



ISTITUTO ITALIANO
DI TECNOLOGIA
PATTERN ANALYSIS
AND COMPUTER VISION

Semantic 3D Scene Understanding: from 3D Point Cloud to Graph representations

Alessio Del Bue

Fondazione Istituto Italiano di Tecnologia (IIT)

September 2 2023

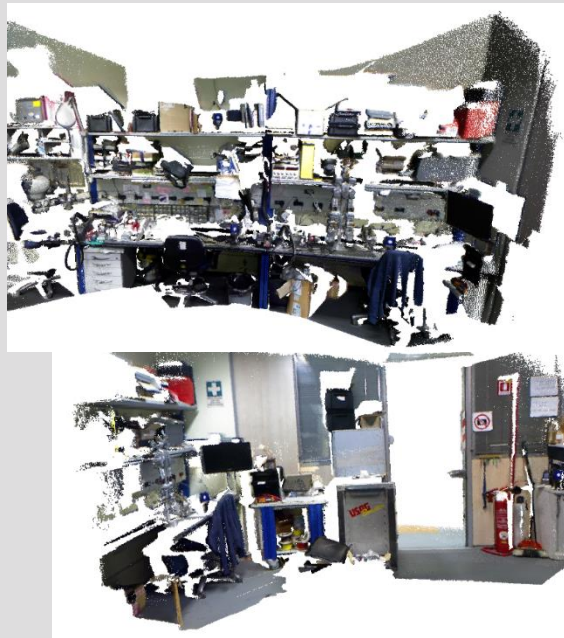
Padova, Italy - September 4-8, 2023

VIISMAC 23

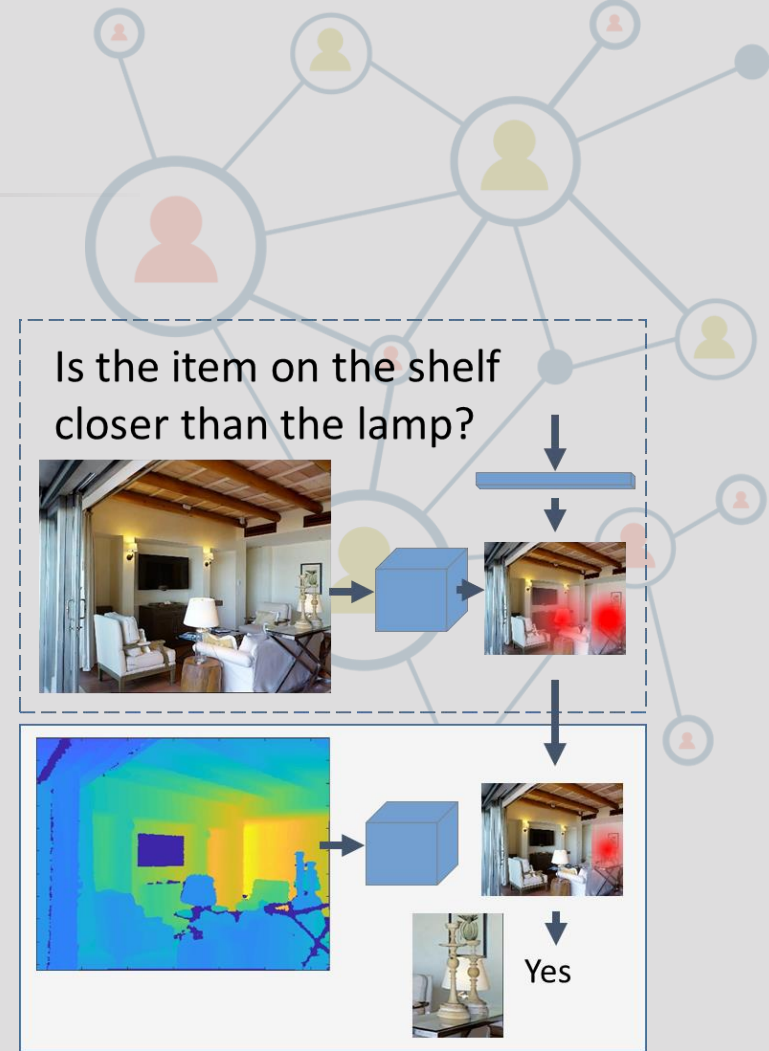
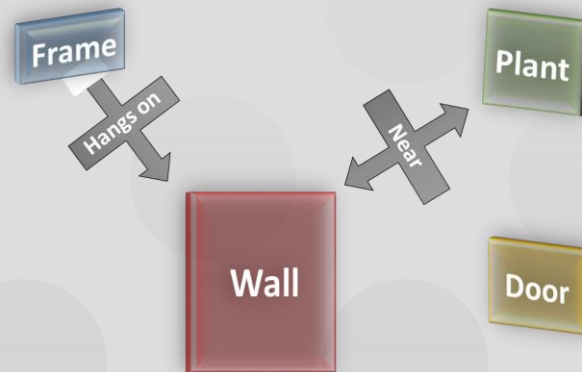
**International Summer School on
Machine Vision**



3D scene understanding



INPUT: 2D object detections in multiple views



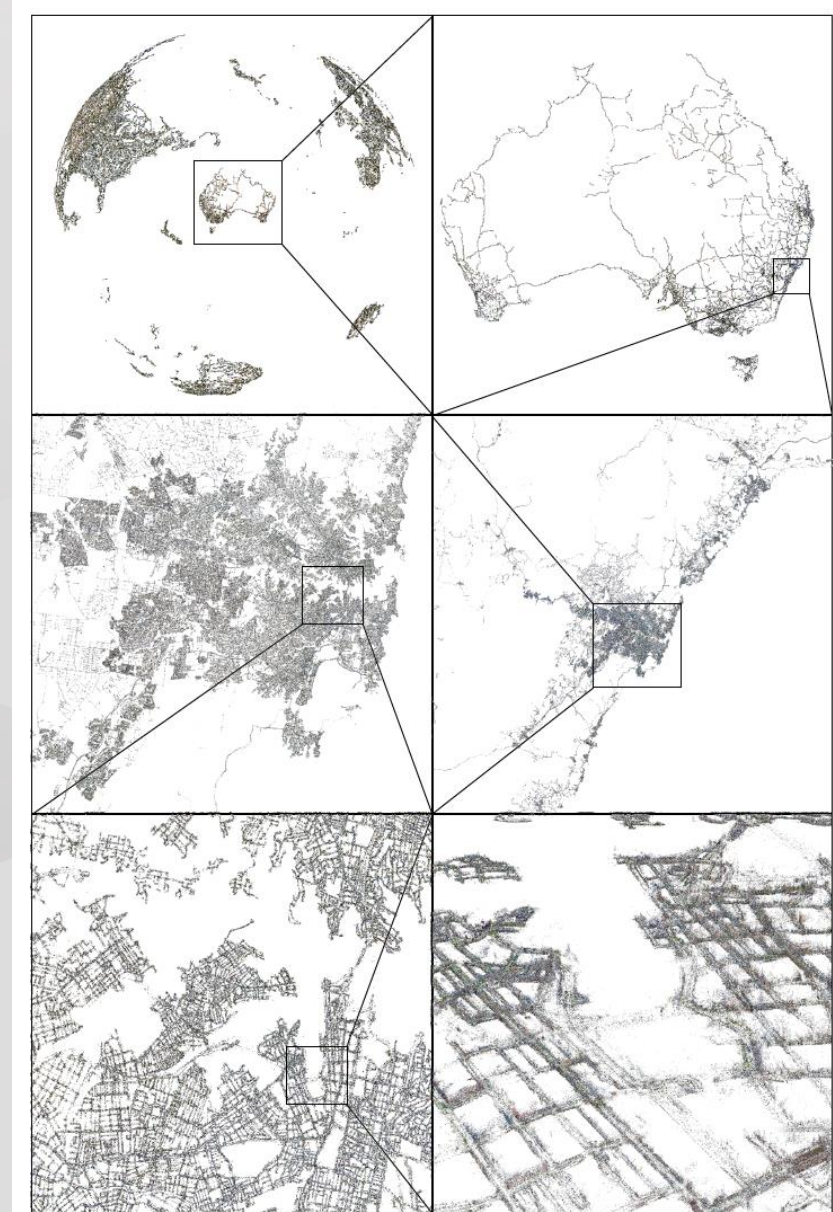
DIGITISATION &
3D MODELLING

REPRESENTATION
(RIGID AND DYNAMIC)

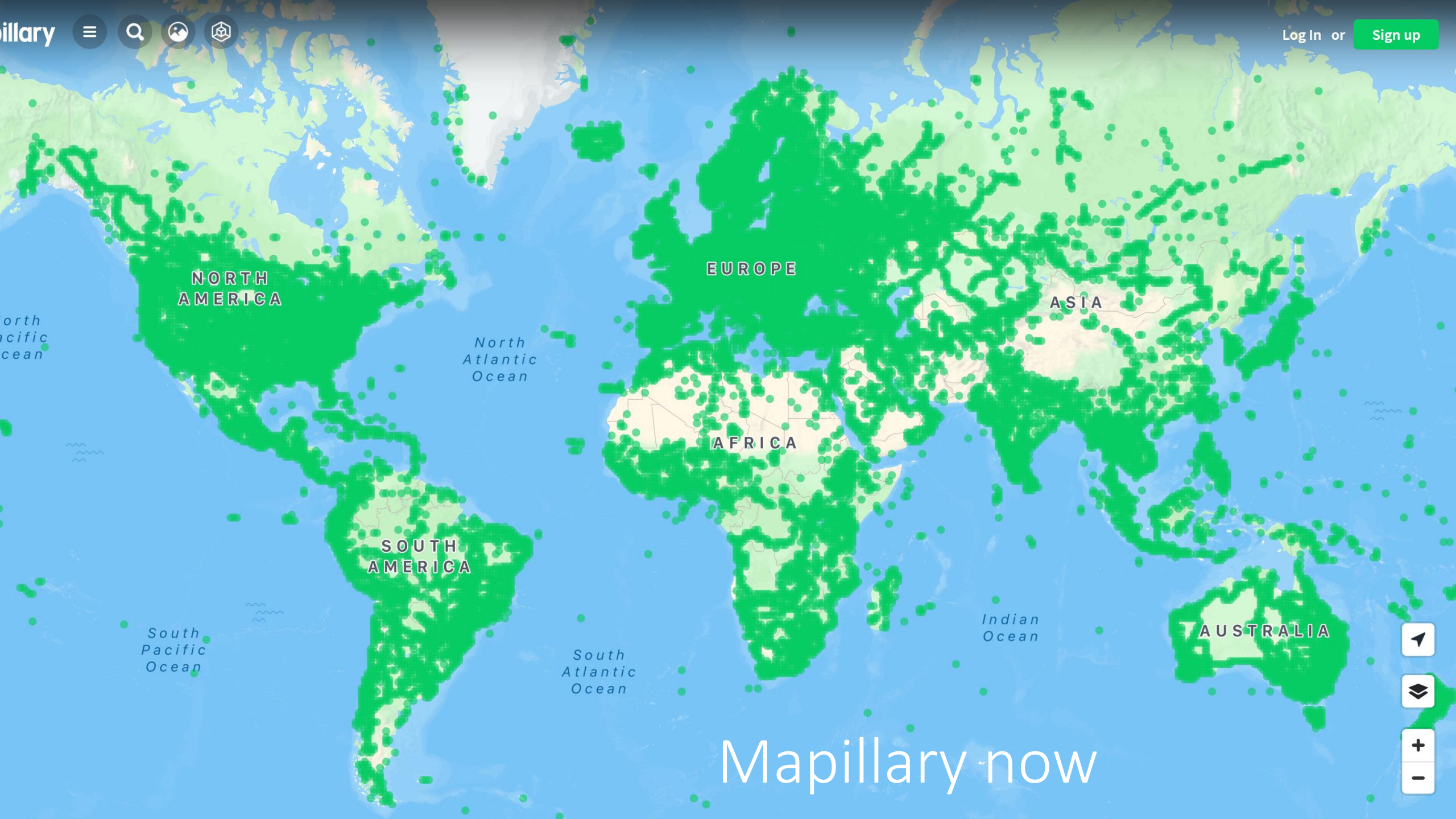
INTERPRETATION &
UNDERSTANDING

Digitisation and 3D modelling so far...

- An appearance-augmented point cloud comprising 404 billion tracked features, computed from Google street-level imagery.
- Every point in the cloud carries its local appearance descriptors from at least three different viewpoints.



Klingner, Bryan, David Martin, and James Roseborough. "Street view motion-from-structure-from-motion." CVPR 2013.



NORTH AMERICA

EUROPE

ASIA

AFRICA

SOUTH AMERICA

AUSTRALIA

North Pacific Ocean

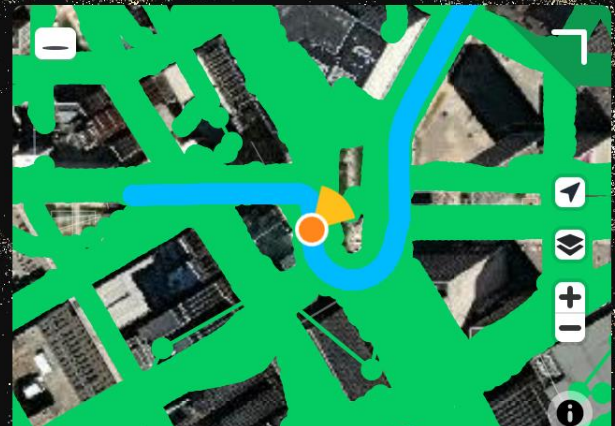
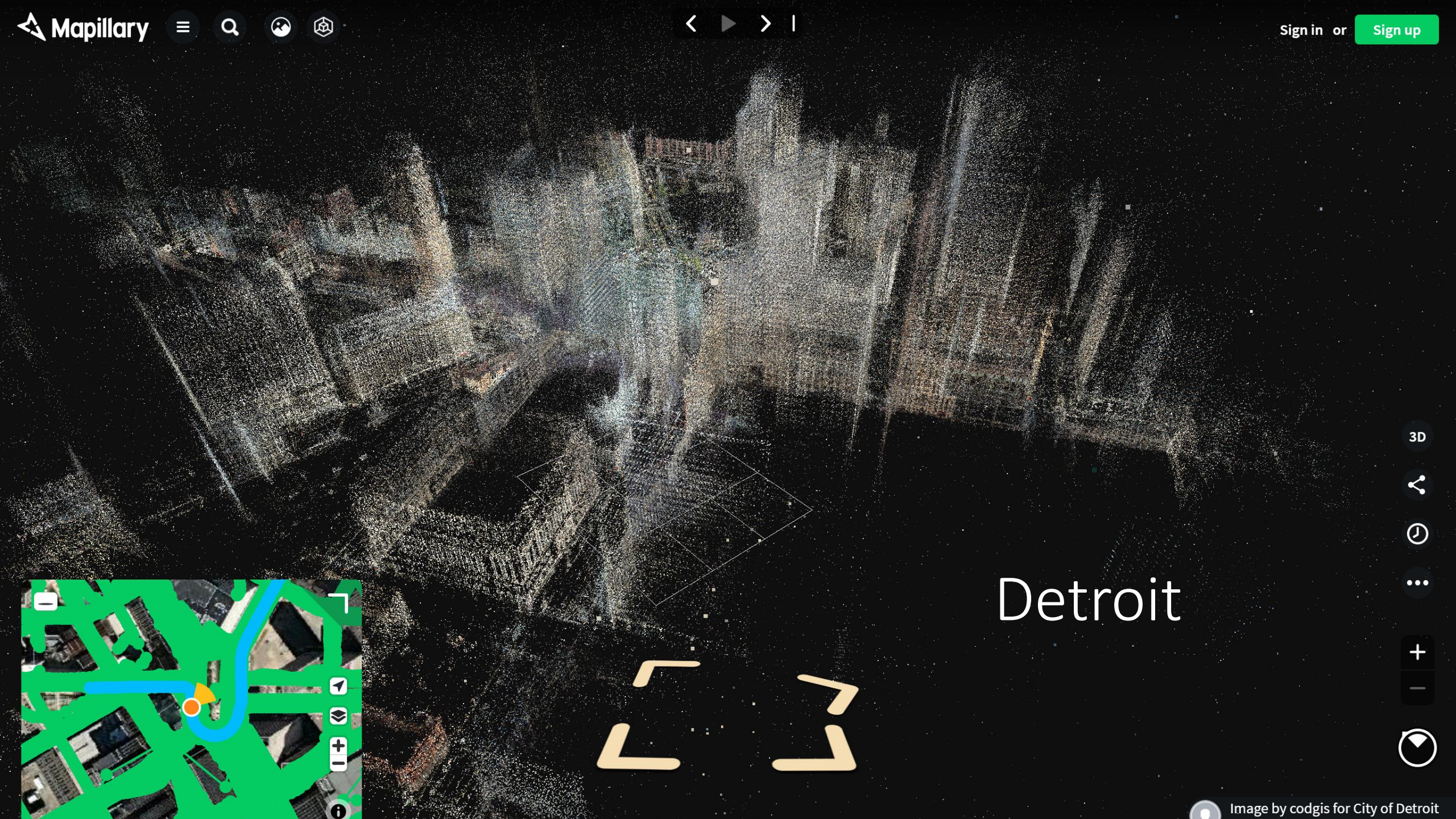
North Atlantic Ocean

South Pacific Ocean

South Atlantic Ocean

Indian Ocean

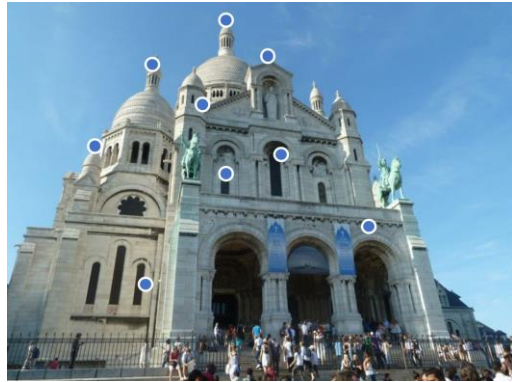
Mapillary now



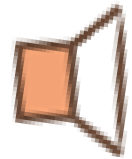
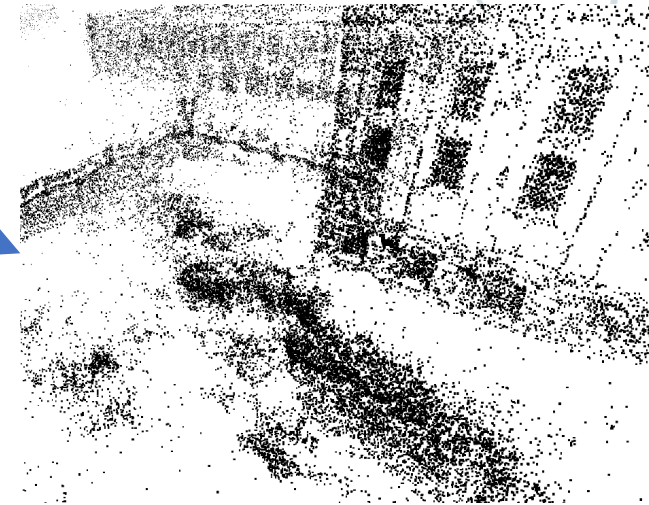
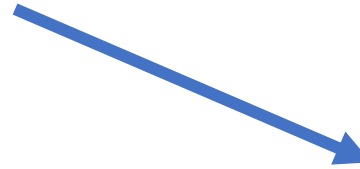
Detroit

- 3D
- 
- 
- 
- 
- 
- 

Multi-view 3D model (projective)



$$= K [R | t] \begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix}$$



3 x 3 matrix with intrinsic parameters (focal length, aspect ratio, principal point)



3 x 4 matrix of extrinsic parameters (camera rotation and translation)

4 x 1 homogeneous coordinates of the 3D point



Camera projection

From projective geometry to image plane coordinates (non-linear operator) $\mathbb{R}^3 \mapsto \mathbb{R}^2$

$$\Pi \begin{pmatrix} u \\ v \\ \lambda \end{pmatrix} = \begin{pmatrix} \frac{u}{\lambda} \\ \frac{v}{\lambda} \end{pmatrix} \Rightarrow \begin{pmatrix} u' \\ v' \end{pmatrix} = \Pi \left(K [R | t] \begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix} \right)$$

$$\begin{pmatrix} u \\ v \\ \lambda \end{pmatrix} = P \begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix}$$

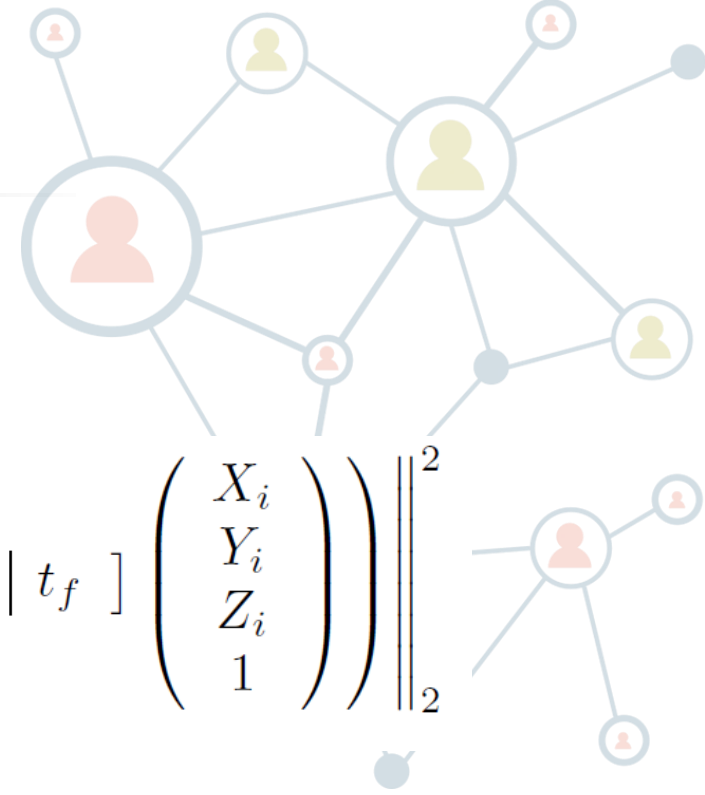
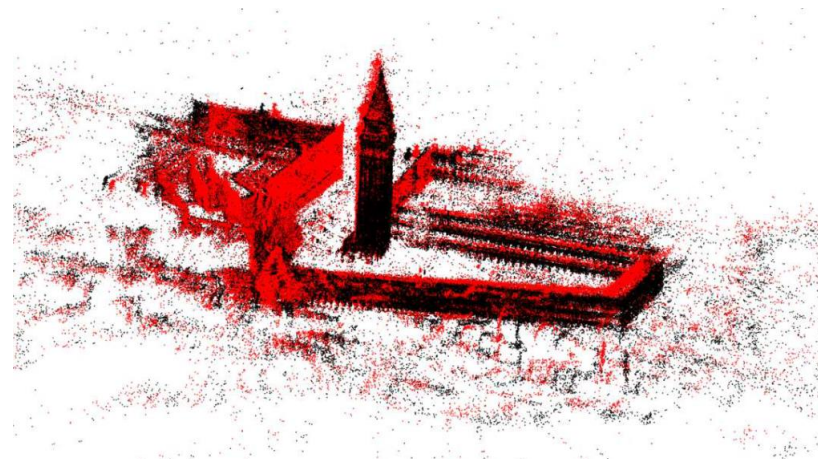
Other option is to estimate 3D structure and motion in the projective space and then upgrade later the reconstruction to metric.



