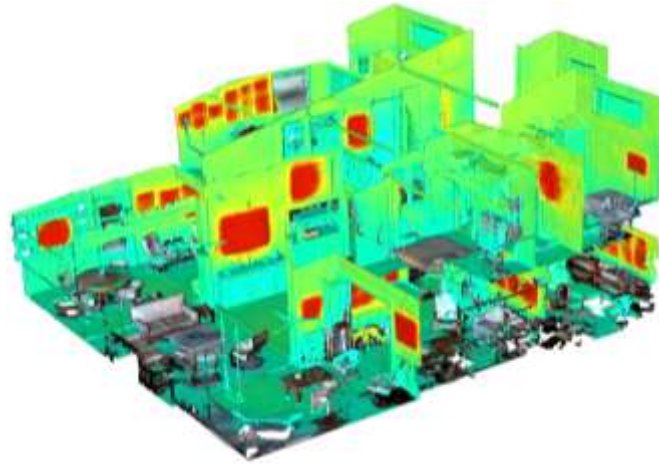# Updatable Scene Graphs (e.g., REACT)
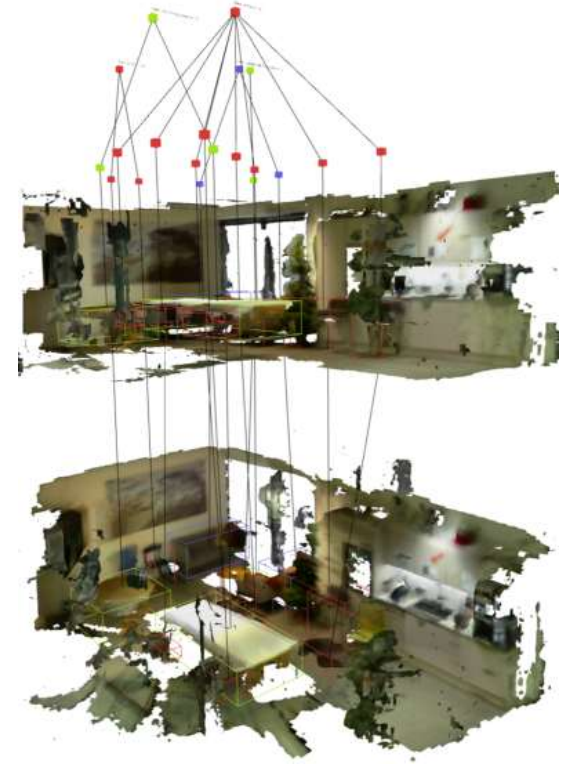
# Event-Grounding Graphs (EGG)

# Takeaways

A!

# Semantic mapping is evolving rapidly



2020

2023

2025

A!

# Trends



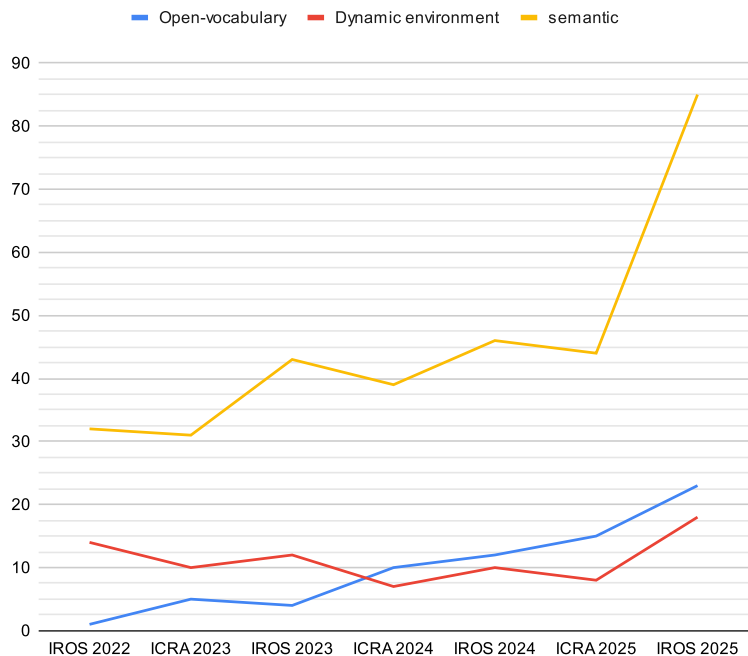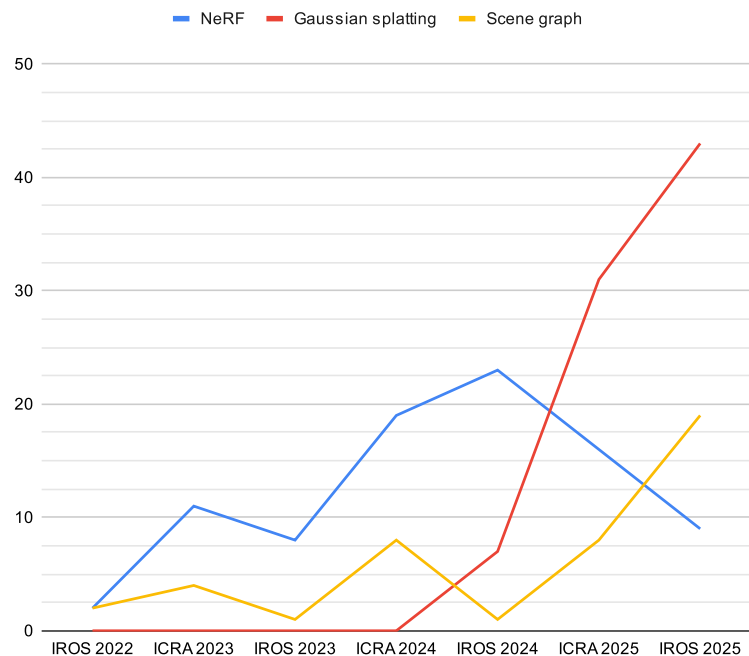**Trends in foundation models**

LLM — VLM — foundation model

**Trends in problems**

Open-vocabulary — Dynamic environment — semantic

**Trends in mapping**

NeRF — Gaussian splatting — Scene graph

A!

# Challenges and open problems

- **Lifelong mapping:** memory, forgetting, and map aging

- **Domain shift and generalization:** self-supervised and foundation models, VLMs, LLMs...

- **Multi-robot semantic mapping and map merging**

- **Task-specific maps:** sub-graph selection, planning domain generation

- **Dynamic scenes:** moving objects, time-dependent scene graphs

**A!**

# Thank you

**Francesco Verdoja**

fverdoja.github.io

A!