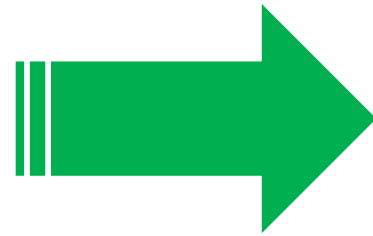
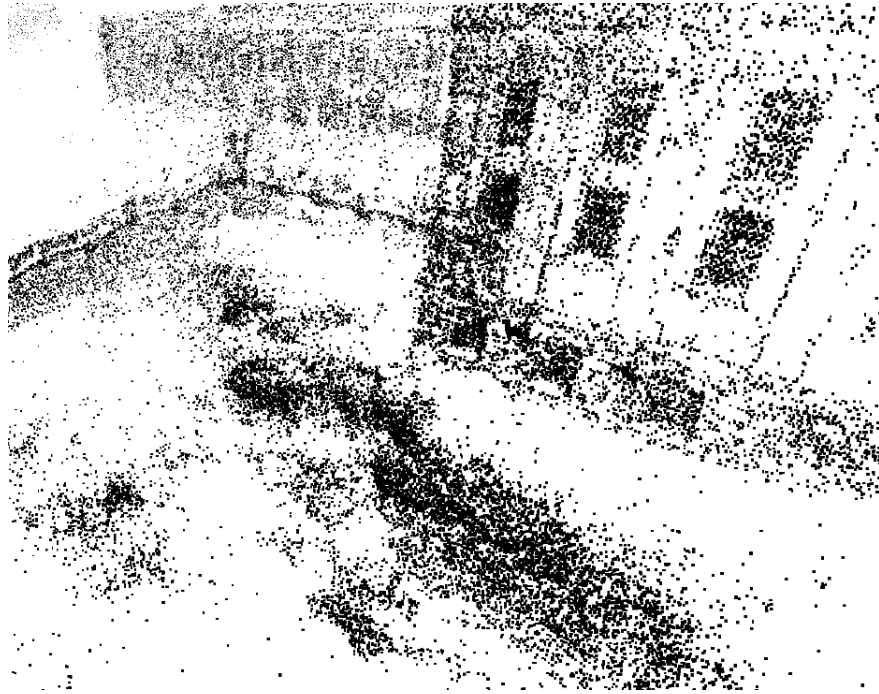
Cost (loss) function

For every 3D point i and every camera f optimise for the model parameters:

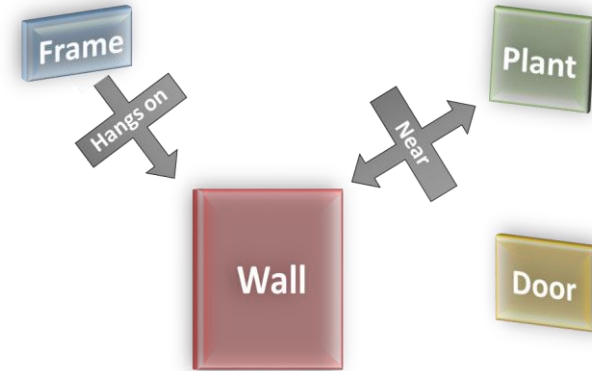
$$L_{\theta} = \sum_{i,f} \left\| \begin{pmatrix} u'_{i,f} \\ v'_{i,f} \end{pmatrix} - \Pi \left(K_f \begin{bmatrix} R_f & | & t_f \end{bmatrix} \begin{pmatrix} X_i \\ Y_i \\ Z_i \\ 1 \end{pmatrix} \right) \right\|_2^2$$



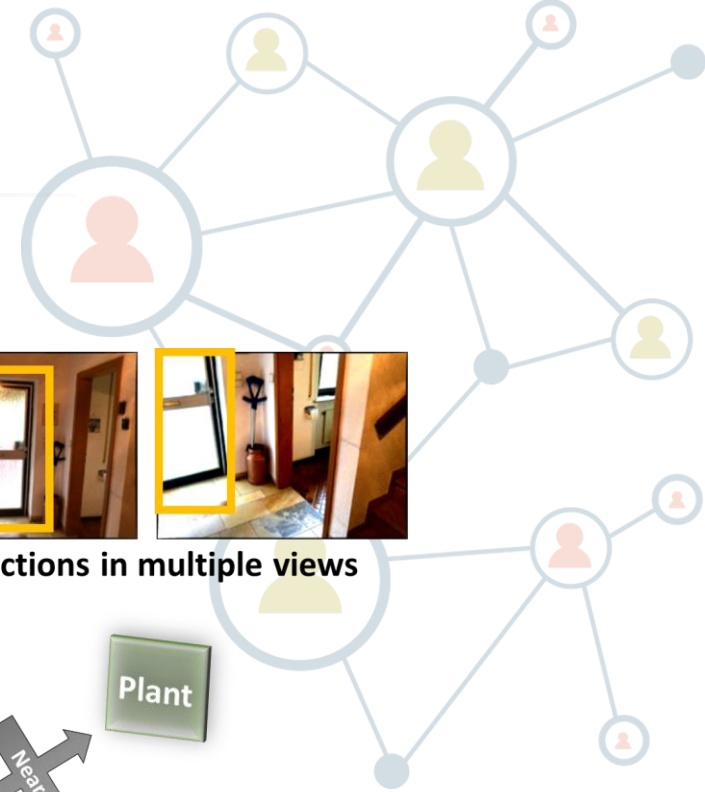
From point cloud to semantic



INPUT: 2D object detections in multiple views



REPRESENTATION
(RIGID AND DYNAMIC)



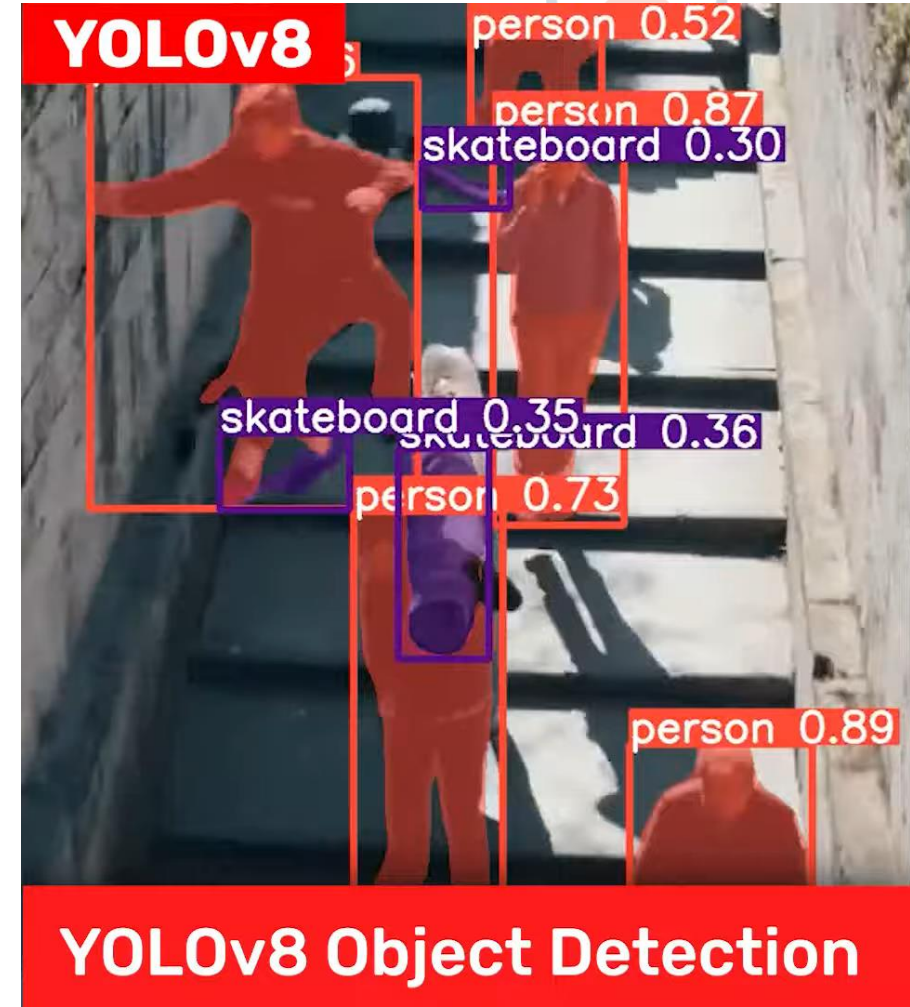
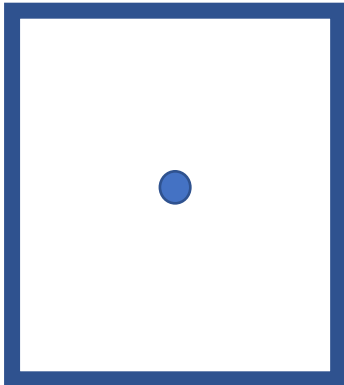
Observation: from points to detections



Source: Nvidia <https://nvidia.ws/2mCdYgx>

Semantic Structure from Motion

- Bounding boxes (bbx) are geometrical instances but they convey semantics
- Degenerate bbx are just image points as in SfM

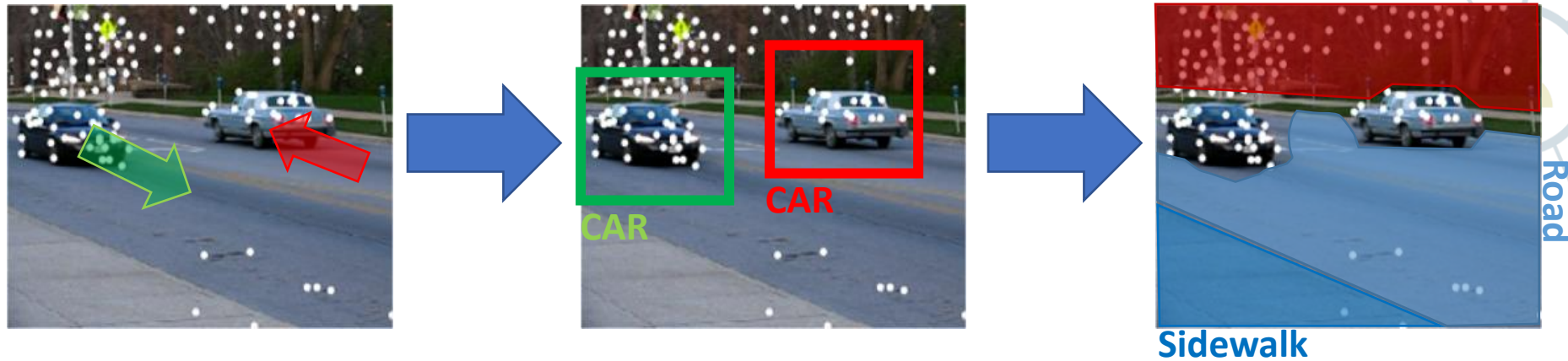


Redmon, Joseph, and Ali Farhadi. "Yolov3: An incremental improvement." *arXiv preprint* 2018.

YOLOv8: <https://github.com/ultralytics/ultralytics>

What can we gain from semantic?

How **high level semantic in 3D reconstruction** can help?



Generic scene

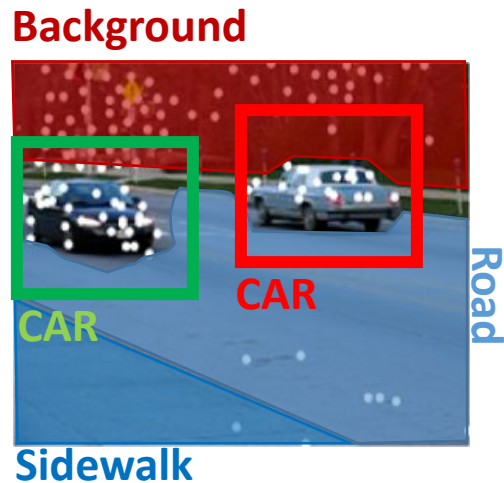
Points to objects
relations

Points to regions
relations

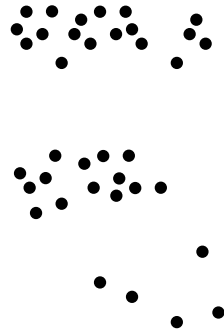
Including these relations **may help** to solve classical Motion Segmentation and Structure from Motion (SfM) problems.

How to bridge the gap?

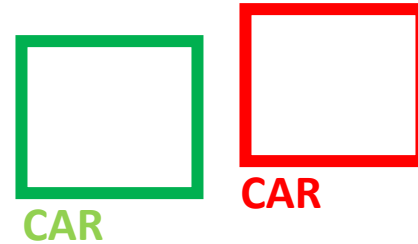
Different data that lives in different spaces (points, bbx, labels, etc.)



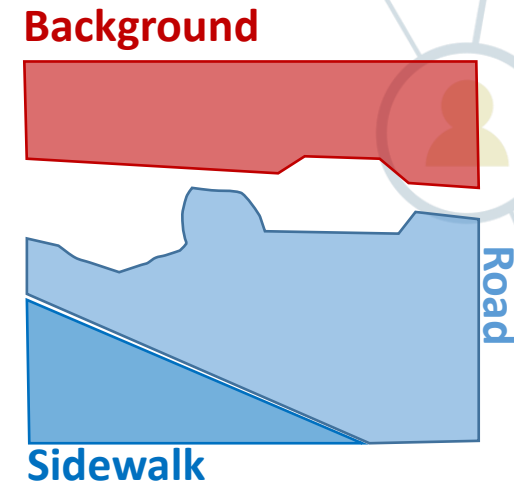
The scene



2D points



Bounding Boxes



Regions&Labels

How to blend this information available from deep learning in a single computational approach? Is there a viable solution (i.e. efficient)?

Now semantic is ready to use



Genoa Porto Antico – Instance segmentation



Image by [javedsial91](#) licensed under [CC-BY-SA](#).

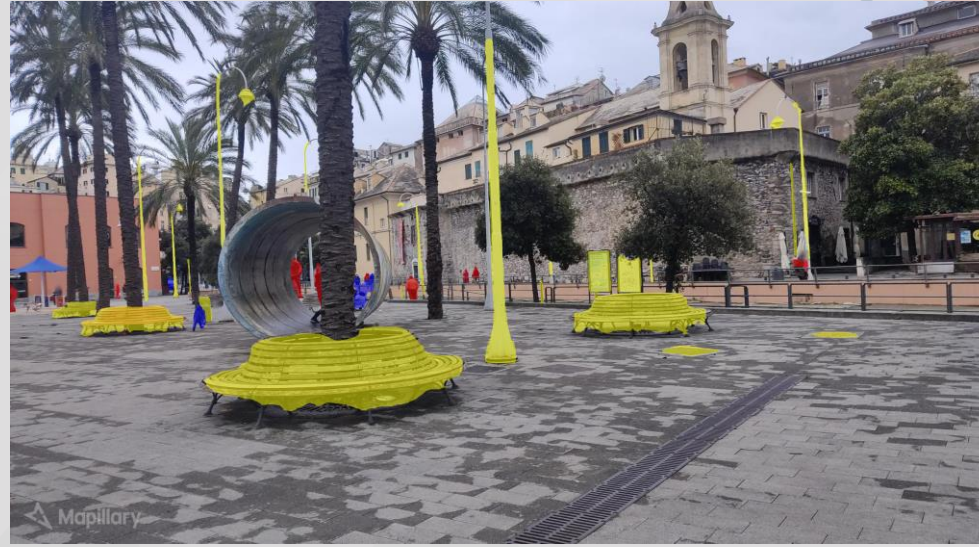


Image by [javedsial91](#) licensed under [CC-BY-SA](#).



Image by [javedsial91](#) licensed under [CC-BY-SA](#).

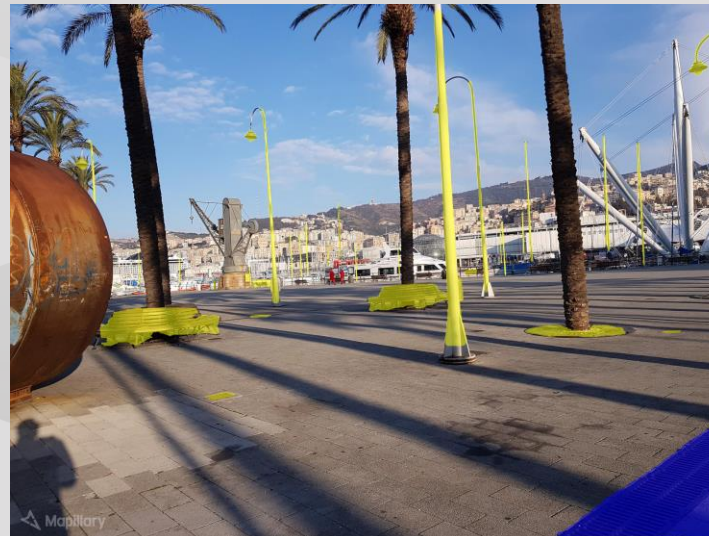
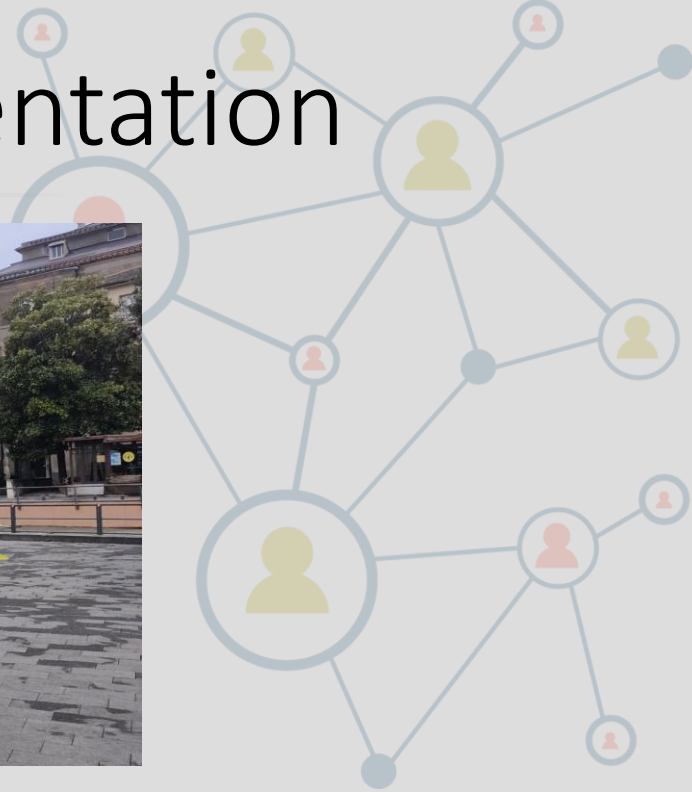


Image by [ale_zena_it](#) licensed under [CC-BY-SA](#).

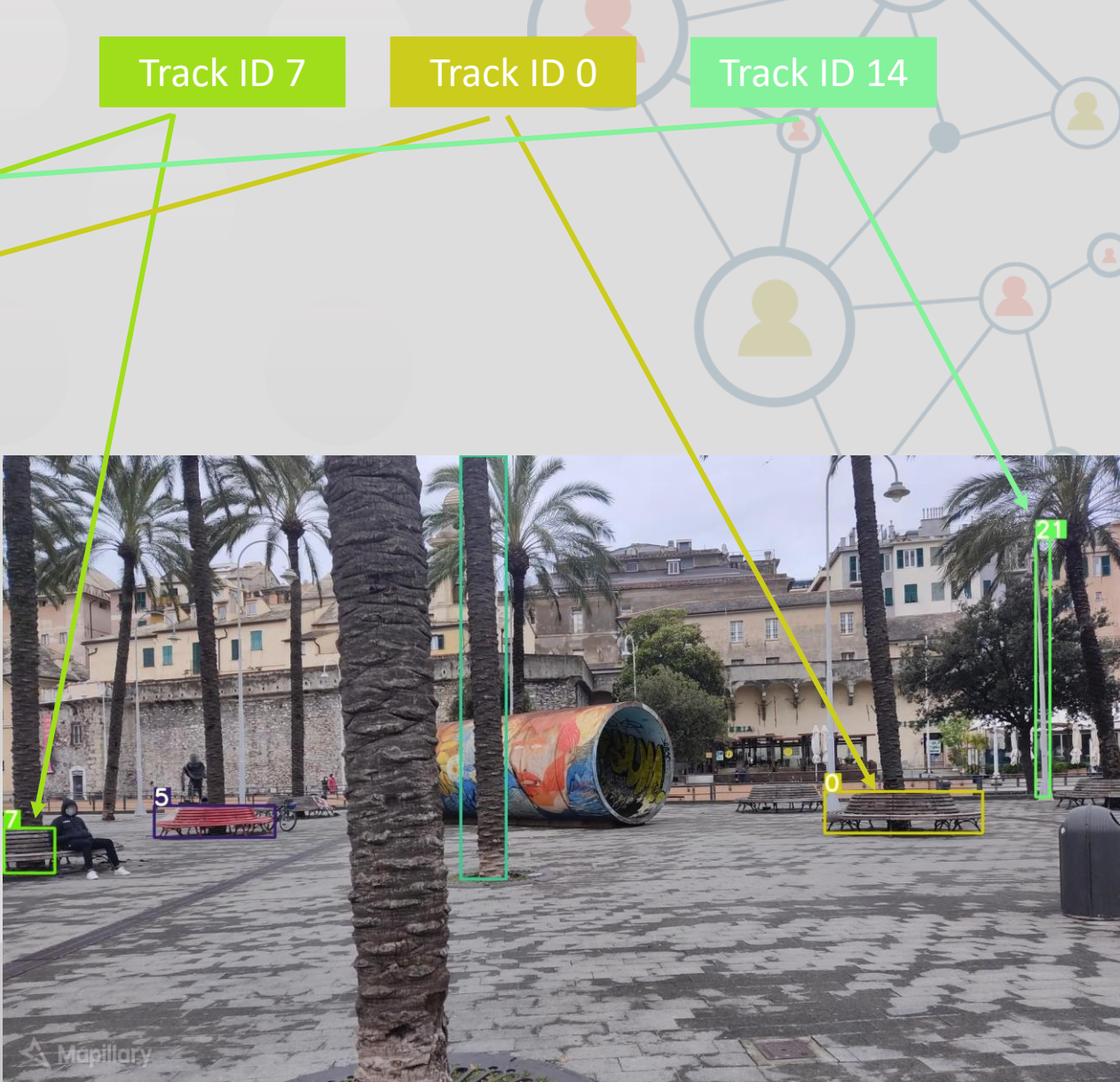


- Red – people, animals
- Blue – vehicles
- Yellow – static objects:
 - Benches
 - Poles
 - Manholes
 - Banners
 - Street lights
 - ...

Genoa Porto Antico – Matched detections

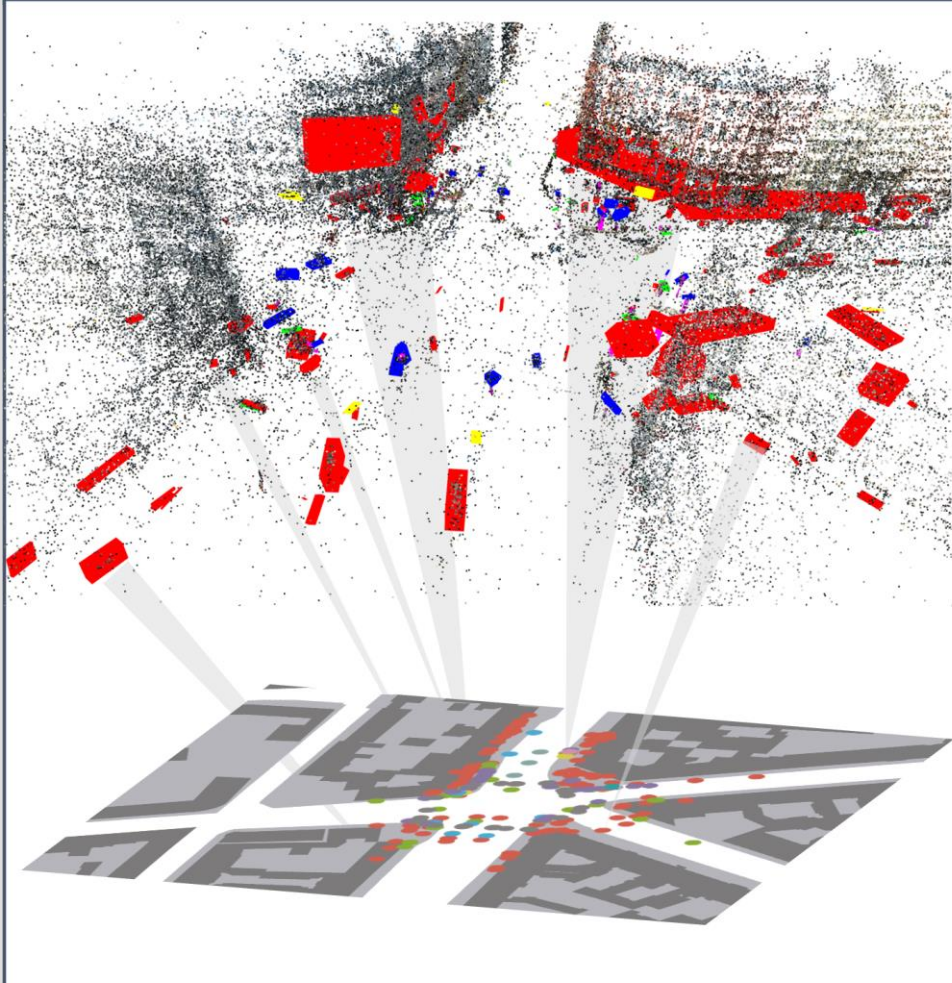


Image by [javedsial91](#) licensed under [CC-BY-SA](#).



Revisiting Semantic Structure from Motion

📍 Berlin



Crocco M, Rubino C, Del Bue A. Structure from motion with objects. CVPR 2016.

Rubino C, Crocco M, Del Bue A. 3d object localisation from multi-view image detections. TPAMI 2017

Gay P, Rubino C, Bansal V, Del Bue A. Probabilistic structure from motion with objects (psfmo). ICCV 2017.

Gay P, Stuart J, Del Bue A. Visual graphs from motion (vgfm): Scene understanding with object geometry reasoning. ACCV 2018.

Giuliari et al. Spatial commonsense graph for object localisation in partial scenes. CVPR 2022.

Taiana et al. "PoserNet: Refining Relative Camera Poses Exploiting Object Detections". ECCV 2022

Toso et al. You are here! Finding position and orientation on a 2D map from a single image: The Flatlandia localization problem and dataset. arXiv 2023.

Key

- | | | |
|------------------|-----------------|----------------|
| ● Arrow Marking | ● Streetlight | ● Bench |
| ● Bicycle Symbol | ● Support Pole | ● Bike Rack |
| ● Sign | ● Traffic Light | ● Catch Basin |
| ● Traffic-Sign | ● Manhole | ● CCTV Camera |
| ● Trash can | ● Junction Box | ● Fire Hydrant |

Including object detection in 3D vision

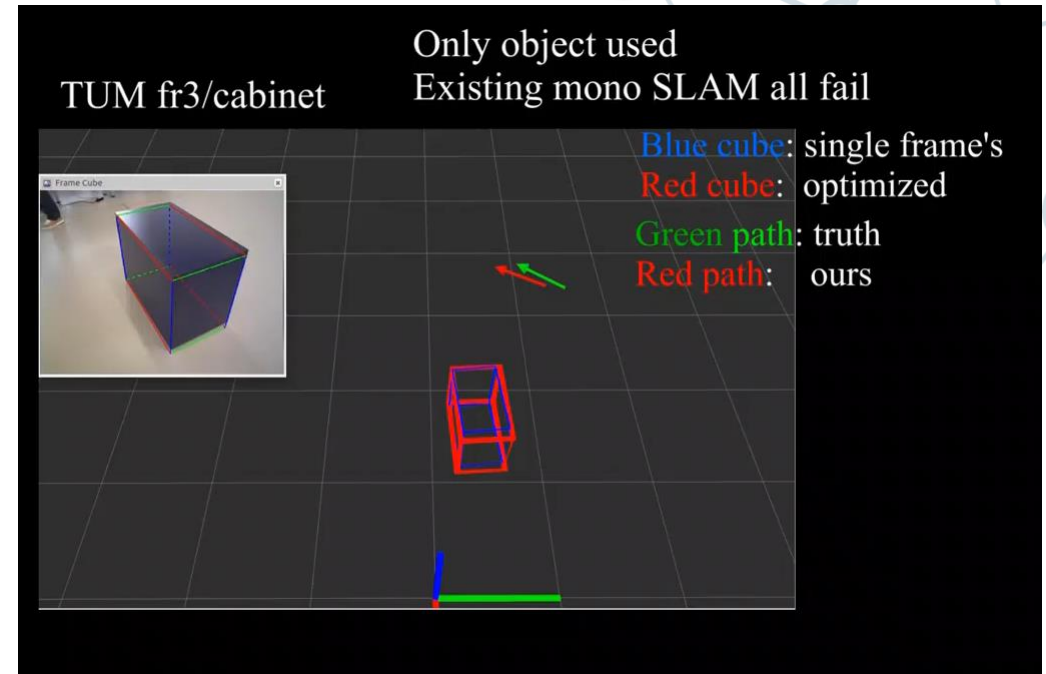
Real-time Monocular Object SLAM

Dorian Gálvez-López, Marta Salas, Juan D. Tardós, J. M. M. Montiel

*Robotics, Perception and Real Time Group
Departamento de Informática e Ingeniería de Sistemas
Instituto de Investigación en Ingeniería de Aragón
Universidad de Zaragoza, Spain*



Gálvez-López, Dorian, et al. "Real-time monocular object slam." *RAS* 2016



Yang, Shichao, and Sebastian Scherer. "Cubeslam: Monocular 3-d object slam." *TRO* 2019

Including object detection in 3D vision



QuadricSLAM: Constrained Dual Quadrics from Object Detections as Landmarks in Object-oriented SLAM

Lachlan Nicholson, Michael Milford, Niko Sünderhauf



ARC Centre of Excellence for Robotic Vision

Nicholson, Lachlan, Michael Milford, and Niko Sünderhauf. "Quadricslam: Dual quadrics from object detections as landmarks in object-oriented slam." *RA-L* 2018



Including object detection in 3D vision

OA-SLAM

(localization mode)

(the scene has been previously mapped and augmented)



ORB-SLAM2

(localization mode)



Camera tracking lost !!!

*all videos are shown at original speed

Matthieu Zins, Gilles Simon, Marie-Odile Berger. "OA-SLAM: Leveraging Objects for Camera Relocalization in Visual SLAM." *ISMAR 2022*

More examples on...

☰ README.md

Awesome Object SLAM

A curated list of Object SLAM papers and resources, inspired by [awesome-implicit-representations](#).

Disclaimer

This list *does not aim to be exhaustive*. In this list, we consider papers that **jointly** optimize robot (camera) and object states, where object states typically include object poses and object shape parameters.

For more general SLAM papers, please refer to [awesome-visual-SLAM](#) and [Awesome-SLAM](#).

This repo is maintained by [Ziqi Lu](#) and [Akash Sharma](#). You are very welcome to contribute to this repo. If you spot anything wrong or missing, please feel free to [submit a pull request](#) or contact the maintainers.

Table of Contents

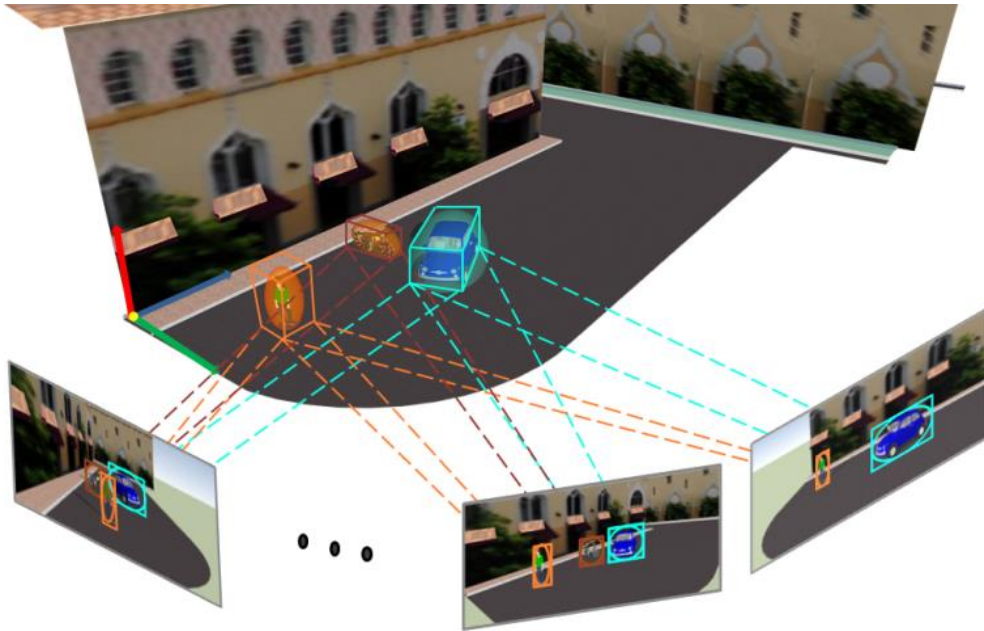
- [What is Object SLAM](#)
- [Papers](#)
 - [Parametric Object Representation](#)
 - [Field Object Representation](#)
 - [Inference Methods for Object SLAM](#)
 - [Reviews](#)
- [Resources](#)
 - [Datasets](#)



<https://github.com/520xyxyzq/awesome-object-SLAM>

Localisation from Detections (LfD)

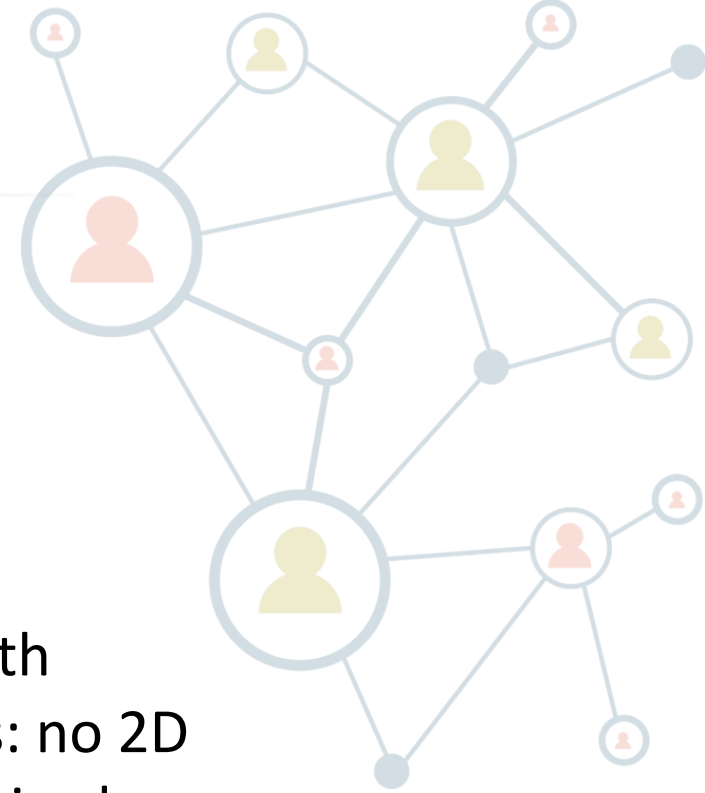
Main goal: to recover the **3D occupancy** of a set of objects given the **2D bounding boxes** from detections at each image frame.



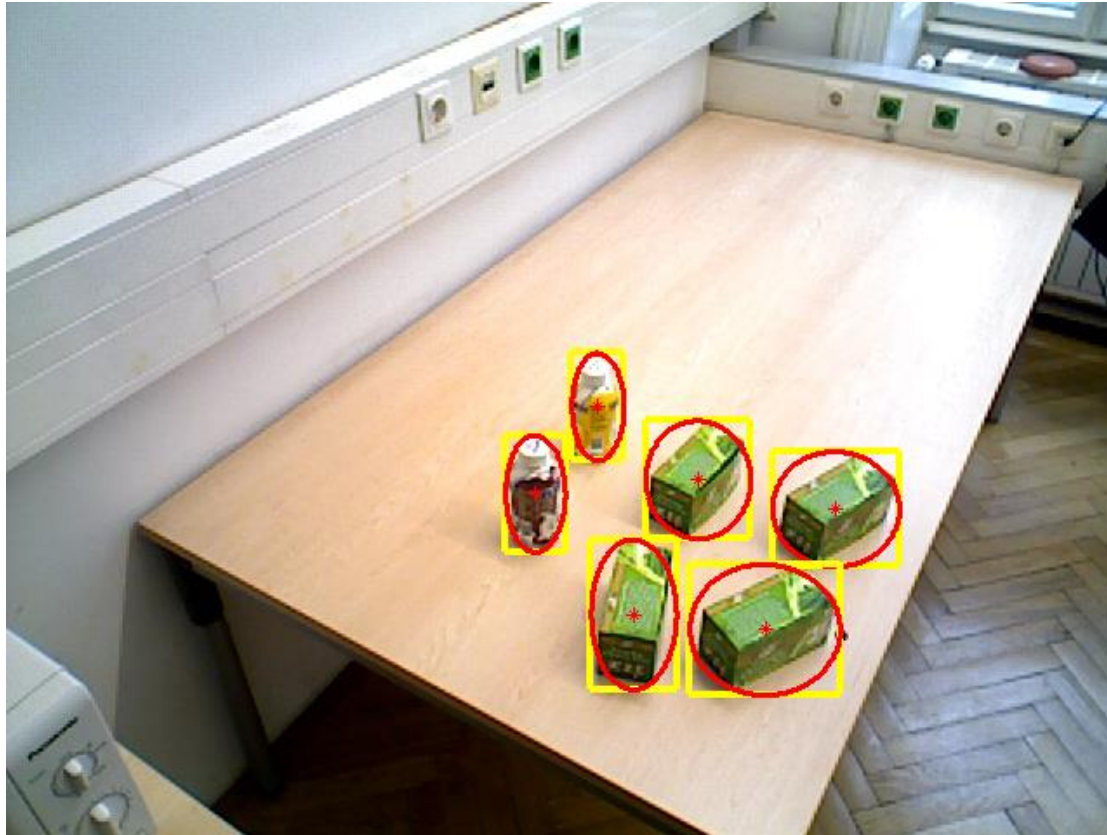
Features:

- Method working with **texture-less** objects: no **2D interest points** required
- Just 2D bounding boxes position, size and aspect ratio are required

Rubino, Cosimo, Marco Crocco, and Alessio Del Bue. "3d object localisation from multi-view image detections." *TPAMI* 2017.



3D ellipsoids from bounding boxes



3D bounding box from set of **2D bounding boxes** (one per frame)

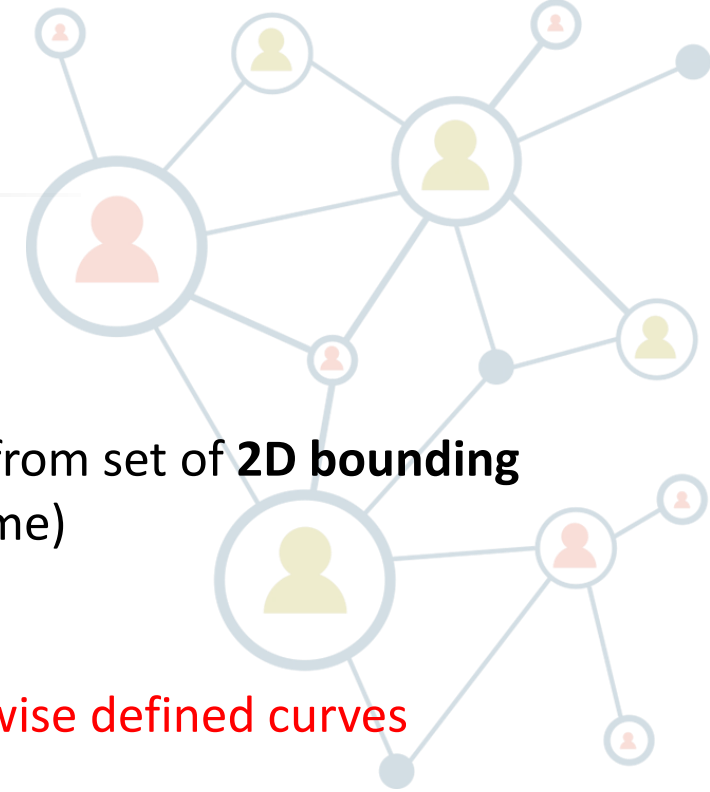


Not simple: piece-wise defined curves

3D ellipsoid from set of **2D ellipses** fitted to 2D bounding box

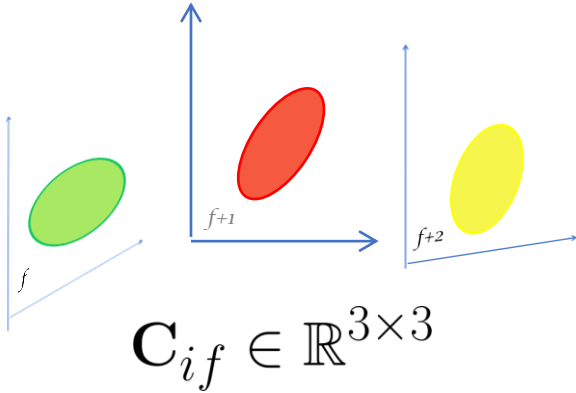


Algebraic solution available in closed form



Multiple conic to quadric reconstruction (1)

Conics from 2D bounding boxes at frames $1 \dots F$ related to object i

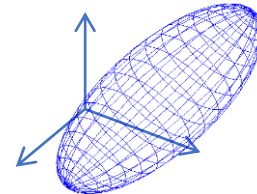


Projection matrices for frames $1 \dots F$: can be estimated from the image background by standard multi-view relations

$$\mathbf{P}_f \in \mathbb{R}^{3 \times 4}$$

$\mathbf{Q}_i \in \mathbb{R}^{4 \times 4}$

Estimated Quadric (ellipsoid) representing 3D space occupancy of object i

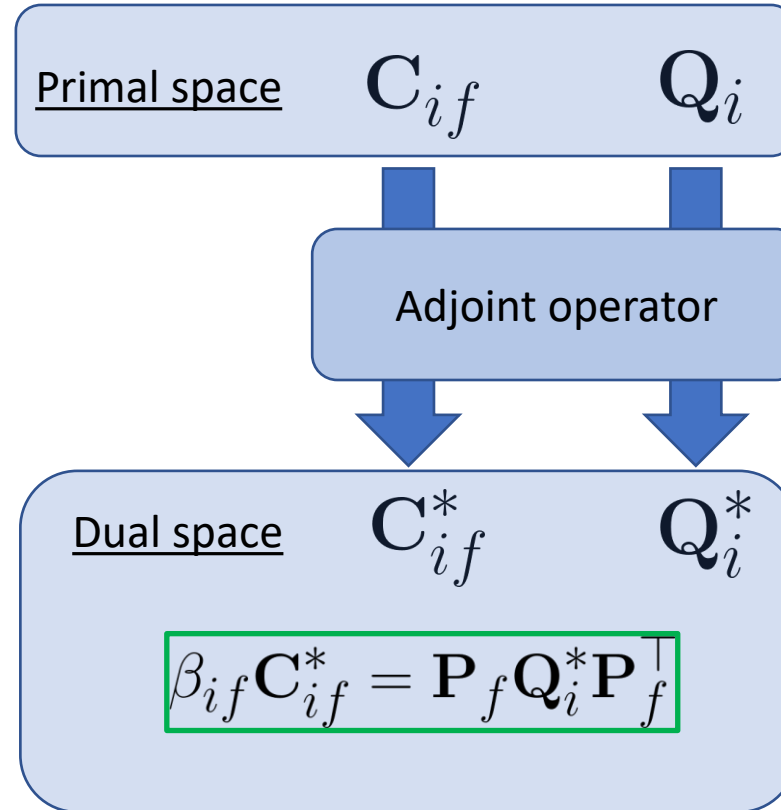


- [1] De Ma, S., Chen, X.: Reconstruction of quadric surface from occluding contour. In: Pattern Recognition, 1994. Vol. 1-Conference A: Computer Vision & Image Processing., Proceedings of the 12th IAPR International Conference on. Volume 1., IEEE (1994) 27–31
- [2] G. Cross and A. Zisserman. Quadric reconstruction from dual-space geometry. ICCV 1998.

Multiple conic to quadric reconstruction (2)

Space of points in 2D and 3D

Space of lines and planes

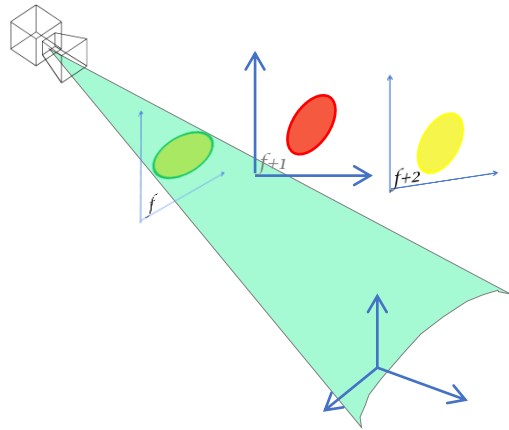


No simple relation between quadric and conics in primal space

Linear relation between quadric and conics in dual space



Multiple conic to quadric reconstruction (3)



$$\beta_{if} \mathbf{C}_{if}^* = \mathbf{P}_f \mathbf{Q}_i^* \mathbf{P}_f^T$$

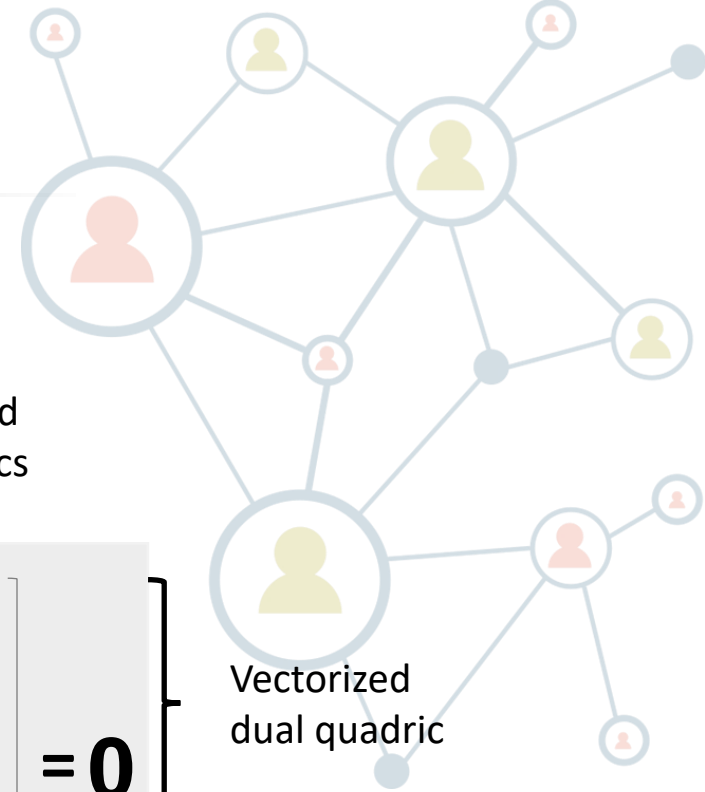
Vectorize dual conics and quadrics and stack over f

Projection matrices Vectorized dual conics

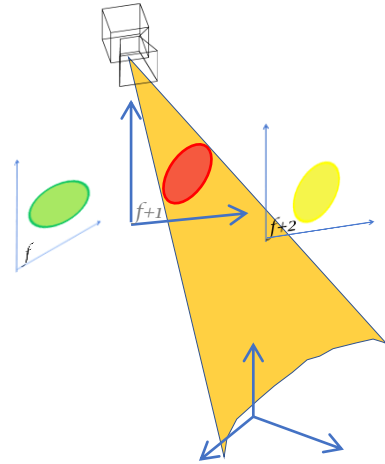
$$\begin{bmatrix} \text{Projection matrices} & \text{Vectorized dual conics} \\ \vdots & \vdots \end{bmatrix} = \mathbf{0}$$

Scale factors β_{if}

$$\mathbf{M}_i \mathbf{w}_i = \mathbf{0}_{6F}$$

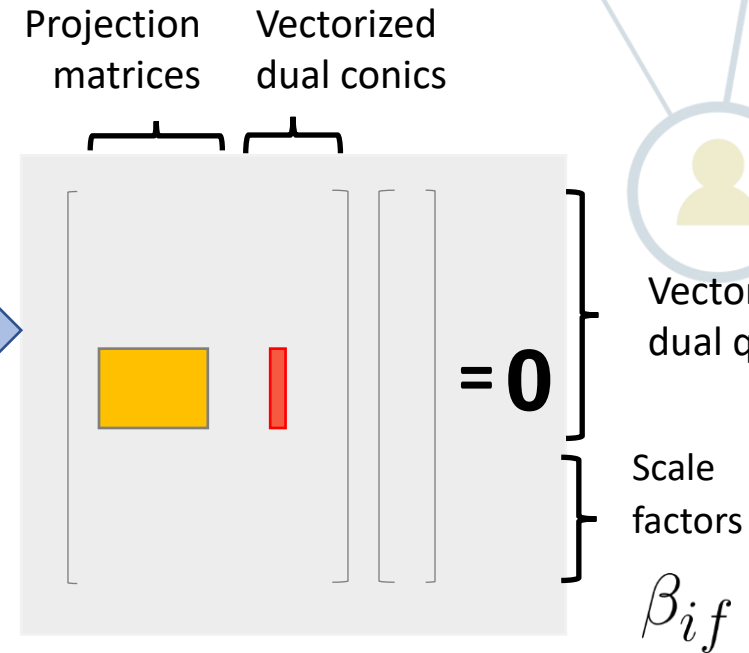


Multiple conic to quadric reconstruction (3)

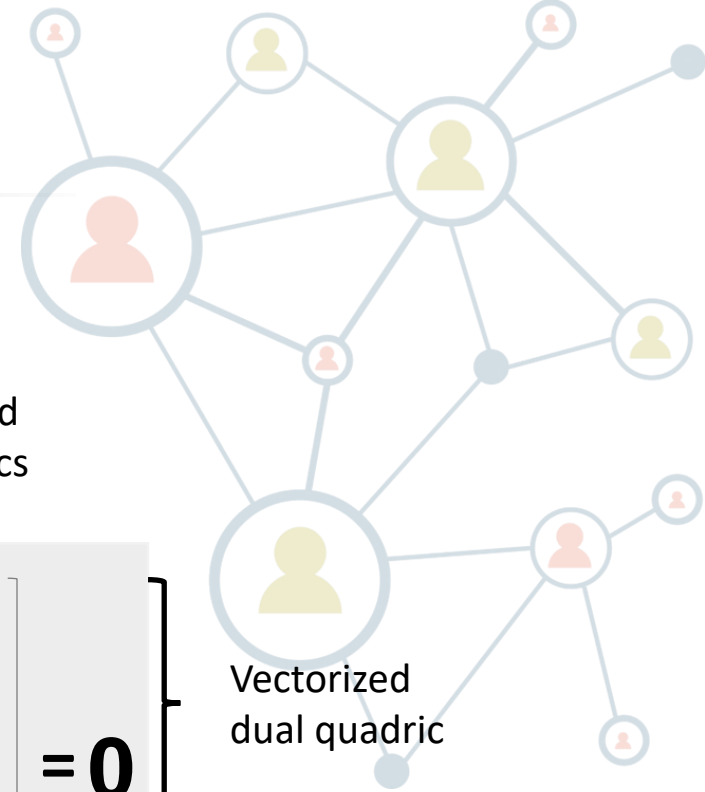


$$\beta_{if} \mathbf{C}_{if}^* = \mathbf{P}_f \mathbf{Q}_i^* \mathbf{P}_f^T$$

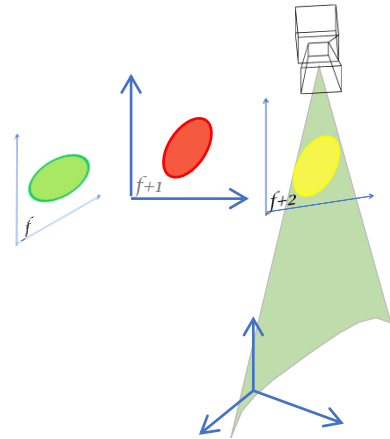
Vectorize dual conics and quadrics and stack over f



$$\mathbf{M}_i \mathbf{w}_i = \mathbf{0}_{6F}$$



Multiple conic to quadric reconstruction (3)



$$\beta_{if} \mathbf{C}_{if}^* = \mathbf{P}_f \mathbf{Q}_i^* \mathbf{P}_f^T$$

Vectorize dual conics and quadrics and stack over f

Projection matrices Vectorized dual conics

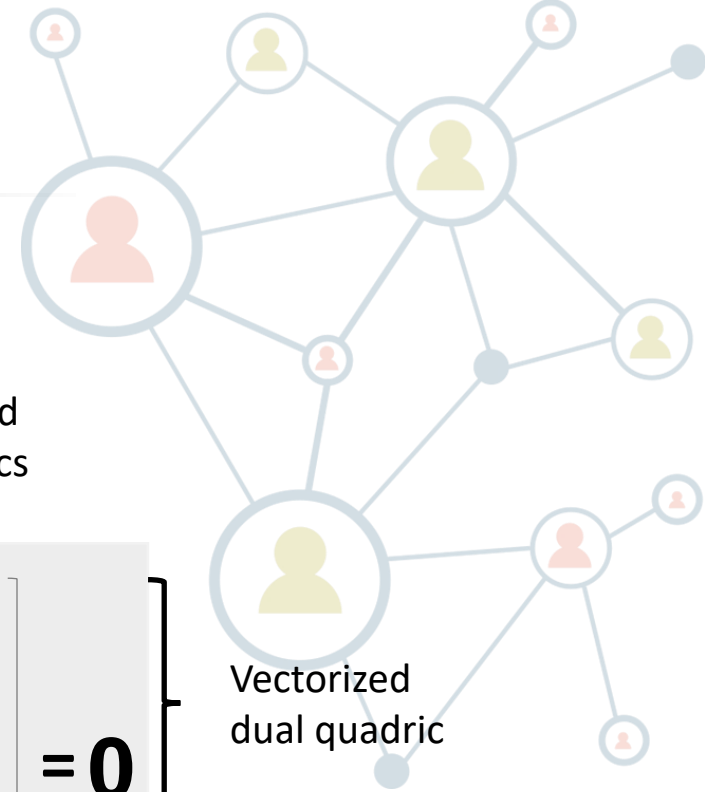
$$\begin{bmatrix} \text{Projection matrices} & \text{Vectorized dual conics} \\ \hline \text{Scale factors} \end{bmatrix} = \mathbf{0}$$

Vectorized dual quadric

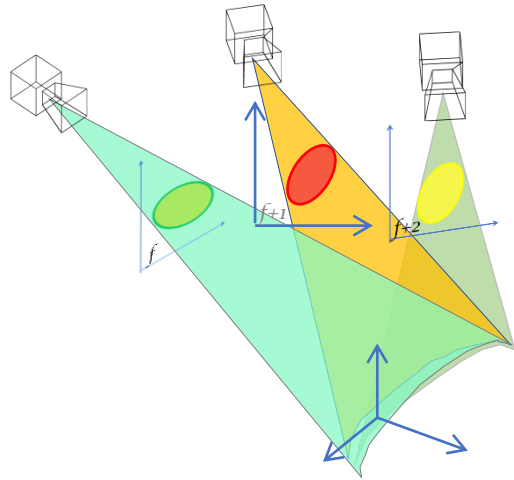
Scale factors

β_{if}

$$\mathbf{M}_i \mathbf{w}_i = \mathbf{0}_{6F}$$

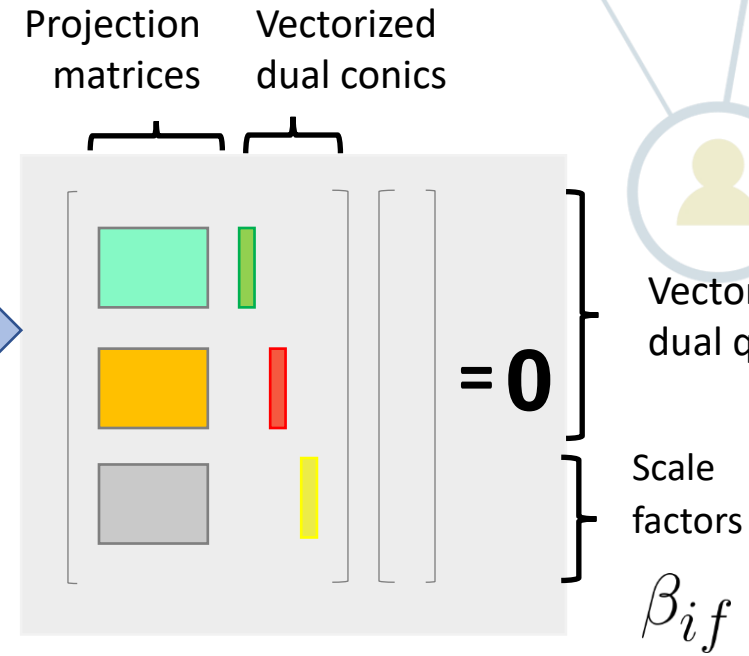


Multiple conic to quadric reconstruction (3)

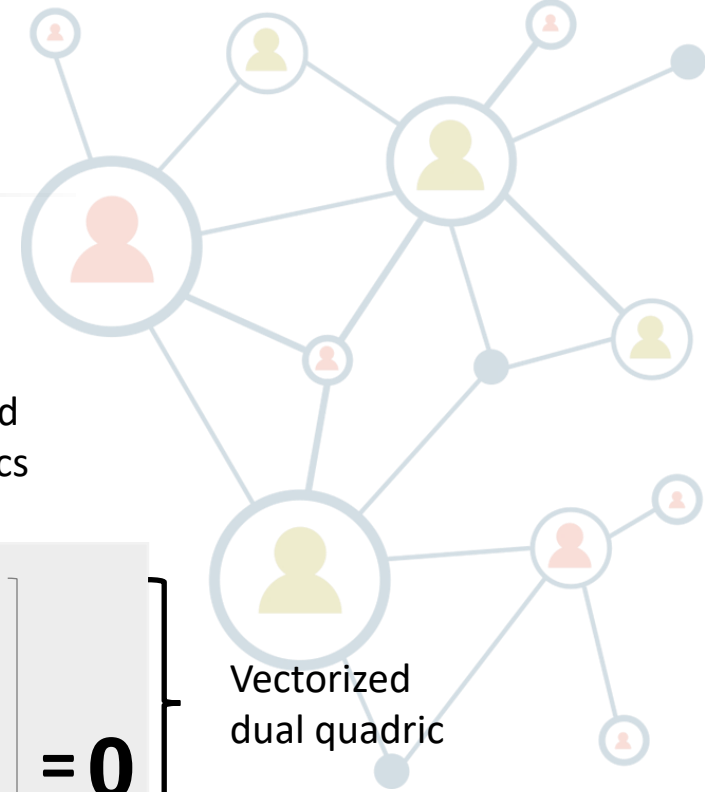


$$\beta_{if} \mathbf{C}_{if}^* = \mathbf{P}_f \mathbf{Q}_i^* \mathbf{P}_f^T$$

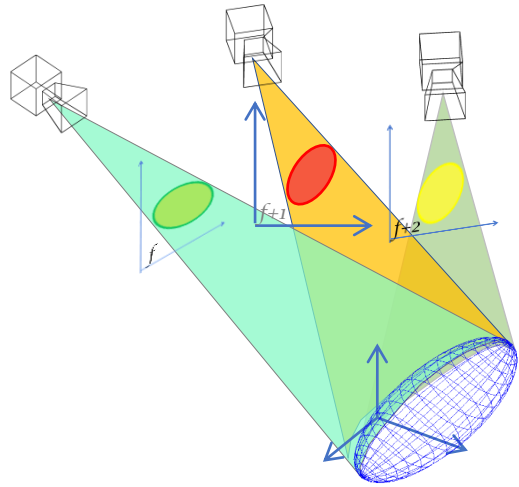
Vectorize dual conics and quadrics and stack over f



$$\mathbf{M}_i \mathbf{w}_i = \mathbf{0}_{6F}$$

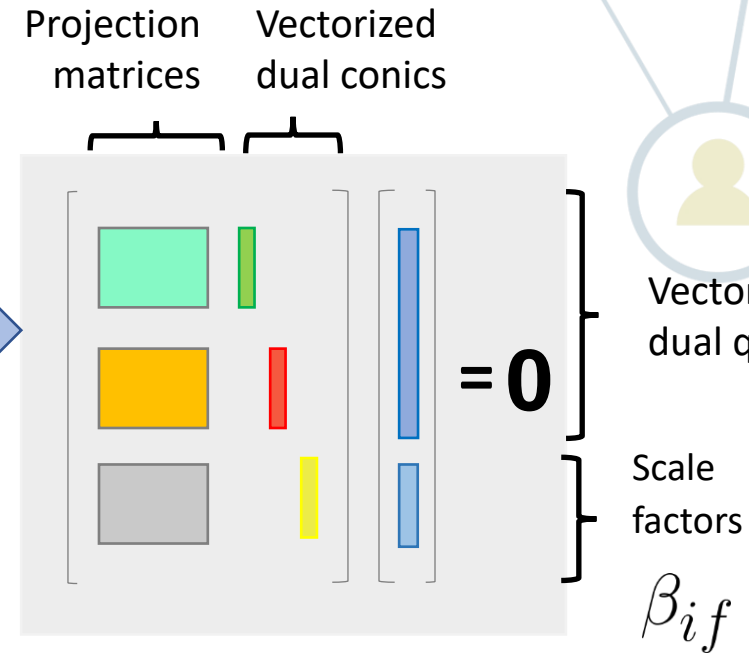


Multiple conic to quadric reconstruction (3)

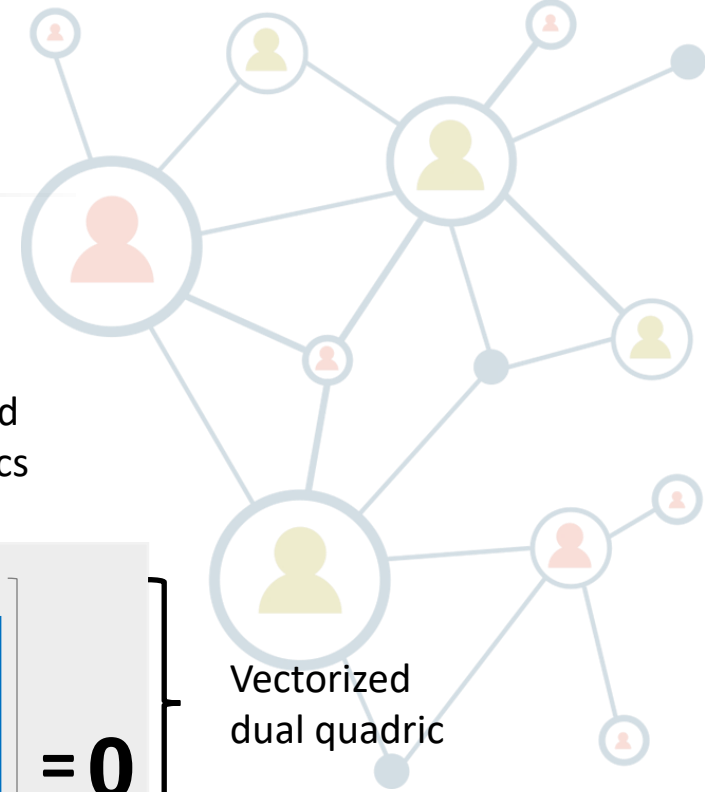


$$\beta_{if} \mathbf{C}_{if}^* = \mathbf{P}_f \mathbf{Q}_i^* \mathbf{P}_f^T$$

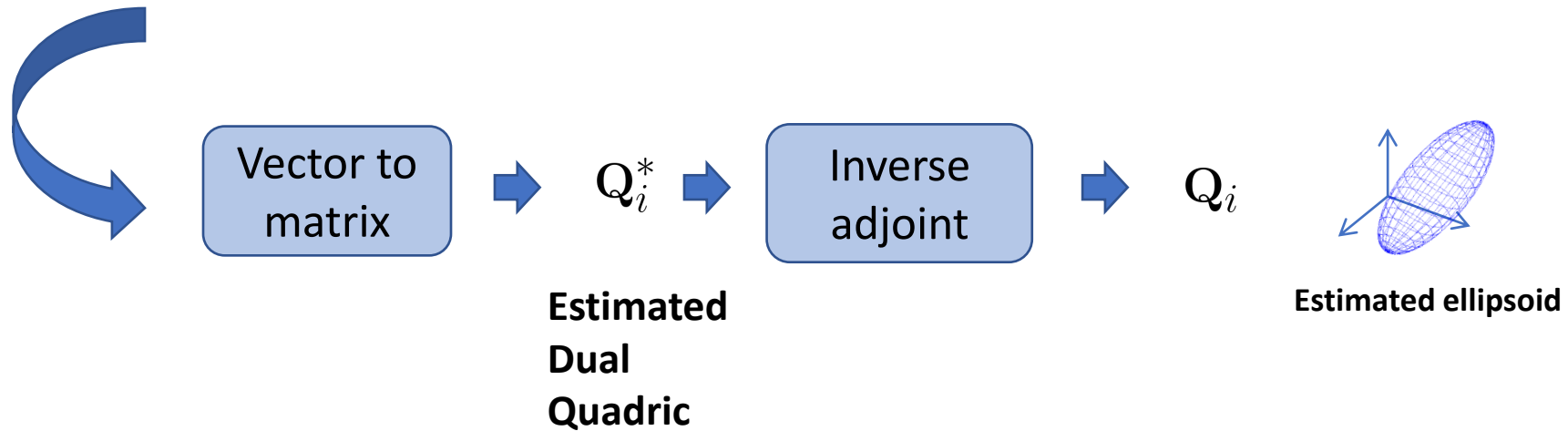
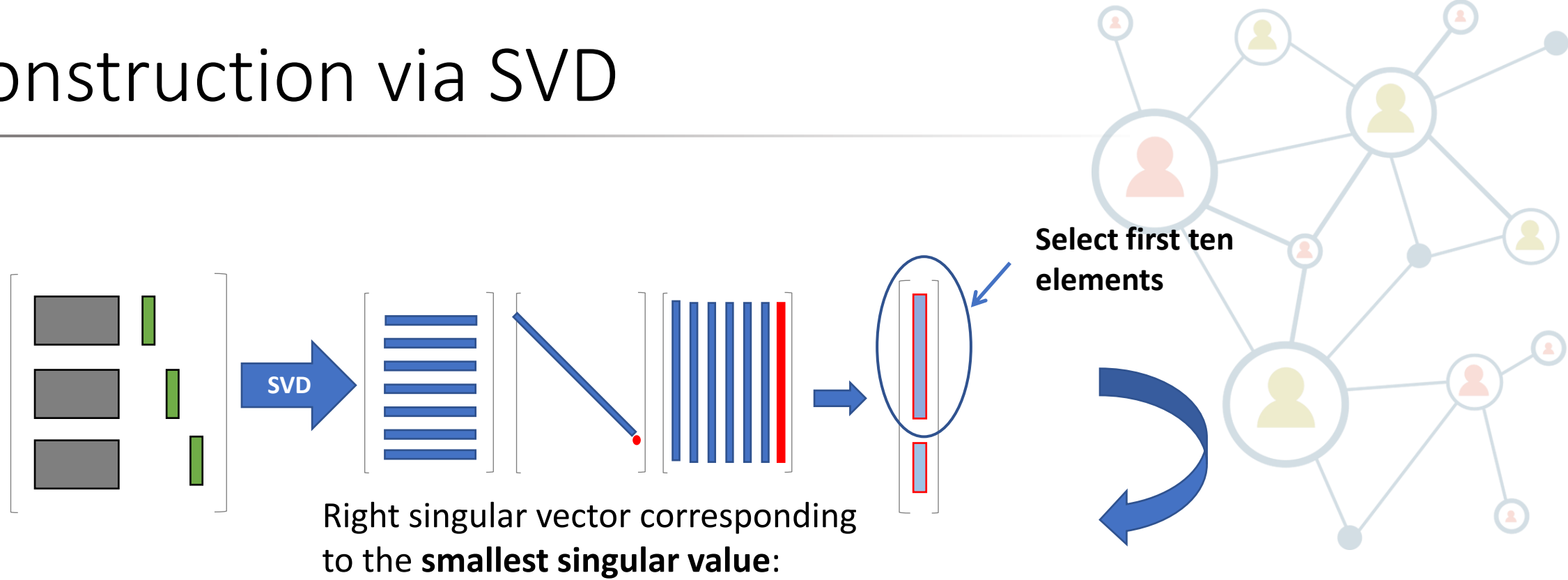
Vectorize dual conics and quadrics and stack over f



$$\mathbf{M}_i \mathbf{w}_i = \mathbf{0}_{6F}$$



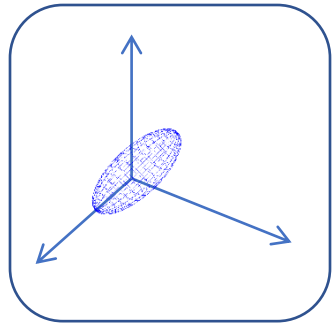
Reconstruction via SVD



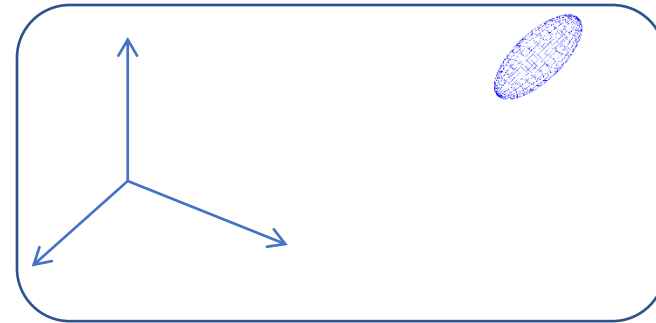
SVD approach: drawbacks

Problem solved in **dual space** may give **ill-conditioned** matrices: reverting to primal space by **matrix inversion** may result in (close to) **degenerate** ellipsoid or even other quadrics (i.e. hyperboloids).

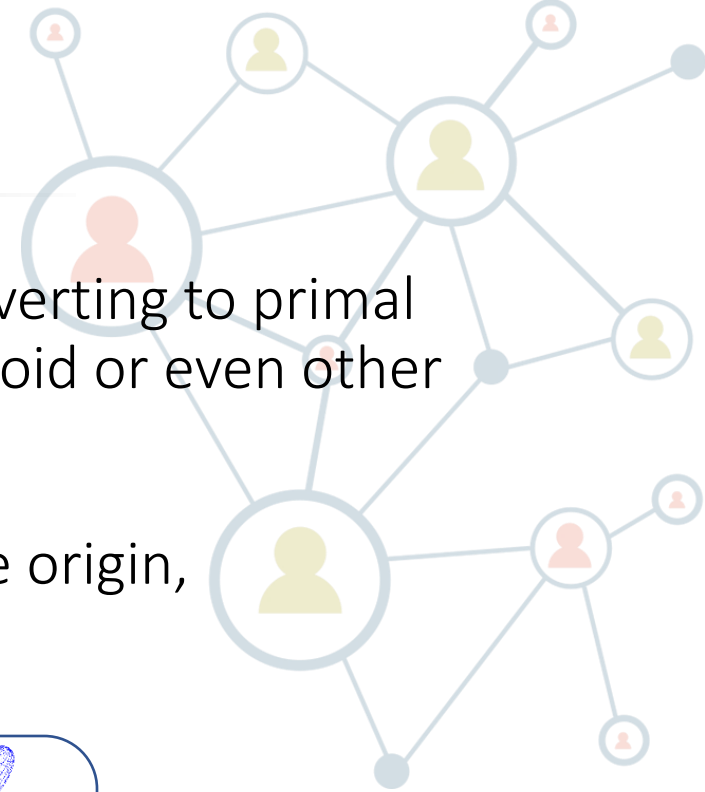
Ill-conditioning is increased for objects **far** from the 3D coordinate origin, where **translation** terms are **dominant** in the quadric.



$$\underline{\mathbf{Rcond}} = \mathbf{0.216} \quad \mathbf{Q}_i^* = \begin{bmatrix} -0.43 & -0.18 & -0.26 & 0.0 \\ -0.18 & -0.43 & -0.26 & 0.0 \\ -0.26 & -0.26 & -0.62 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.25 \end{bmatrix}$$



$$\underline{\mathbf{Rcond}} = \mathbf{0.000037} \quad \mathbf{Q}_i^* = \begin{bmatrix} 24.5 & 24.8 & 24.7 & 2.5 \\ 24.8 & 24.5 & 24.7 & 2.5 \\ 24.7 & 24.7 & 24.3 & 2.5 \\ 2.5 & 2.5 & 2.5 & 0.25 \end{bmatrix}$$



Regularization



Strategy: enforce **prior** on (known) objects **aspect ratio**: penalize the distance of ellipsoids from a **sphere of given radius and center**.

$\mathbf{S}^*(a, t_1, t_2, t_3) = \begin{bmatrix} 1 & 0 & 0 & t_1 \\ 0 & 1 & 0 & t_2 \\ 0 & 0 & 1 & t_3 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} a & 0 & 0 & 0 \\ 0 & a & 0 & 0 \\ 0 & 0 & a & 0 \\ 0 & 0 & 0 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ t_1 & t_2 & t_3 & 1 \end{bmatrix}$

Arbitrary center Axes equal Feasible radius

$a > 0$

Sphere in dual space

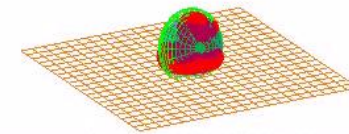
Regularization: cost function

$$\|\mathbf{M}\mathbf{w}_i\|_2^2 + \lambda \|\mathbf{v}_i^* - \text{vech}(\mathbf{S}^*(a, t_1, t_2, t_3))\|_2^2$$

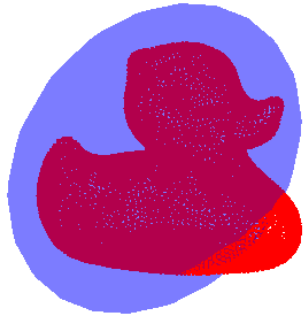
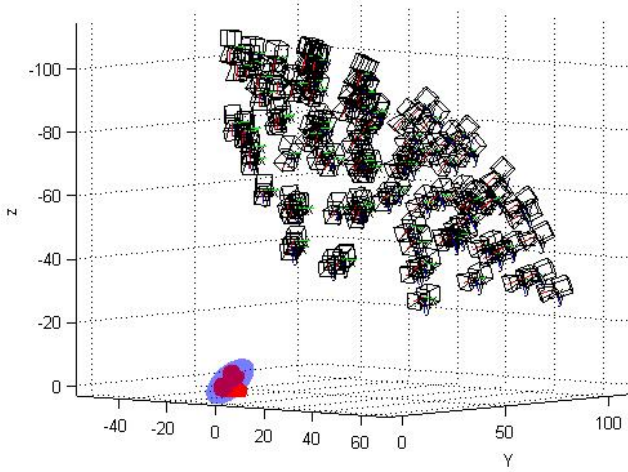
s.t. a > 0

Problem solvable with nonlinear Least Squares with boundary constraints

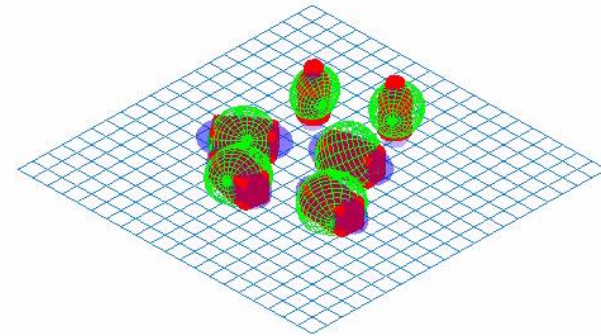
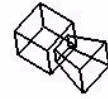
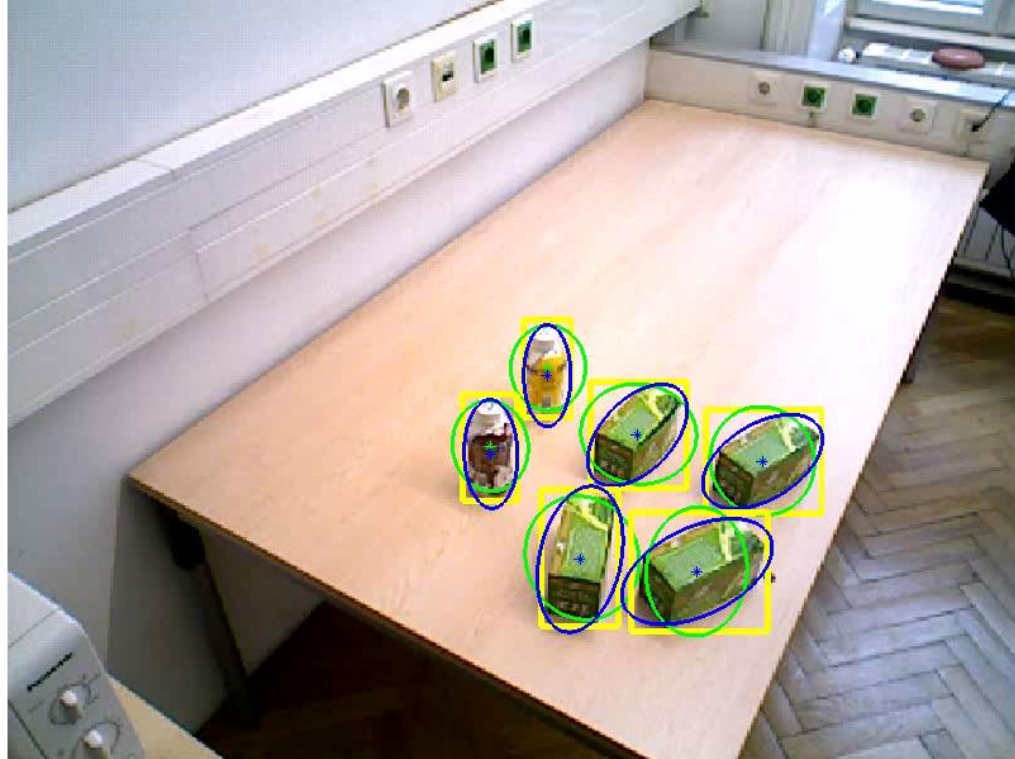
ACCV sequence (Hinterstoißer dataset)



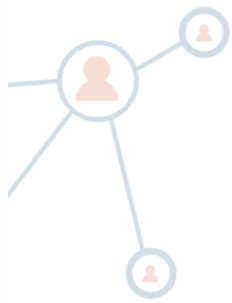
ACCV sequence (Hinterstoißer dataset)



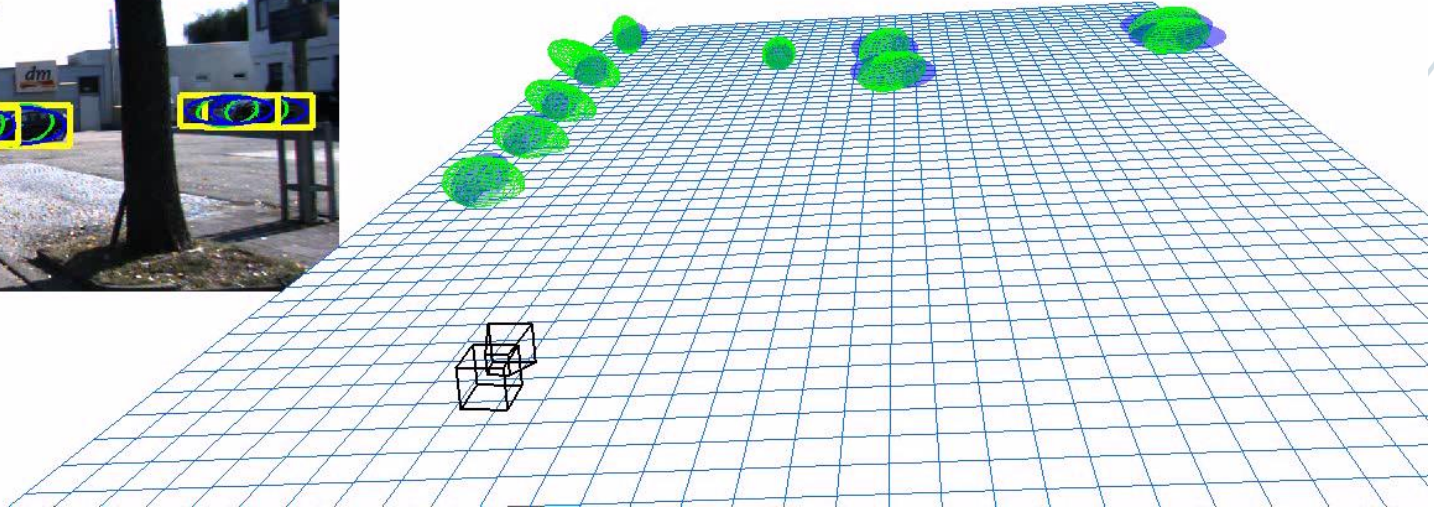
TUW dataset sequence



KITTI sequence



frame 90

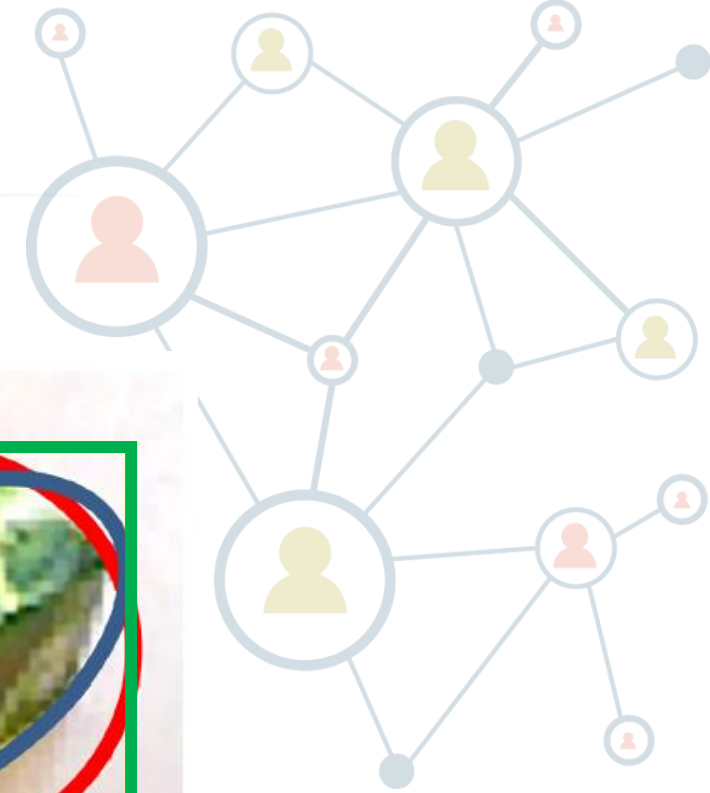


Problem 1: Ellipse from BB mismatch

- Bounding boxes (green) from generic object detectors are not precisely matched with their ground truth -> **ellipses from BB** (red) are not exactly matched as well with their **ground truth** counterpart (blue)
- Also with perfect BBs, ellipse axes are aligned to the image axes -> **high rotation and size error** in respect to GT ellipse if the object is rotated with respect to the image axes
- Further works extend ellipses fitting to objects or use instance segmentation.



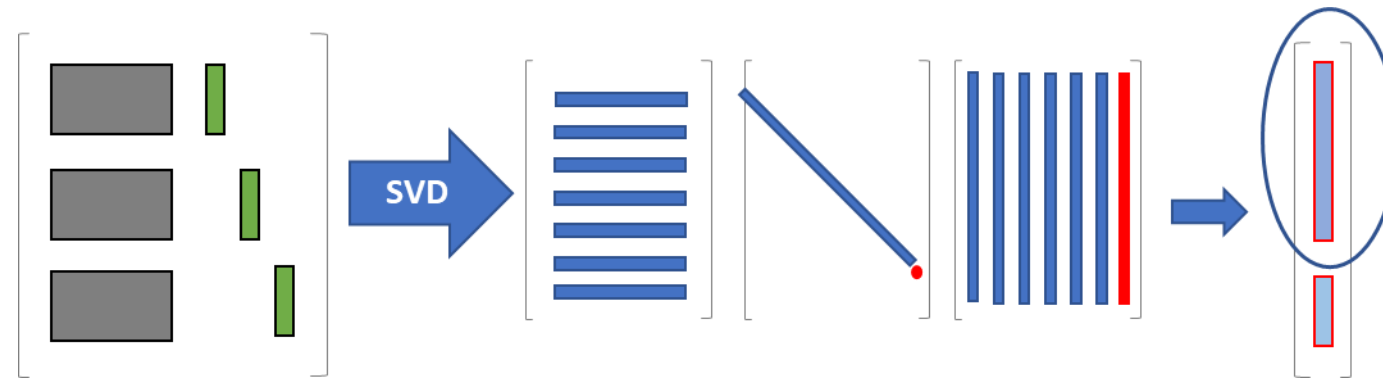
Example of mismatch between ground truth ellipse (blue) and ellipse from bounding box (red).



Dong W, Roy P, Peng C, Isler V. Ellipse R-CNN: Learning to infer elliptical object from clustering and occlusion. IEEE TIP 2021.

Problem 2: Analytical solution

The solution computes the quadric from a collection of ellipses:

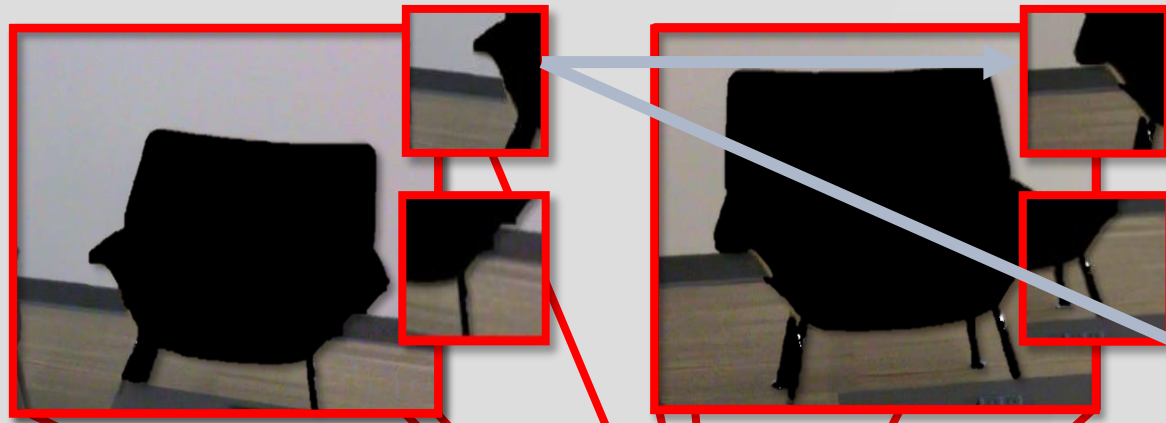


Ideally, as it happens in SfM we achieve better results and we avoid degenerate solution by minimising a reprojection error.

$$\text{However: } O_{2D} = \frac{1}{J} \sum_{f=1}^F \sum_{i=1}^N I(i, f) \frac{C_{if} \cap \tilde{C}_{if}}{C_{if} \cup \tilde{C}_{if}}$$

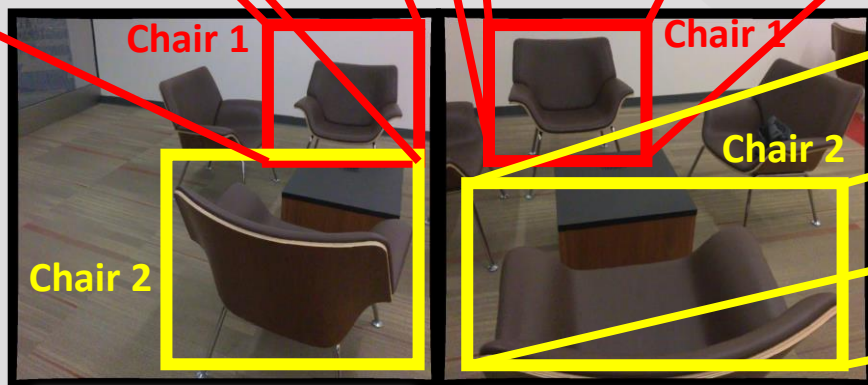


Problem 3: the same object everywhere!



Objective: identifying multiple instances of the same object in rigid scenes

Background helps in identifying the same instance in different views



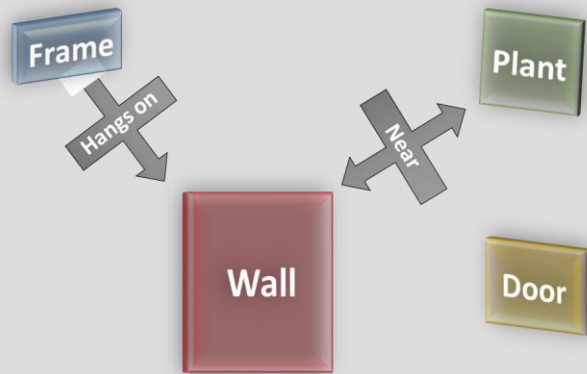
Frame t

Frame $t + n$

3D scene understanding

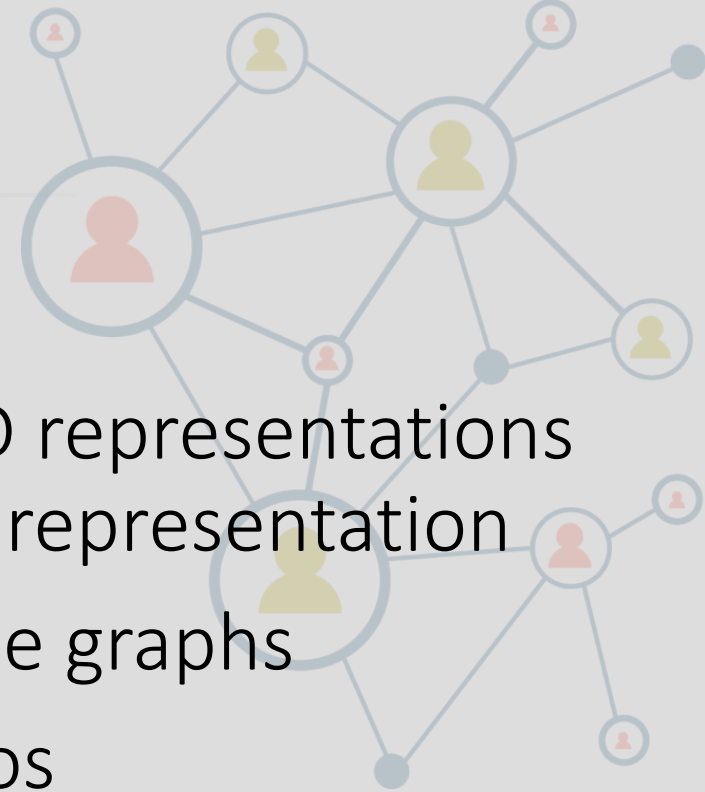


INPUT: 2D object detections in multiple views



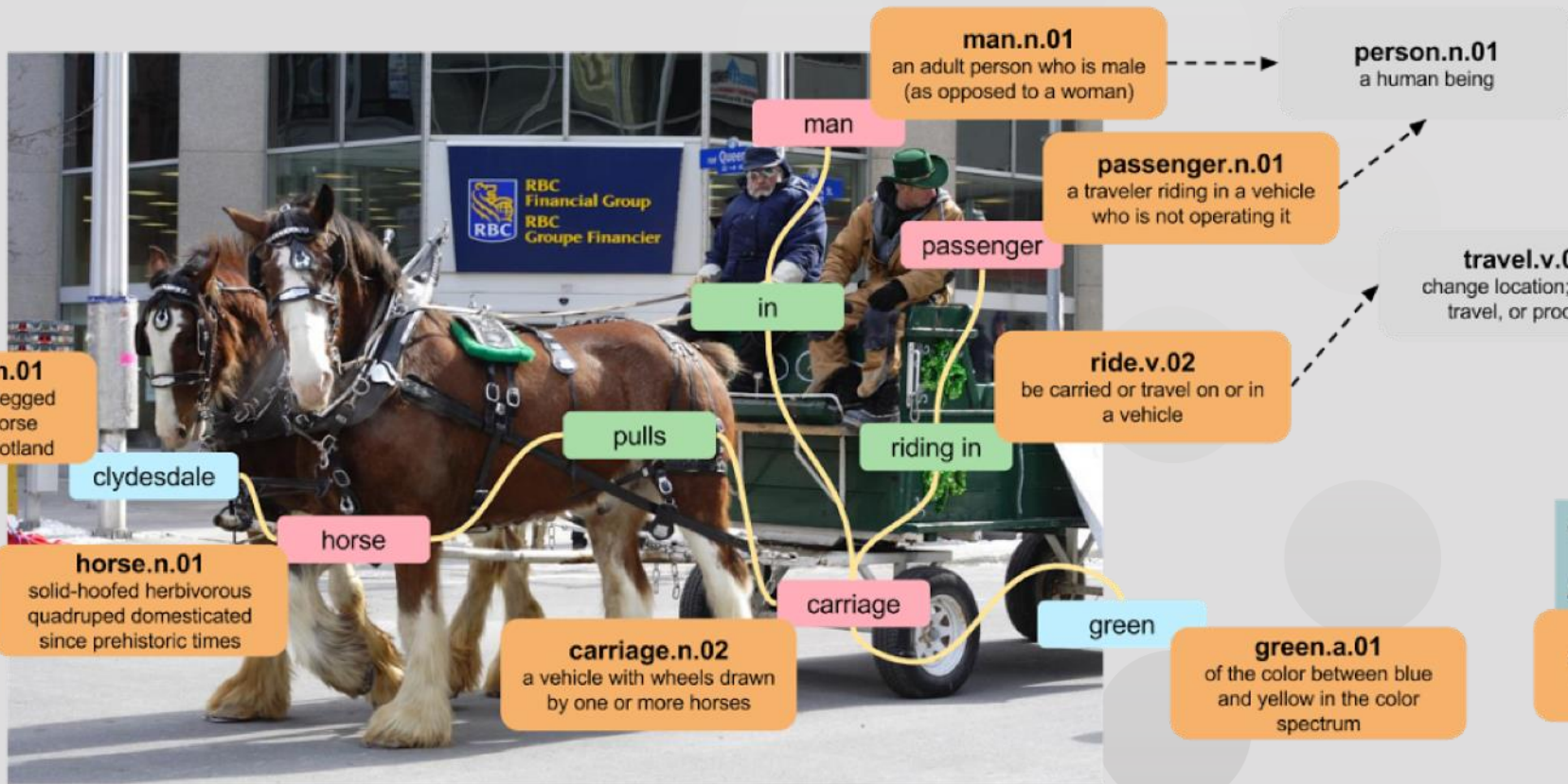
REPRESENTATION
(RIGID AND DYNAMIC)

- How to compute 3D representations with higher level of representation
- From 2D to 3D scene graphs
- Application scenarios



Back to 2D...

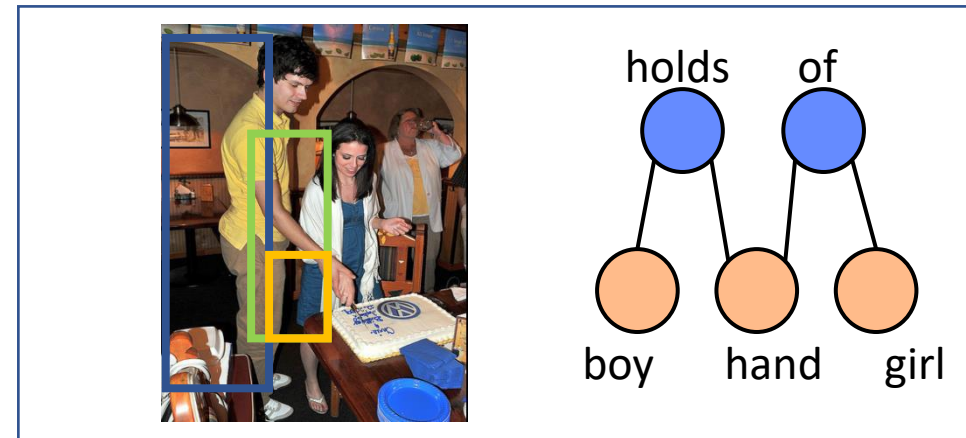
Images can be described as graph encoding objects and their relationships in the scene



Legend: object attribute relationship mapped synset derived synset QA pair extracted NP --> is derived hyponym of

Previous work

- Recent architectures generate scene graphs from single images
- Explicit modelling of object and relations with Graph Neural Network (GNN) [1,2].
 - Global optimisation with message passing on a bi-partite graph.
 - Dealing with chain of relations

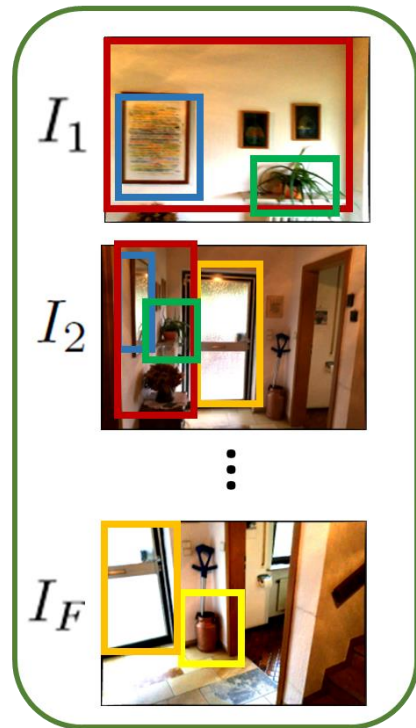


[1] Danfei et al, Scene Graph Generation by Iterative Message Passing, CVPR 2017

[2] Yikang et al, Scene Graph Generation from Objects, Phrases and Region Captions, ICCV 2017

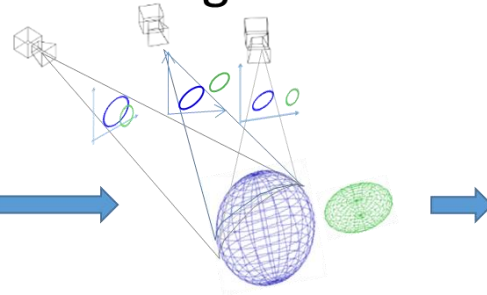
From 2D to 3D scene graphs

Input: object detections in multiple views

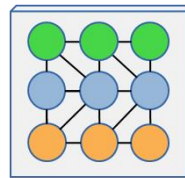


Object Bounding Boxes

3D object position from bounding boxes

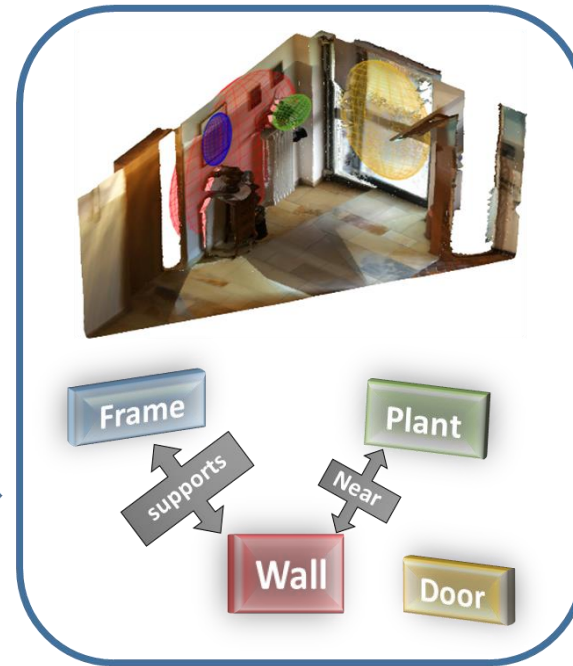


Visual features



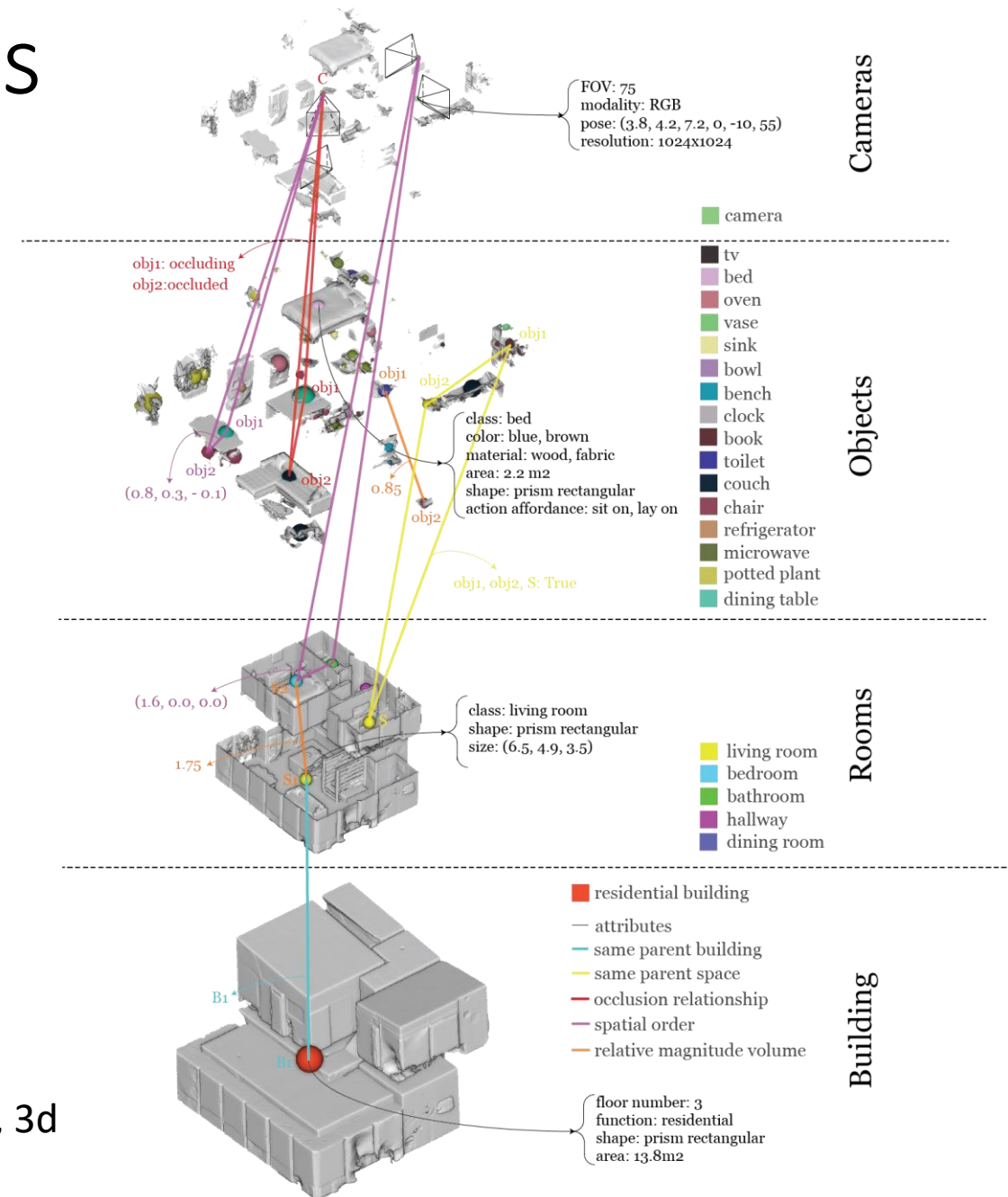
Scene graph Prediction (RNN)

Output: 3D Scene Graph



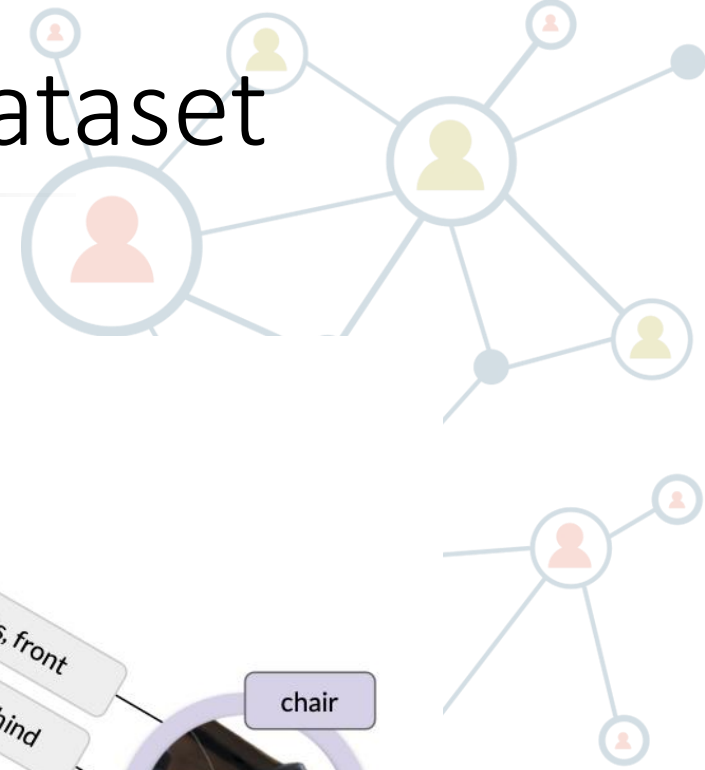
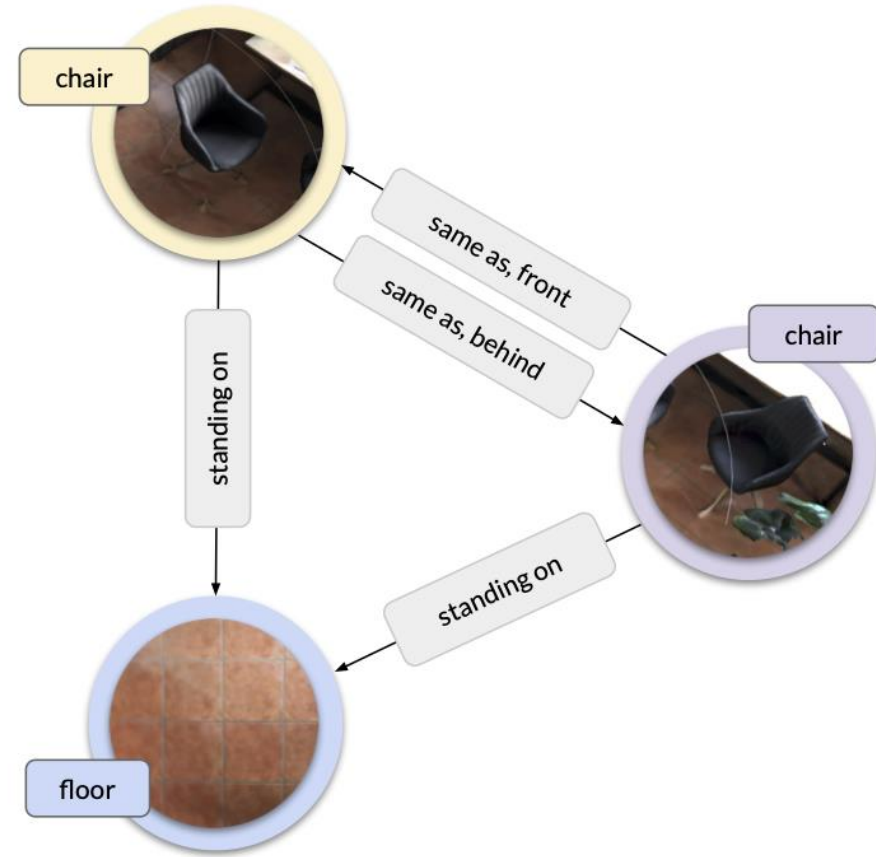
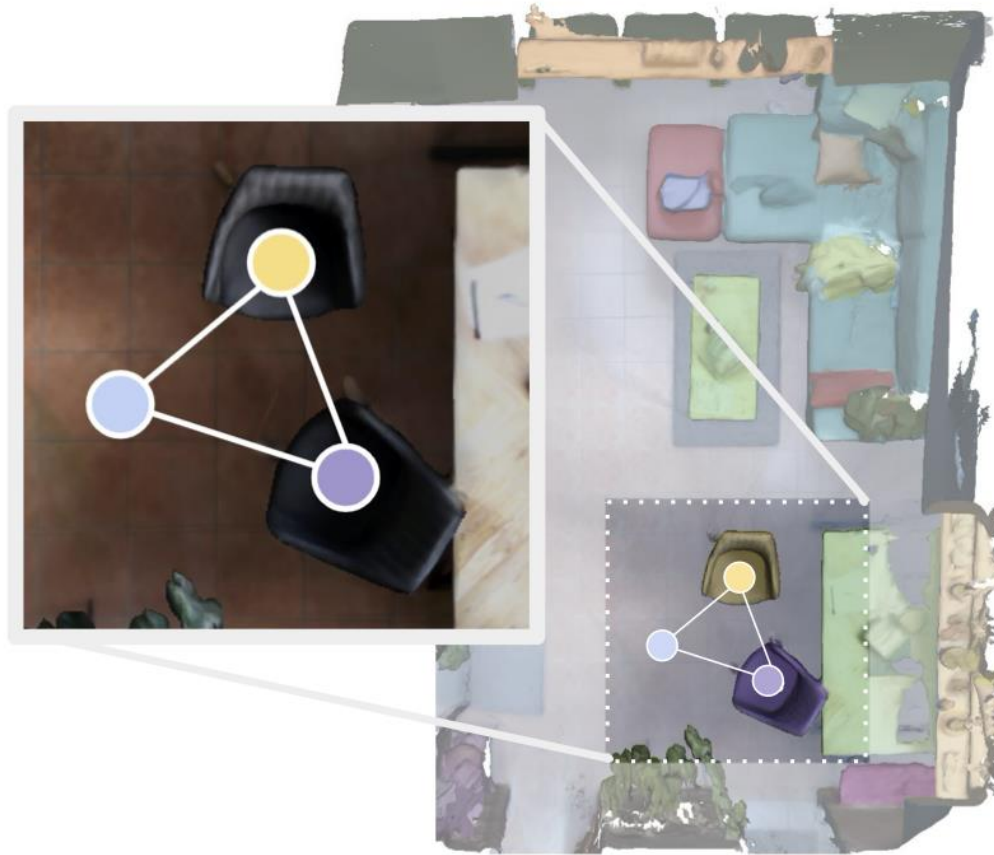
From 2D to 3D scene graphs

- Scene graphs include semantics on objects (e.g., class, material, and other attributes), rooms (e.g., scene category, volume, etc.) and cameras (e.g., location, etc.)
- The representation can have different levels of hierarchy, from a building to a city...

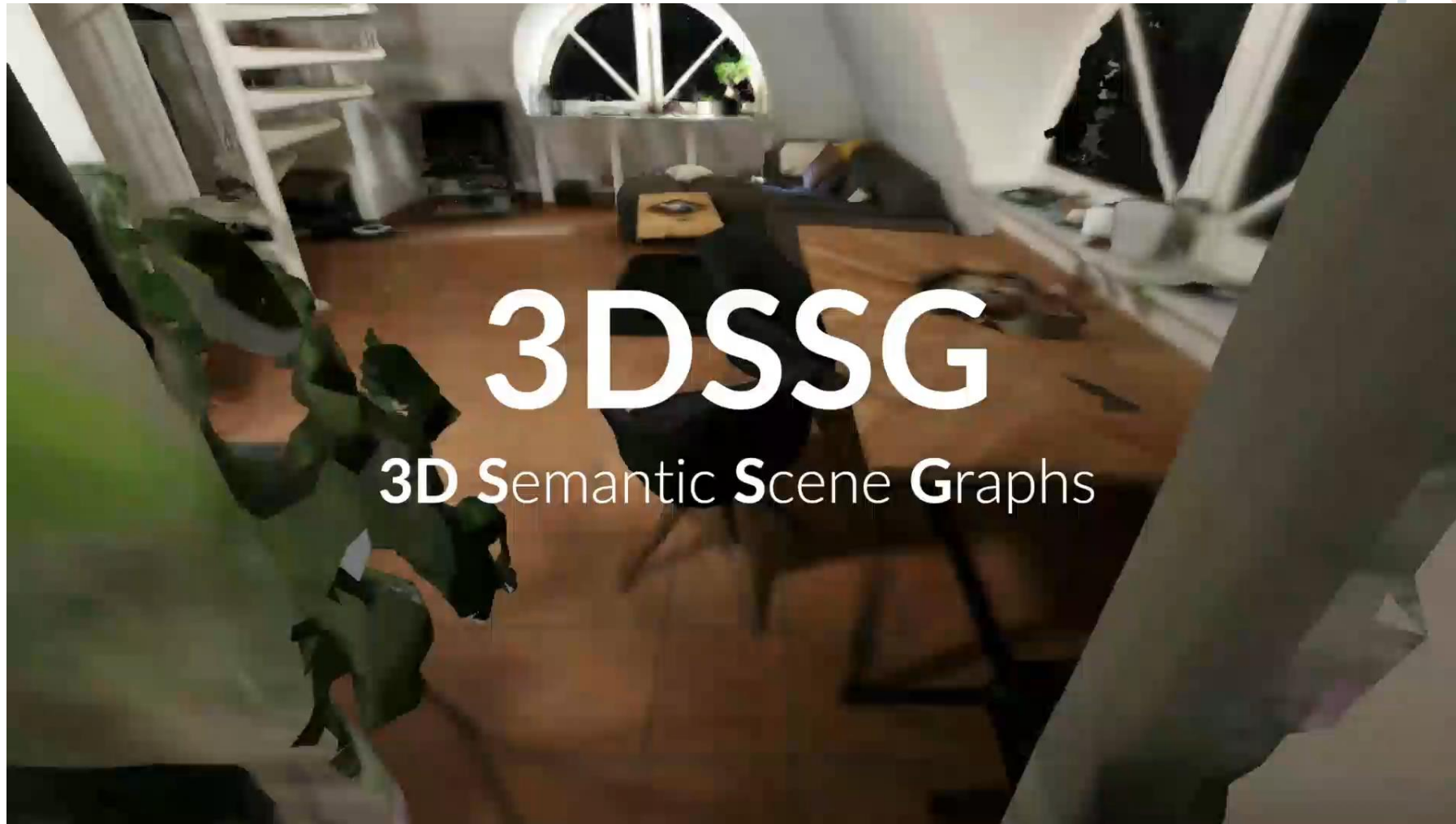


Armeni, Iro, et al. "3d scene graph: A structure for unified semantics, 3d space, and camera." ICCV 2019.

Learning graphs from RGBD indoor dataset



Learning graphs from RGBD indoor dataset



Wald, Johanna, et al. "Learning 3d semantic scene graphs from 3d indoor reconstructions." *CVPR* 2020.

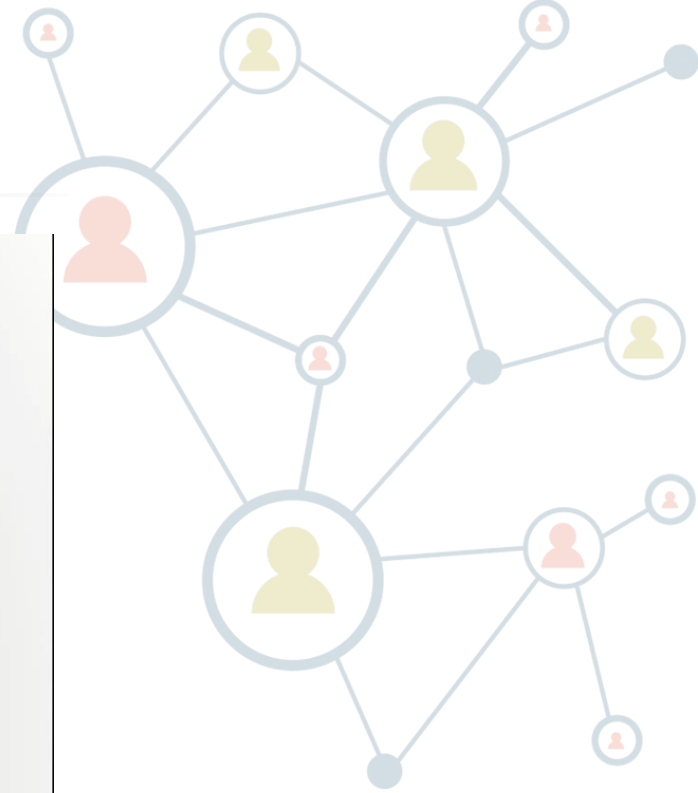
Learning graphs in real-time

Hydra: A Real-time Spatial Perception Engine for 3D Scene Graph Construction and Optimization

Nathan Hughes, Yun Chang, Luca Carlone



Hughes N, Chang Y, Carlone L. Hydra: A real-time spatial perception system for 3D scene graph construction and optimization. RSS 2022.

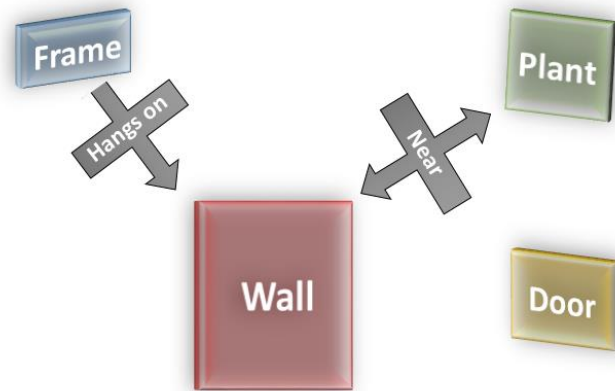


3D scene graph from multi-view

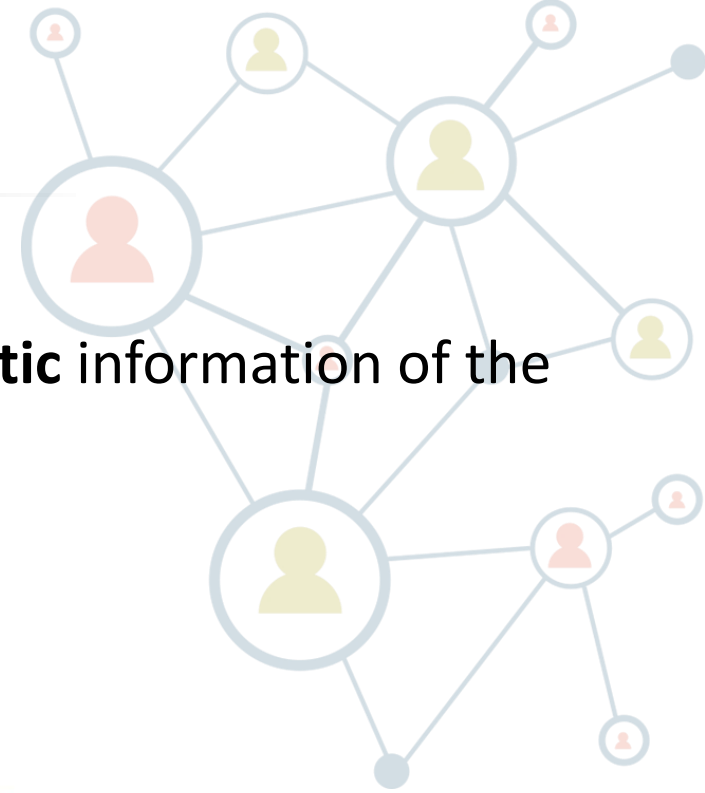
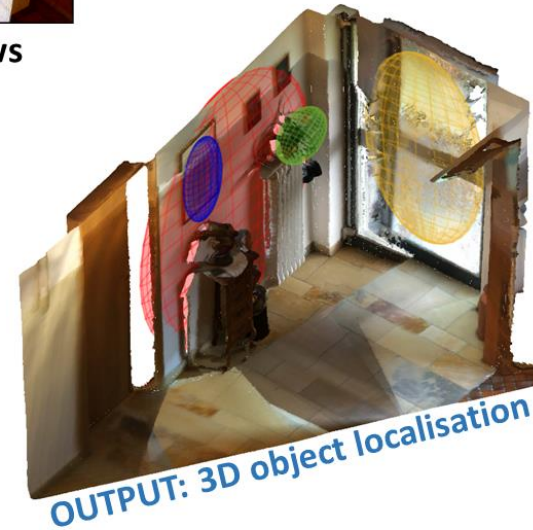
Goal: Build **3D Scene graph** to obtain a **rich** and **grounded semantic** information of the scene using only images and object detections.



INPUT: 2D object detections in multiple views



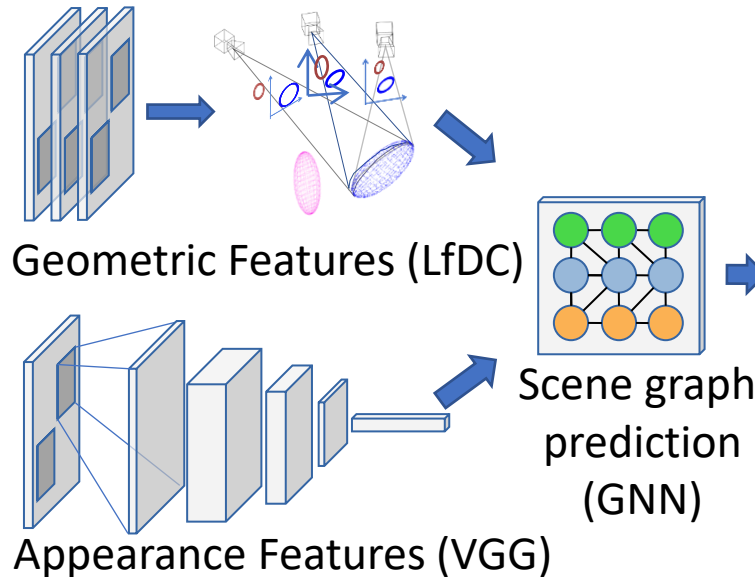
OUTPUT: 3D Scene Graph



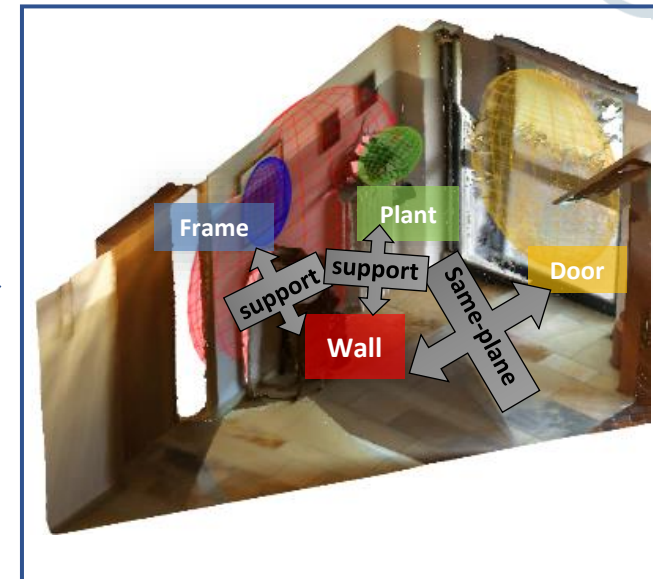
How to build a 3D scene graph?

- Extraction of the observations: Visual and geometric features
 - Ellipsoids computed from multiple views.
 - Visual features extracted from each image.
- Scene graph prediction with a tri-partite Graph Neural Network (GNN)

Input: Image object detections

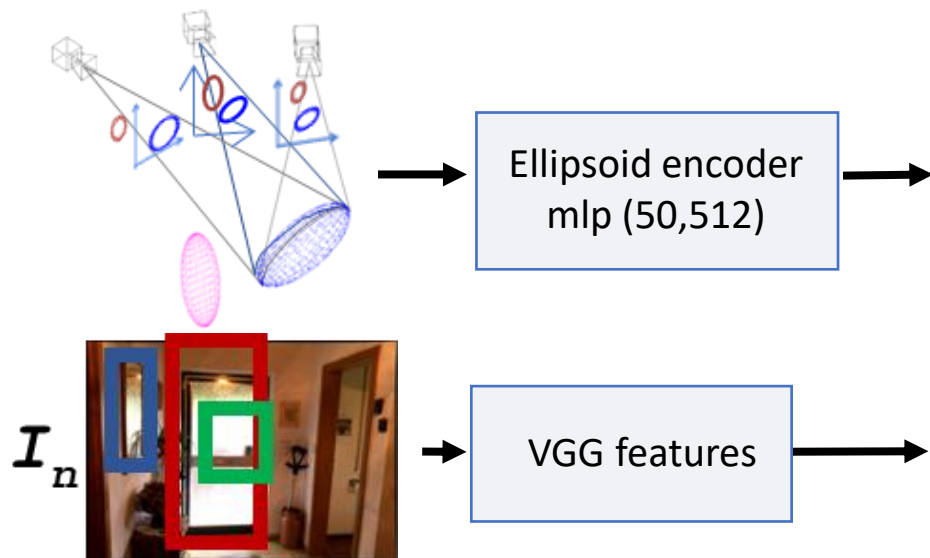


Output: 3D Scene Graph



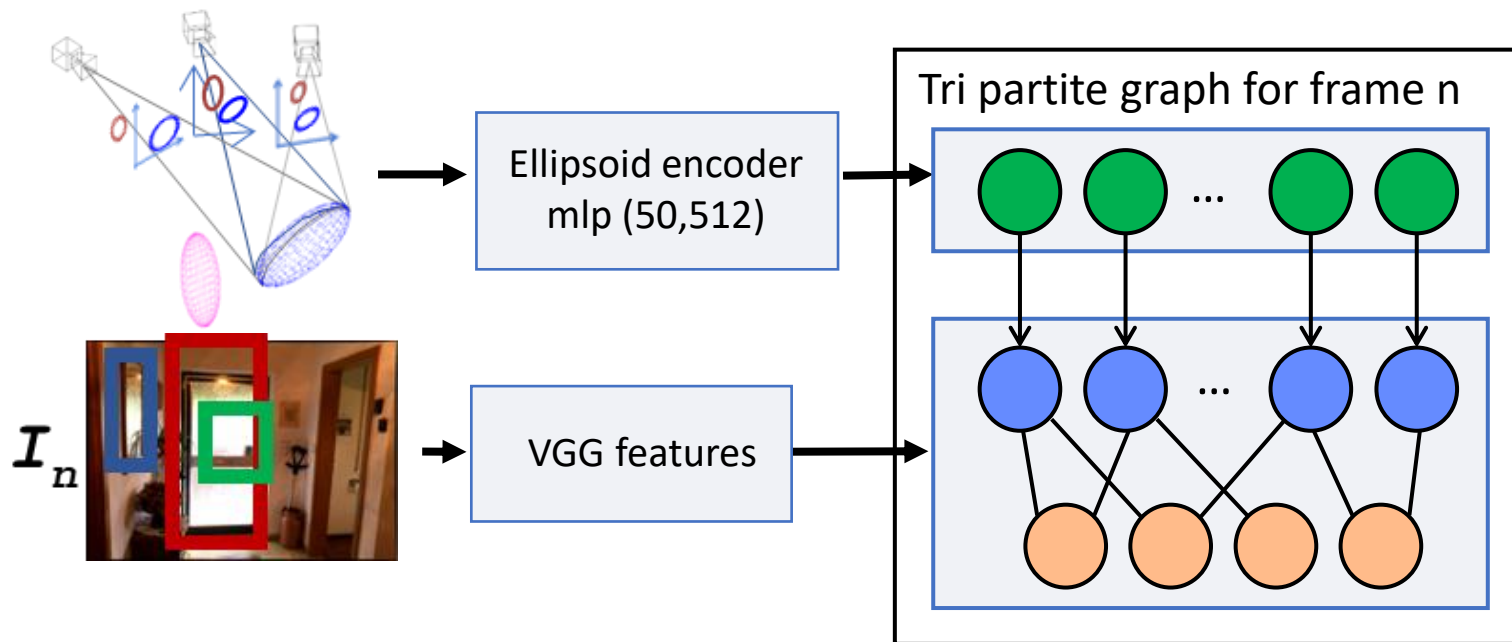
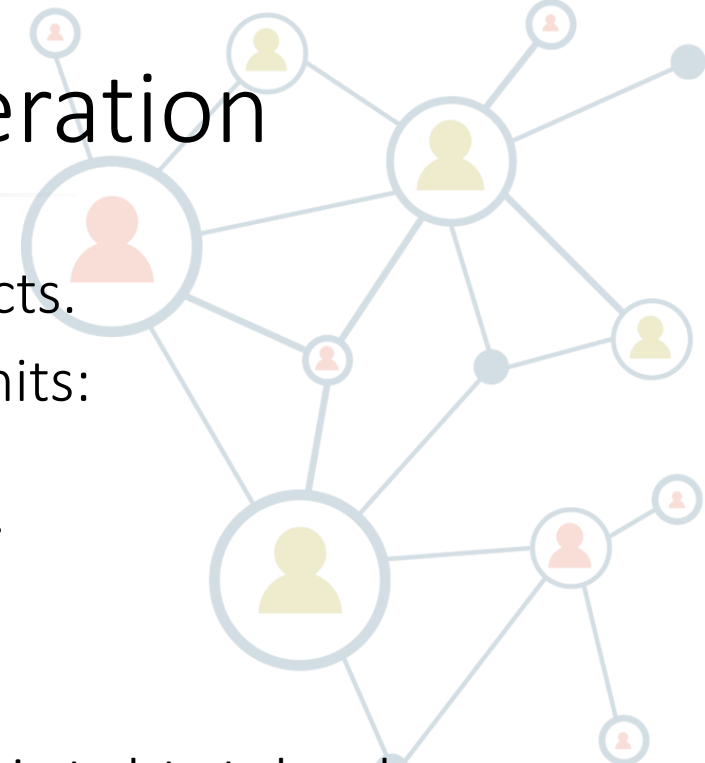
Use of the tri-graph for scene graph generation

- Visual appearance representation : VGG features are extracted from each bounding box.
- Geometric representation: MLP encodes each pair of ellipsoids.
- These features are then used as initialisation of the node states.



Use of the tri-graph for scene graph generation

- Tri-graph containing every possible relation between the N objects.
- Relation (blue) and object nodes (orange) are RNNs with GRU units:
 - Update of the GRU hidden state thanks to incoming messages;
 - Geometric nodes (green) are observations, i.e. they only send messages.



For N objects detected, we have:

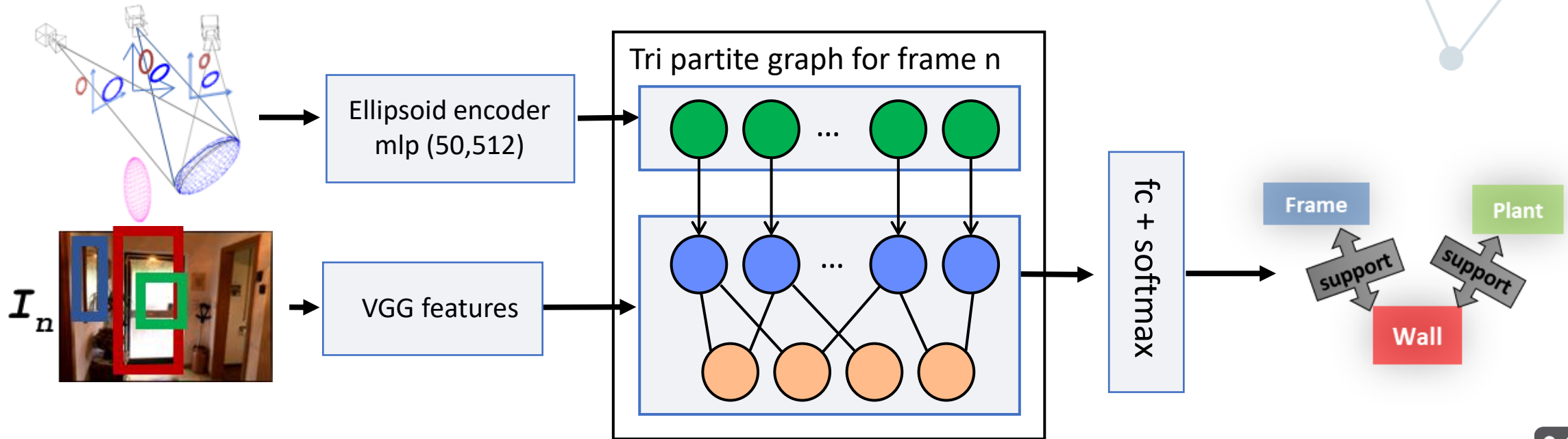
● $N^2 - N$ GRU geometric nodes

● $N^2 - N$ GRU relation nodes

● N GRU object nodes

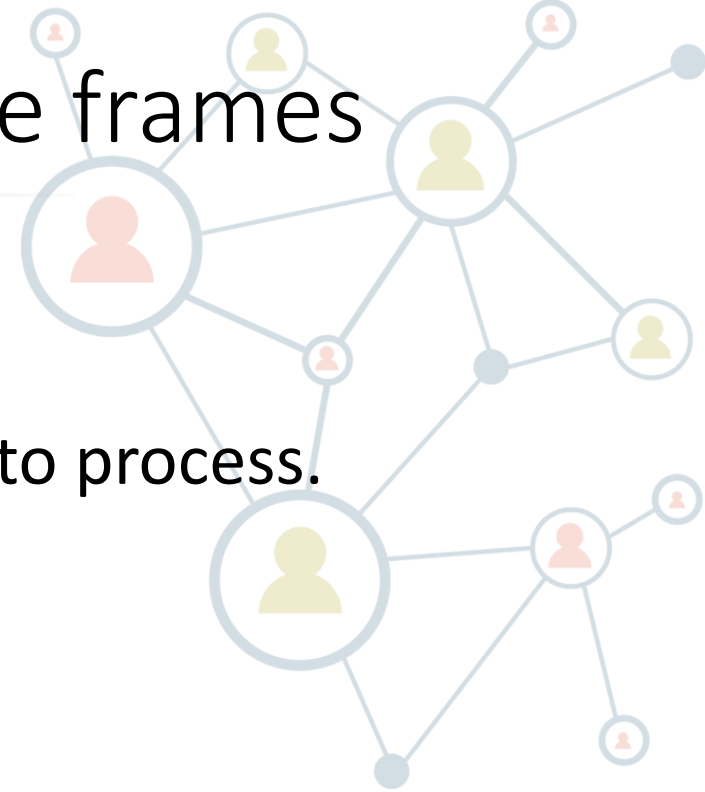
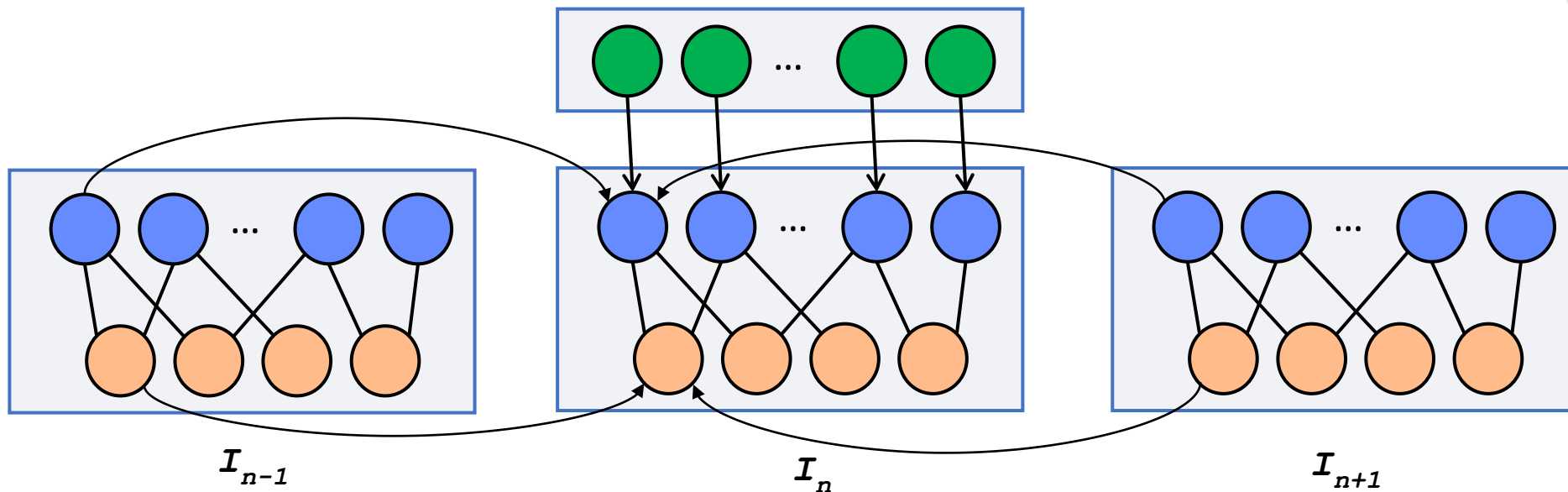
Use of the tri-graph for scene graph generation

- Lastly, a softmax layer gives the graph prediction from the hidden states:
Probability distribution for each relation and object nodes over its set of labels.
- We now have a scene graph given a single image and geometry:
But given a video, how to exploit the visual appearance from the other images?



Fusion of the visual features with multiple frames

- Message passing among multiple frames
 - Fusion of the appearances of the object through the video.
- One tri-graph is instantiated in parallel for each image to process.

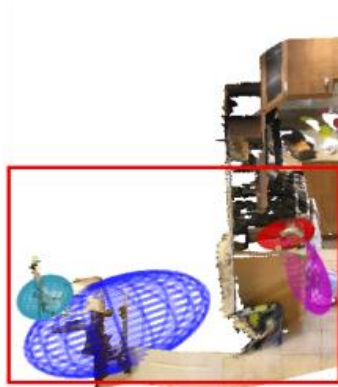


Qualitative results on ScanNet

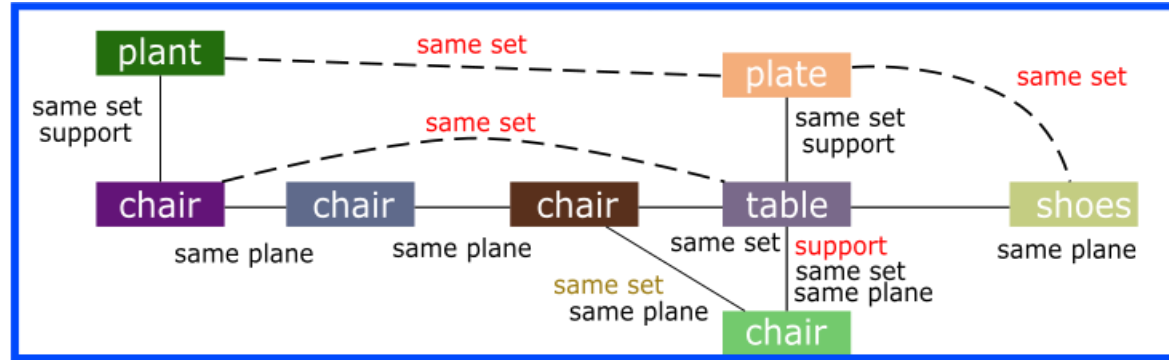
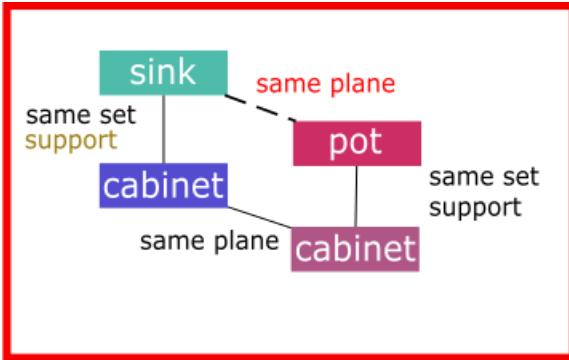
Bounding Boxes



Reconstructed Ellipsoids



Scene graph



Correct
False alarm
Missed

Some final considerations

Considerations:

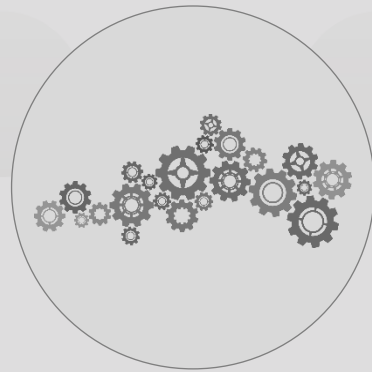
- Structure from Motion and 3D reconstruction methods are mature, but what use can have a trillion of 3D points?
- Scene graphs are an accessible, human interpretable, way to visualise the information.



Ongoing problems:

- Annotating 3D data is already hard, scaling to annotate pairwise relations between objects is too time consuming -> **self-supervision!**
- We might build efficient rigid scene graphs, what happens when the scene is dynamic? (check: Rosinol, A., Gupta, A., Abate, M., Shi, J., & Carlone, L. (2020). 3D Dynamic Scene Graphs: Actionable Spatial Perception with Places, Objects, and Humans. RSS 2020.)

Applications



MEMEX



Linking digital information to 3D and 2D maps

Matteo Taiana, Matteo Toso, Stuart James, Alessio Del Bue
Istituto Italiano di Tecnologia (IIT)



2019-2022



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 870743



Phil Lynott Statue

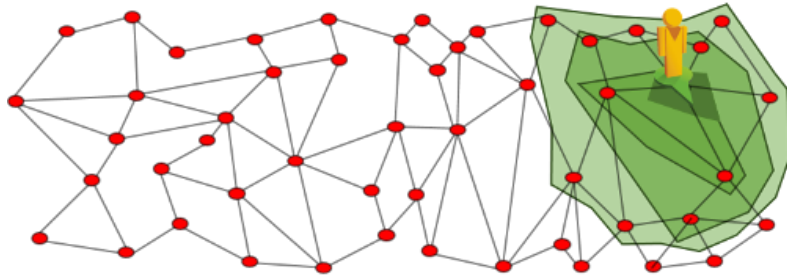
Paper ID: 1116



MEMEX The logo for MEMEX, featuring the word "MEMEX" in a blue, sans-serif font. To the right of the text is a cluster of approximately 12 small circles in red and green, arranged in a roughly circular pattern.

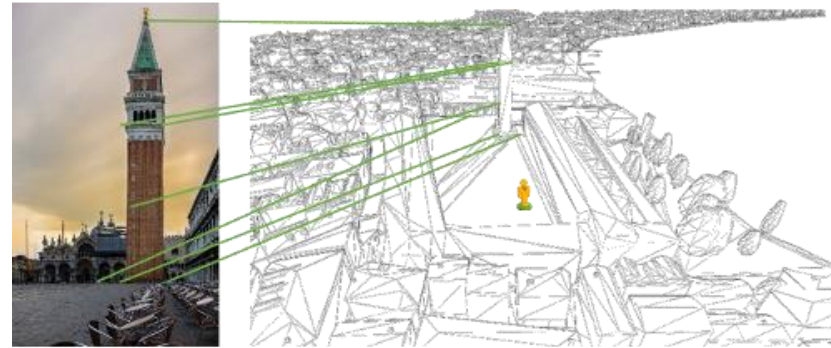
Artificial Intelligence enhances Technology

Researching and developing new technologies to facilitate the increase collaboration and inclusion of communities. MEMEX focuses on three core reusable technologies:



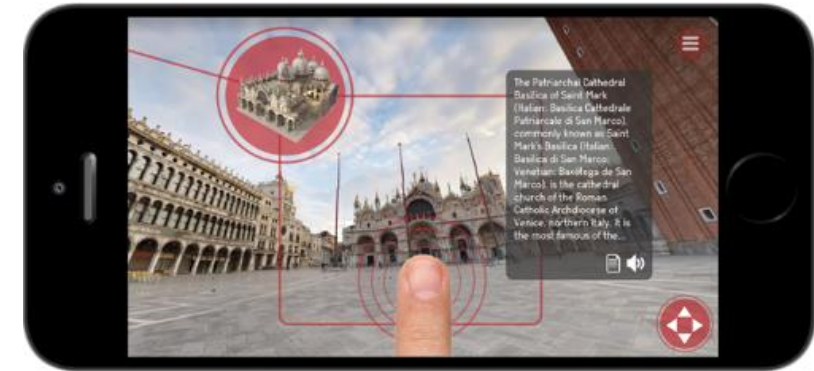
Knowledge Graph

Creating new infrastructure for geolocalised Cultural Heritage to reason on.



Localisation

Computer vision based automatic localisation of users and objects.

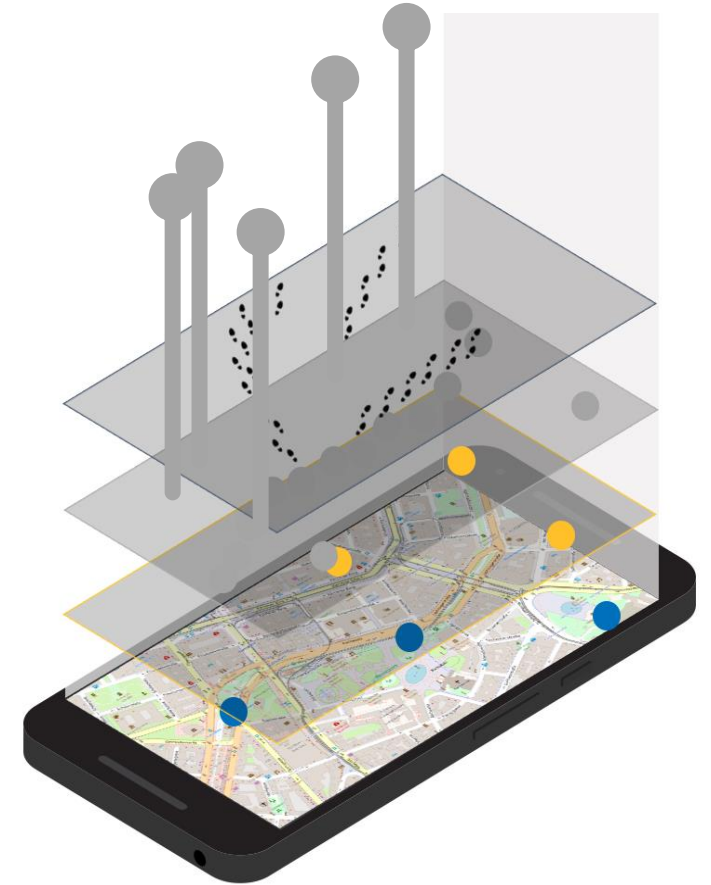


Storytelling through Augmented Reality

Assisted story creation and visualisation using advanced AR technologies.

Combining information from multiple sources

- Open maps
- Digital GLAM's content
- Street level images
- Users stories and objects



Map registration problem

Localised data sources

The image is a horizontal collage of seven panels illustrating localised data sources:

- Wikidata:** A yellow panel featuring the Wikidata logo (vertical bars) and a screenshot of a Wikidata entity page for "Douglas Adams". The screenshot shows fields for label, description, statements, and property, with various values and references listed.
- OpenStreetMap:** A blue panel showing a magnifying glass over a street map, with the "OpenStreetMap" logo below it.
- Mapillary:** A green panel showing a street map with orange lines and dots representing user-uploaded photos, with the "Mapillary" logo on the left.
- Databases:** An orange panel with the text "Databases: Local archives and Museums, +" and the "Europeana think culture" logo.
- MEMEX:** A row of seven small images: a book cover "THE STRANGER'S GUIDE IN LISBON", a building facade with graffiti, a poster for "A FEIRA DO BIRATEJO SANTA RÊM" (22 de Maio a 5 de Junho - 1955), a historical illustration of figures, a portrait of a man in a military uniform, a painting of a man on horseback, and a black and white photograph of a group of people in a room.

Localised data sources



Databases:
Local archives and Museums,

+



label — Douglas Adams (Q42) — item identifier

description — English writer and humorist
Douglas Noel Adams | Douglas Noel Adams — aliases
► In more languages

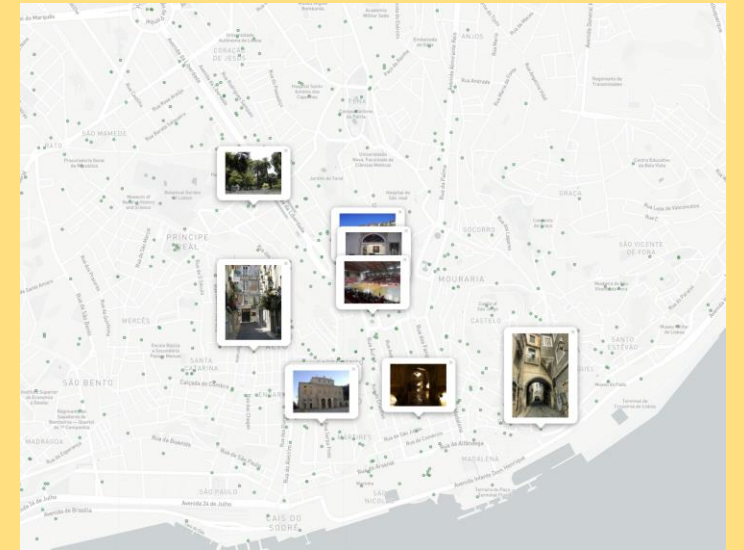
Statements

property — educated at — value

rank —

statement group —

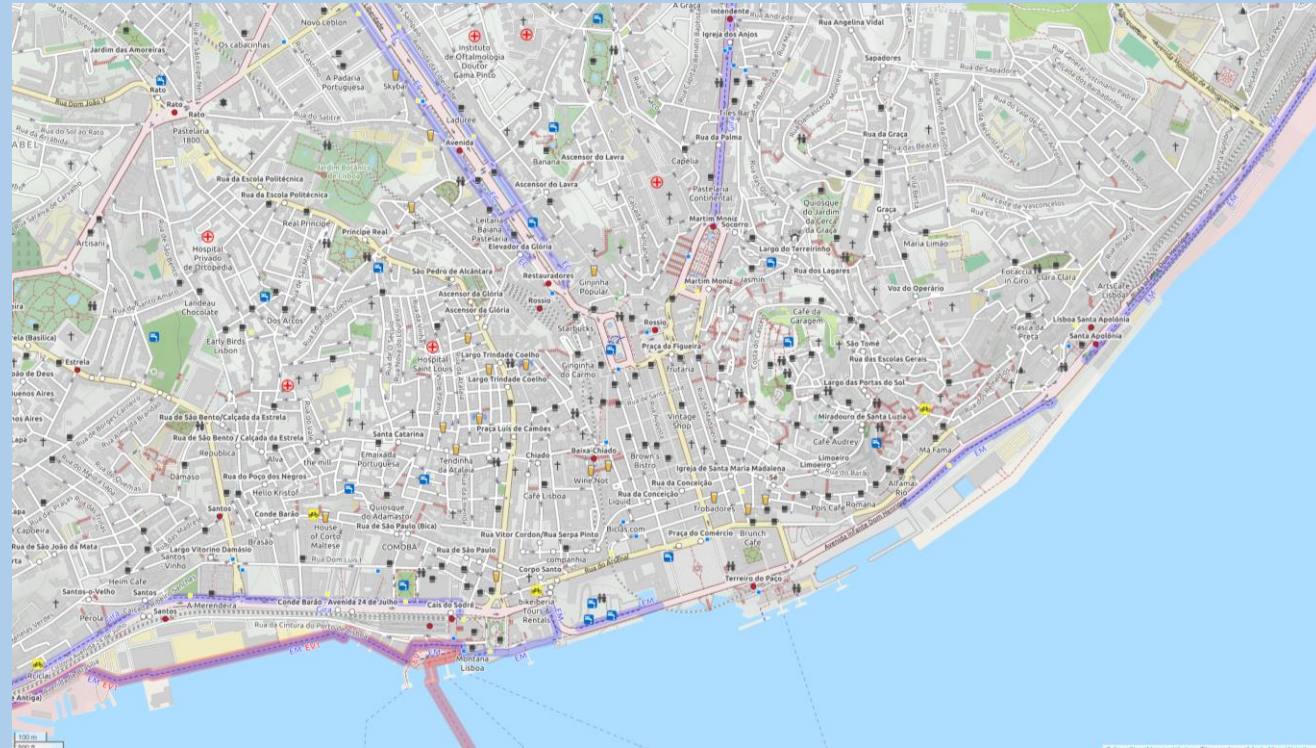
St John's College	value
end time	1974
academic major	English literature
academic degree	Bachelor of Arts
start time	1971
qualifiers	
▼ 2 references	
opened references	
stated in	Encyclopædia Britannica Online
reference URL	http://www.britannica.com/peoples/024/000023662/
original language of work	English
retrieved	7 December 2013
publisher	NNDG
site	Douglas Adams (English)
+ add reference	
Brentwood School	
end time	1970
start time	1969
collapsed reference	
► 0 references	
+ add statement	



Localised data sources



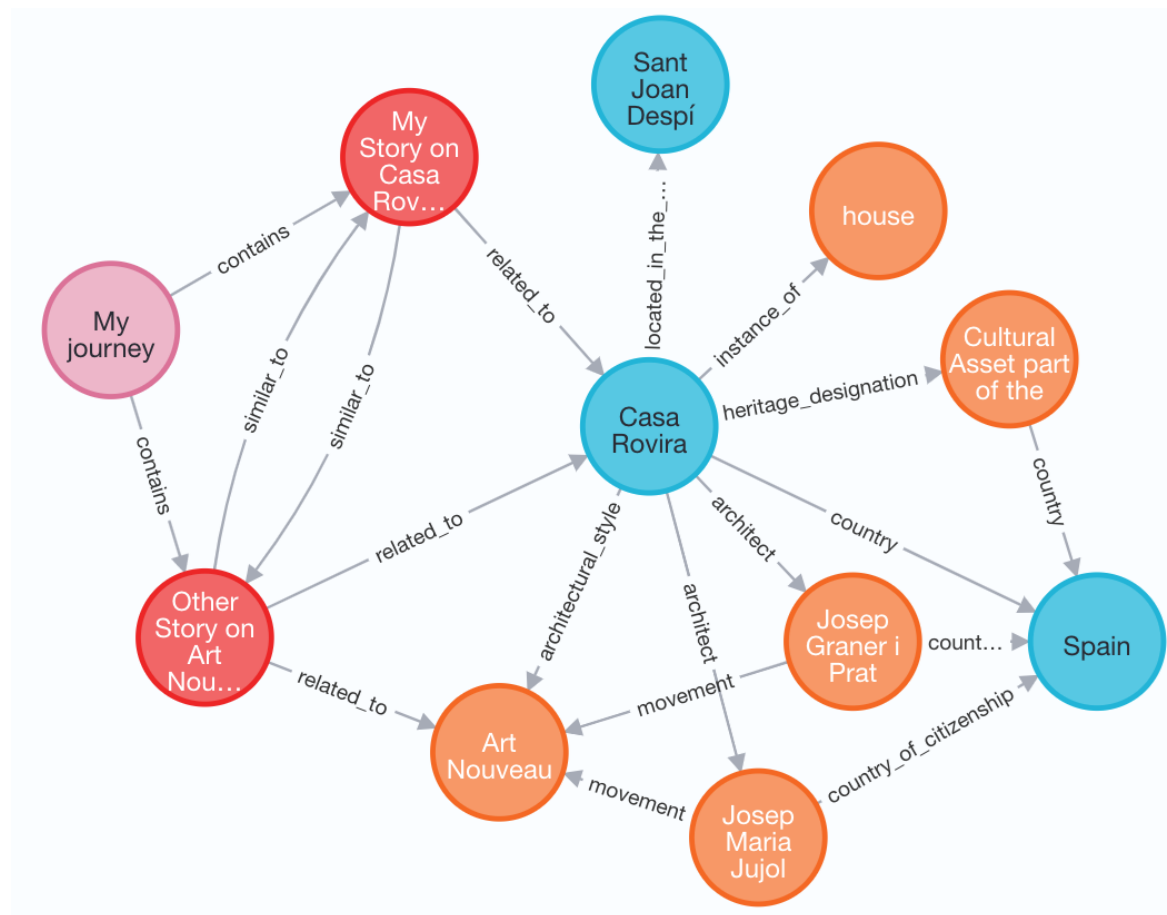
Databases:
Local archives and Museums,
+



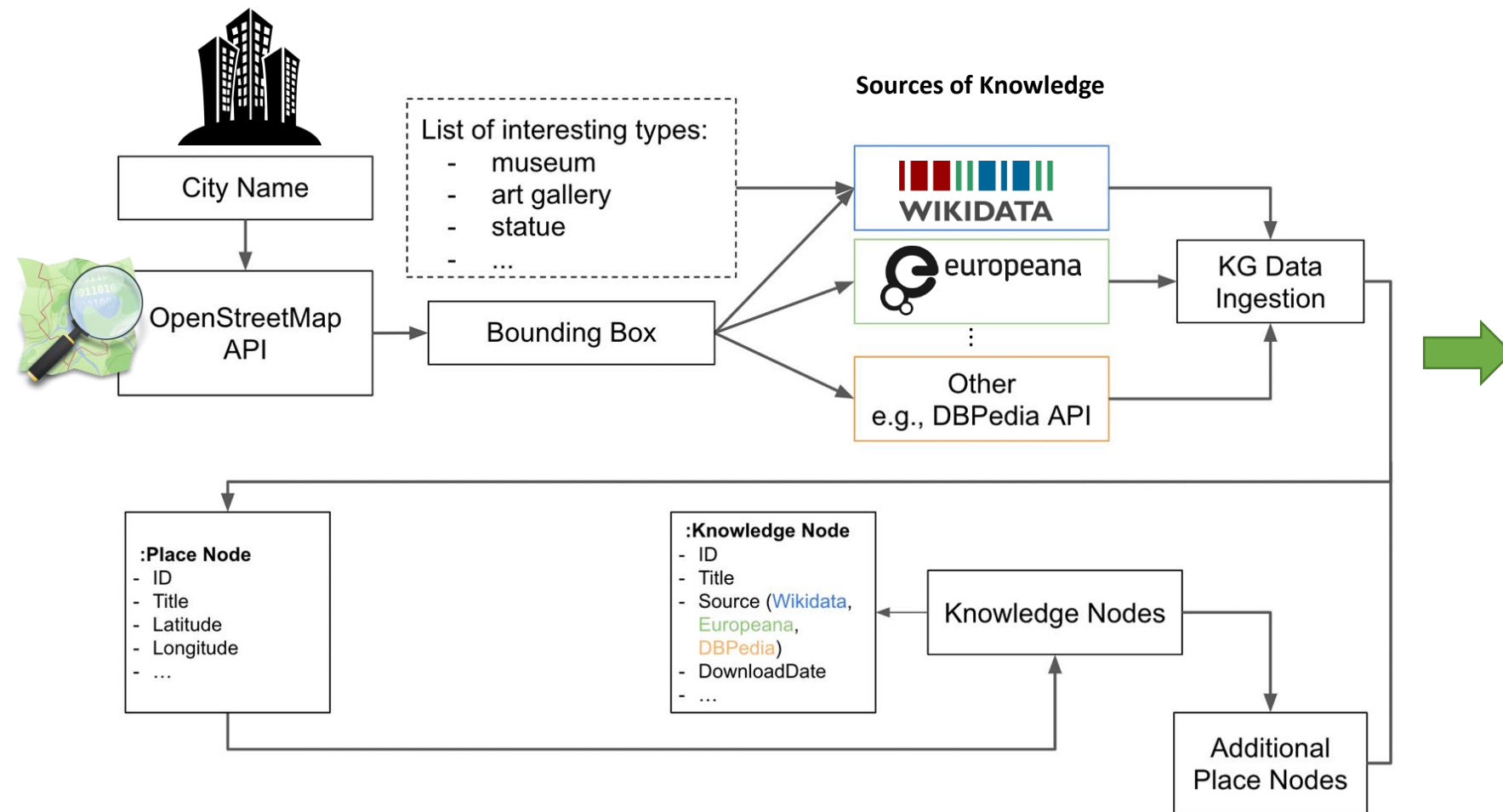
MEMEX

Knowledge Graph construction

The MEMEX-KG Structure:



Knowledge Graph construction



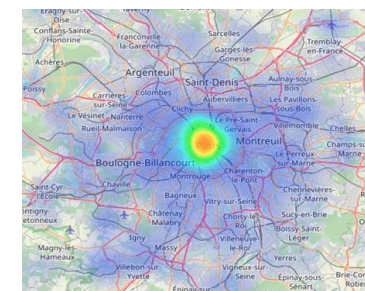
Lisbon

10,209 entities
18,860 relations
326 relation types



Barcelona

16,832 entities
43,934 relations
394 relation types



Paris

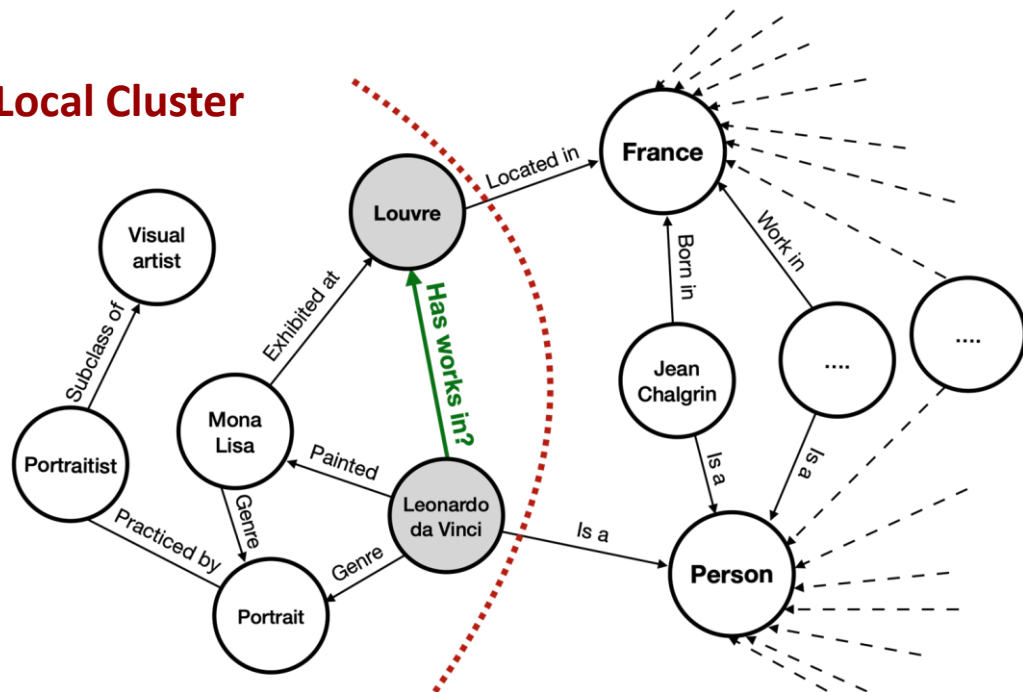
62,211 entities
192,036 relations
562 relation types

Knowledge and connecting knowledge

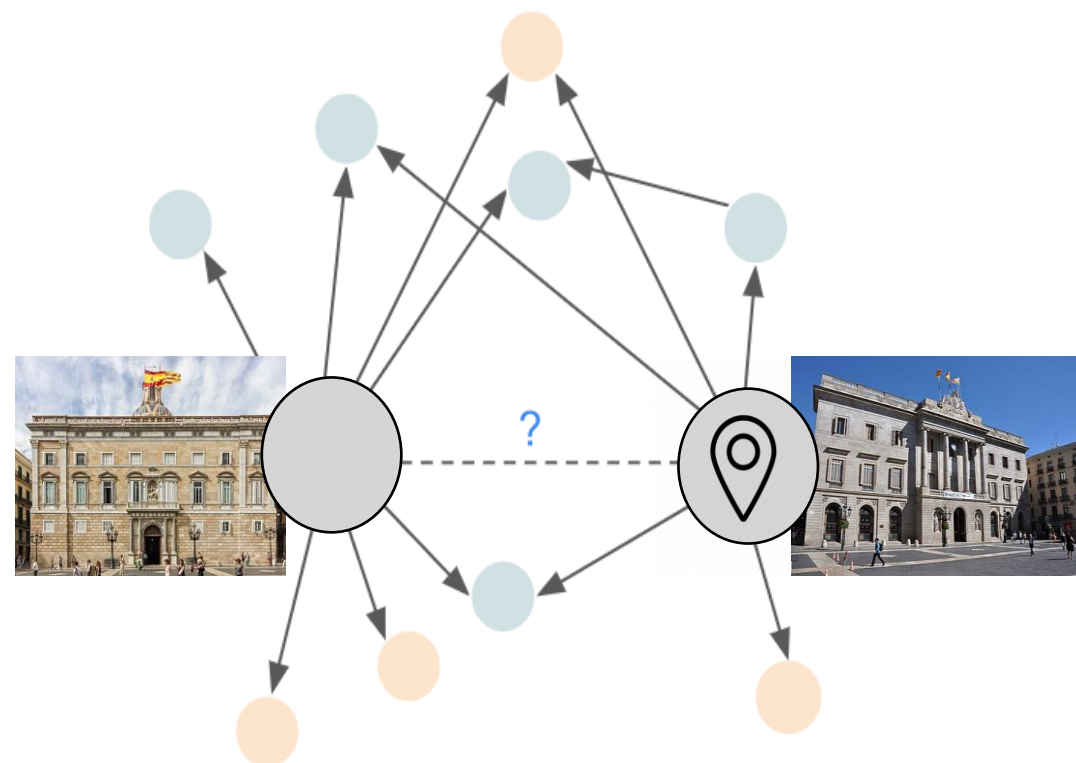
KGs suffer from **incompleteness** - especially when ingesting data from heterogeneous sources of information

Link Prediction:
 inferring missing links between nodes in the KG¹

Local Cluster



Geolocation of CH using Graph Embeddings²

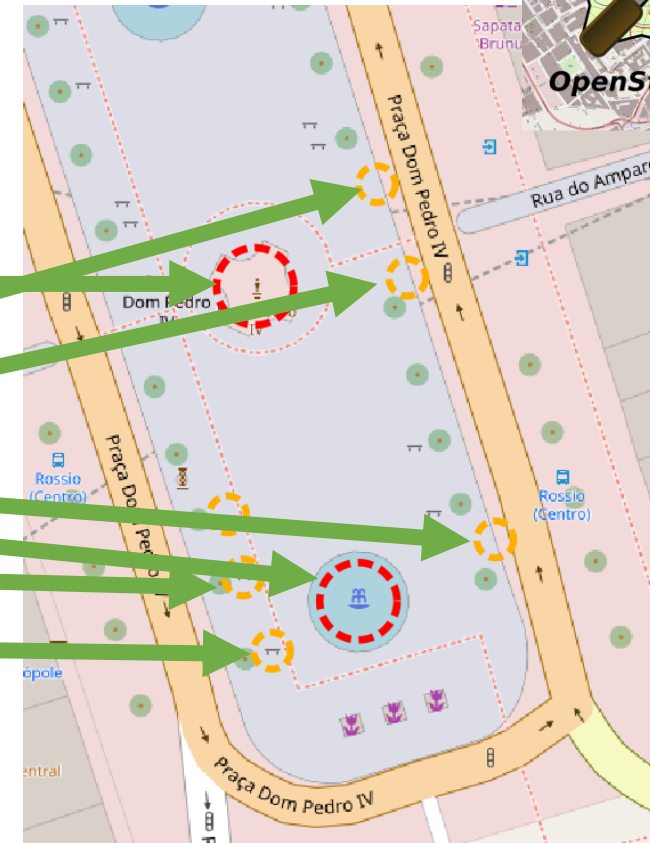
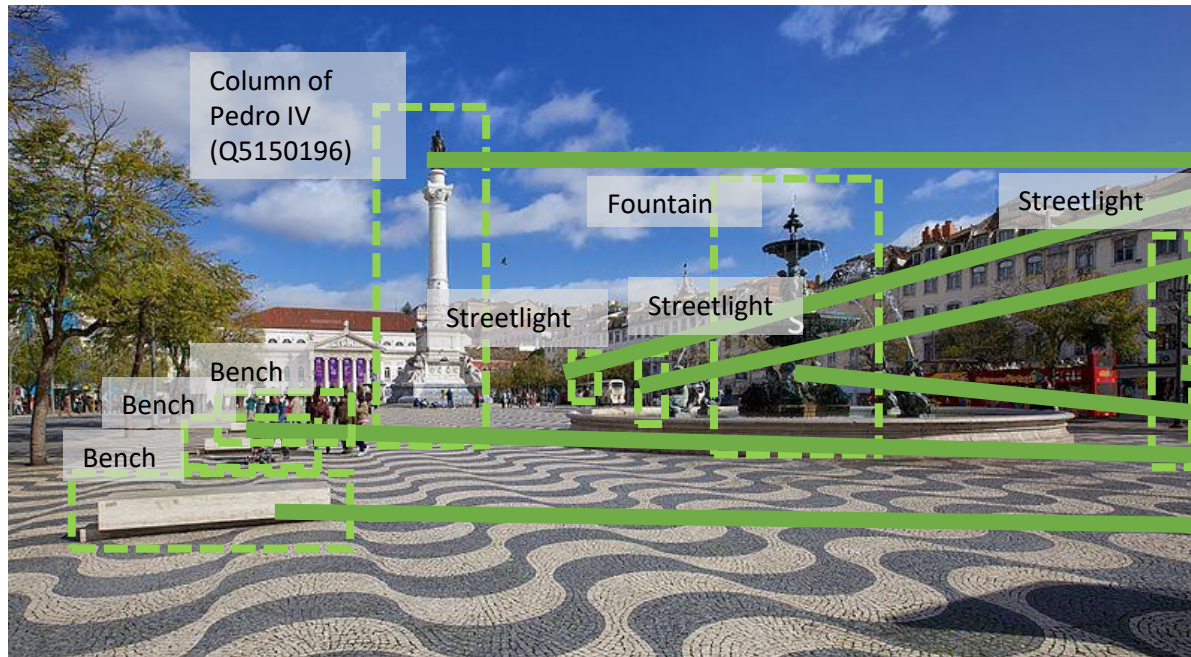


1. H. M. et al. "Locality-aware subgraphs for inductive link prediction in knowledge graphs" - Submitted to PRL2022.

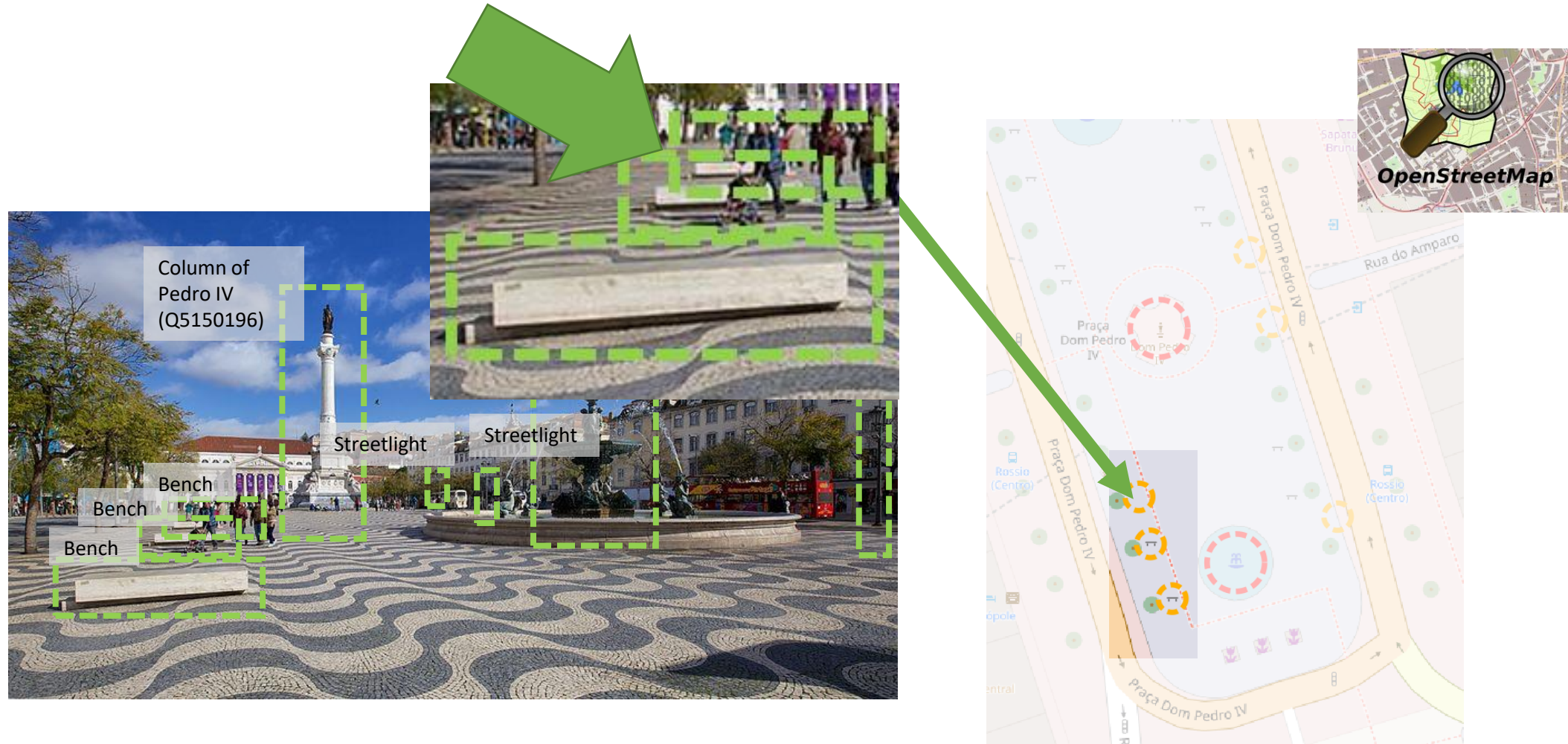
2. H. M. et al. "Geolocation of Cultural Heritage using Multi-View Knowledge Graph Embedding PatReCH 2022

Why do we care about benches?

Target scenario

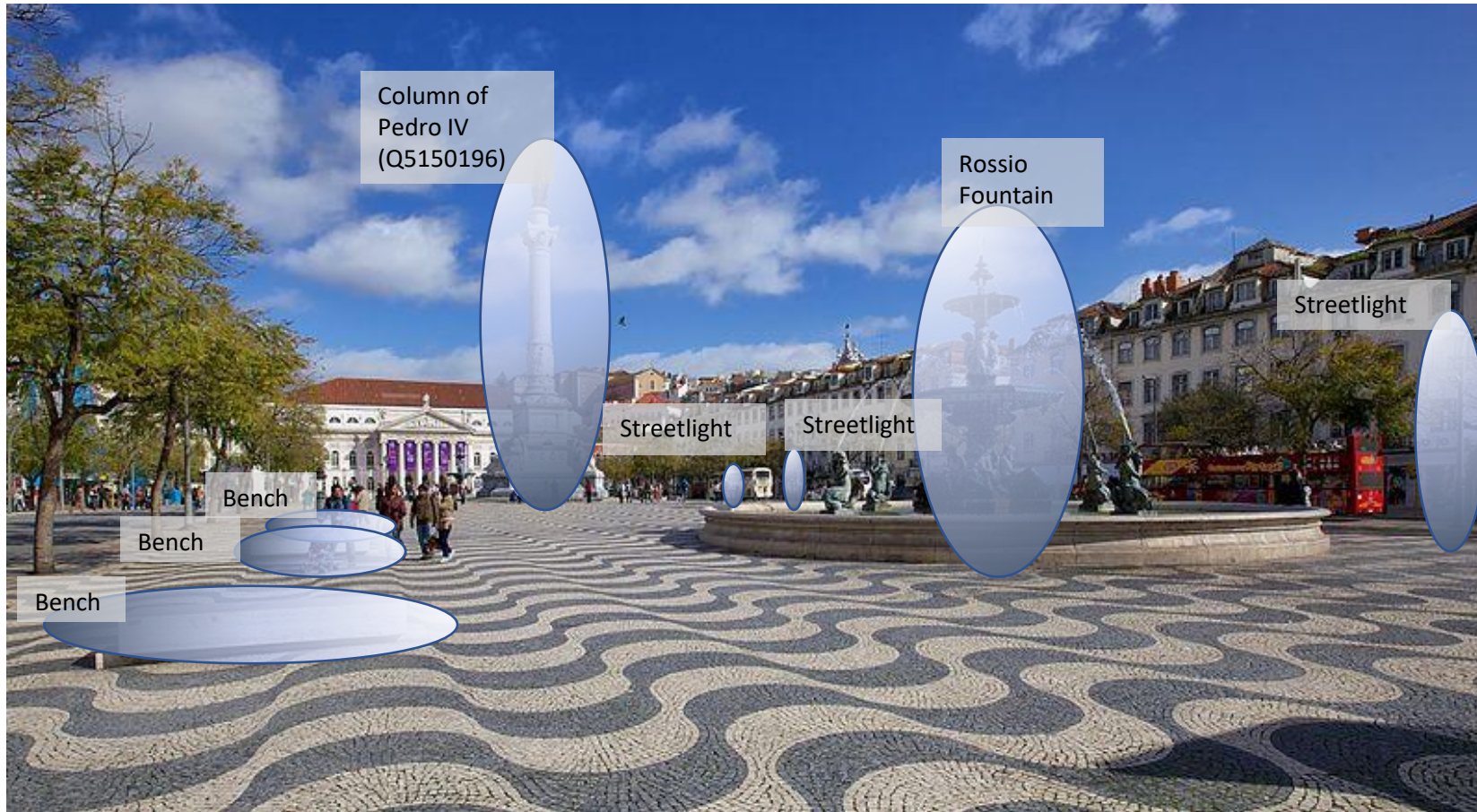


Why do we care about benches?



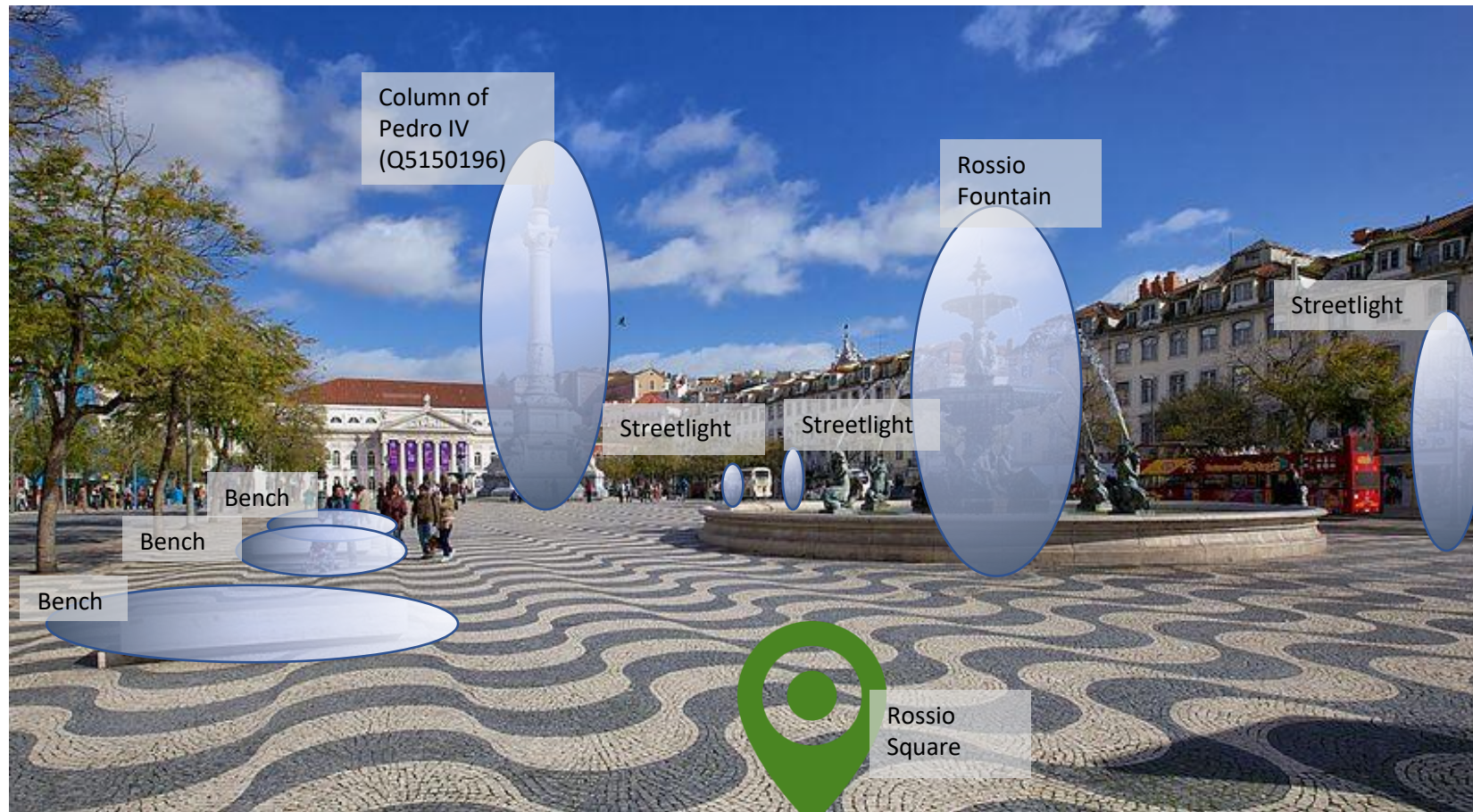
Where does this take us?

Self-localization using common landmarks!



Where does this take us?

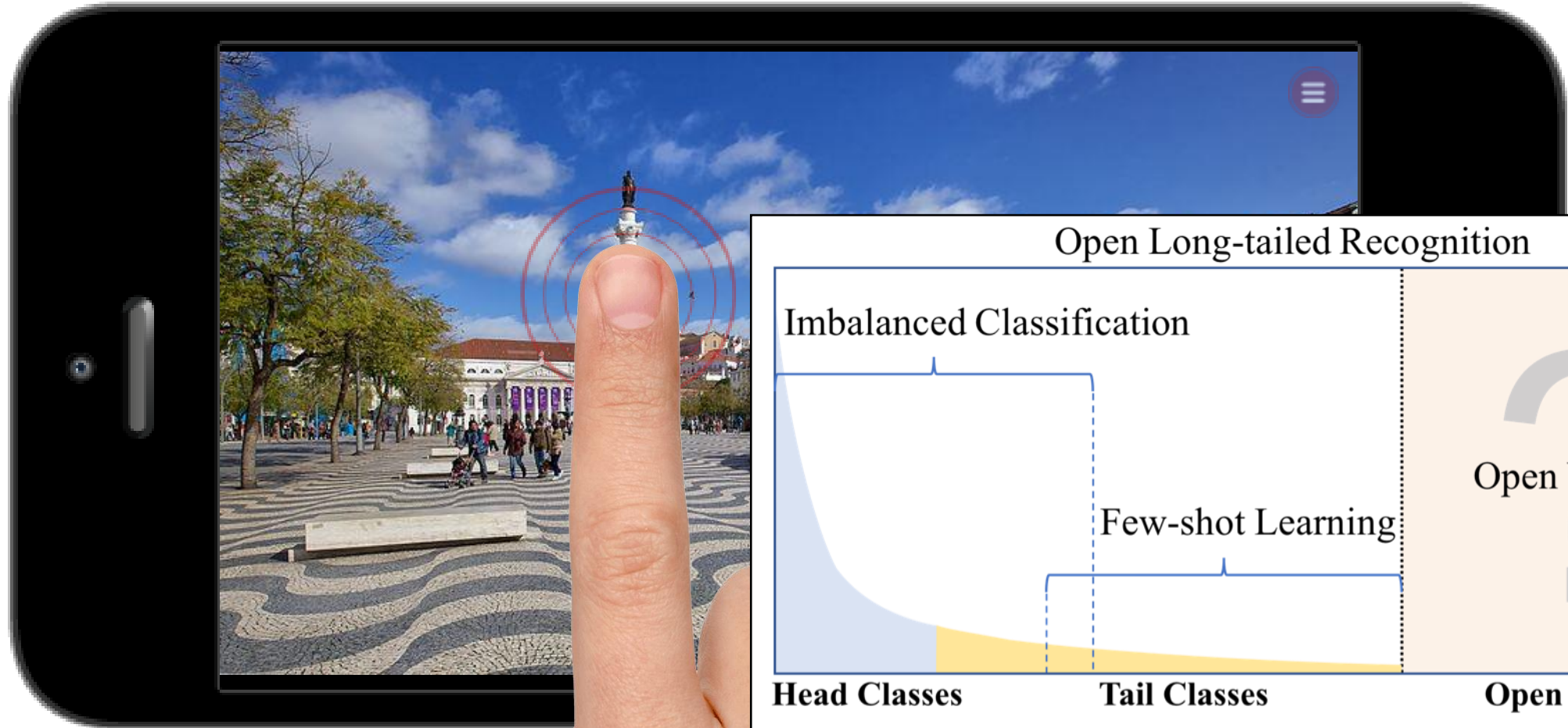
Self-localization using common landmarks!



How do we get to Annotations?



How do we get to Annotations?



Liu et al. "Large-Scale Long-Tailed Recognition in an Open World" CVPR'19

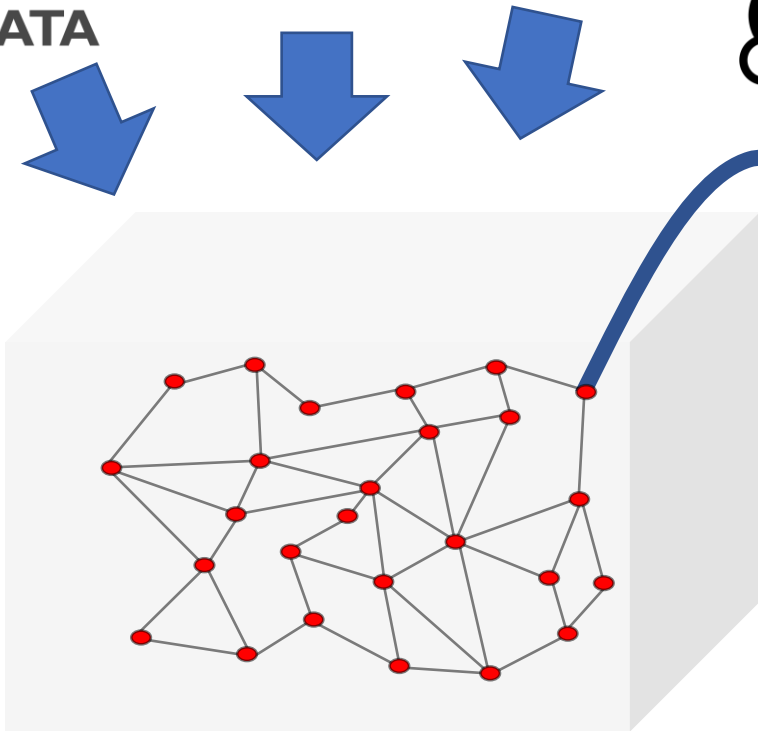
Position-based Object Retrieval



Multi-modal Object Retrieval

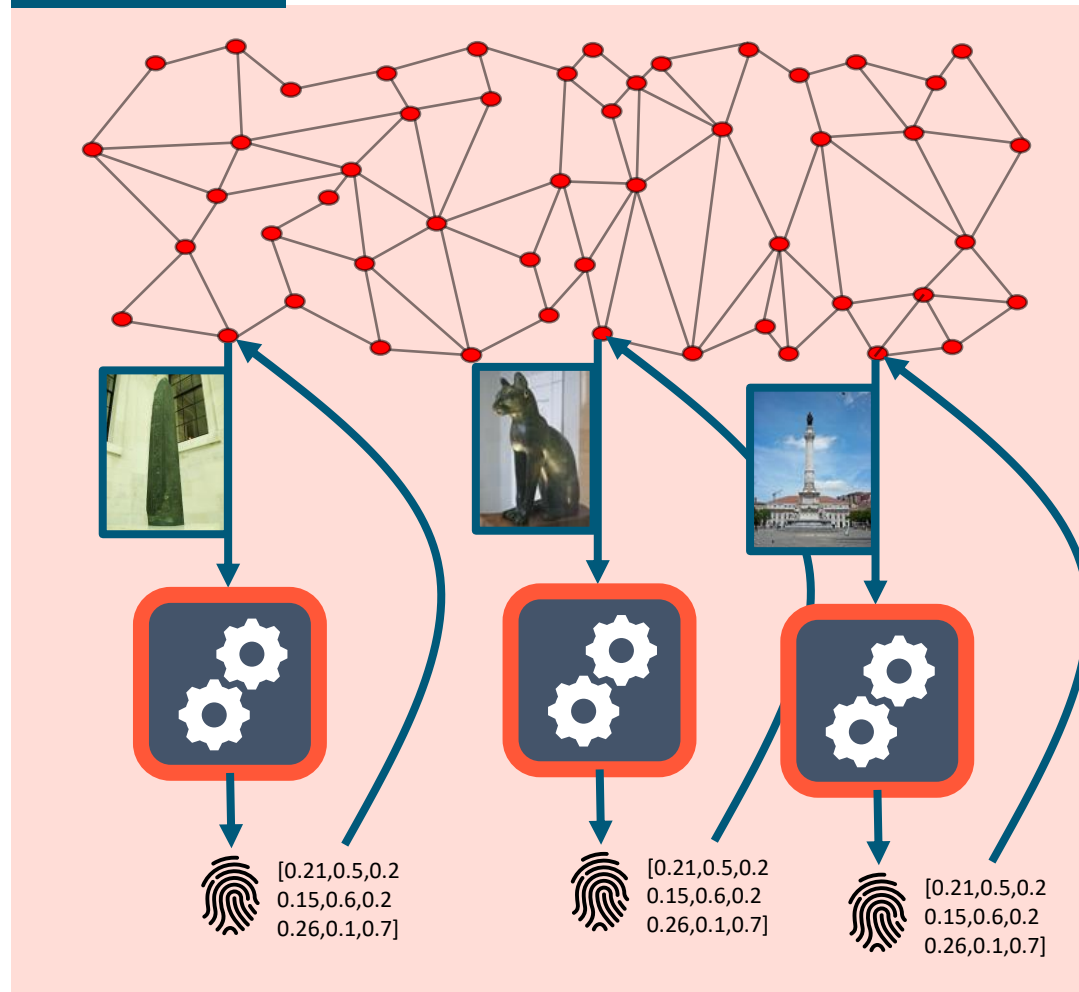


WIKIDATA

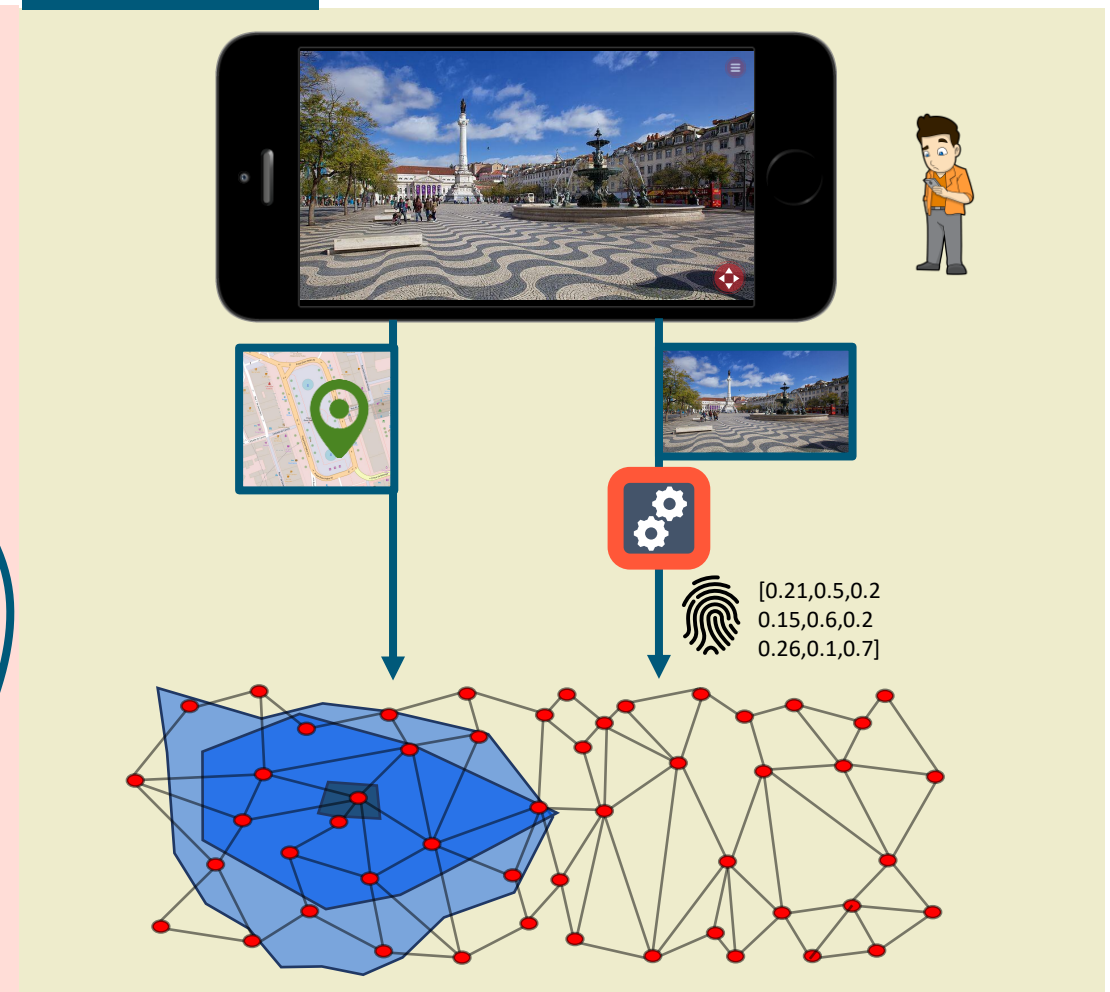


Multi-Modal Retrieval

Offline



At query time



Multi-Modal Retrieval

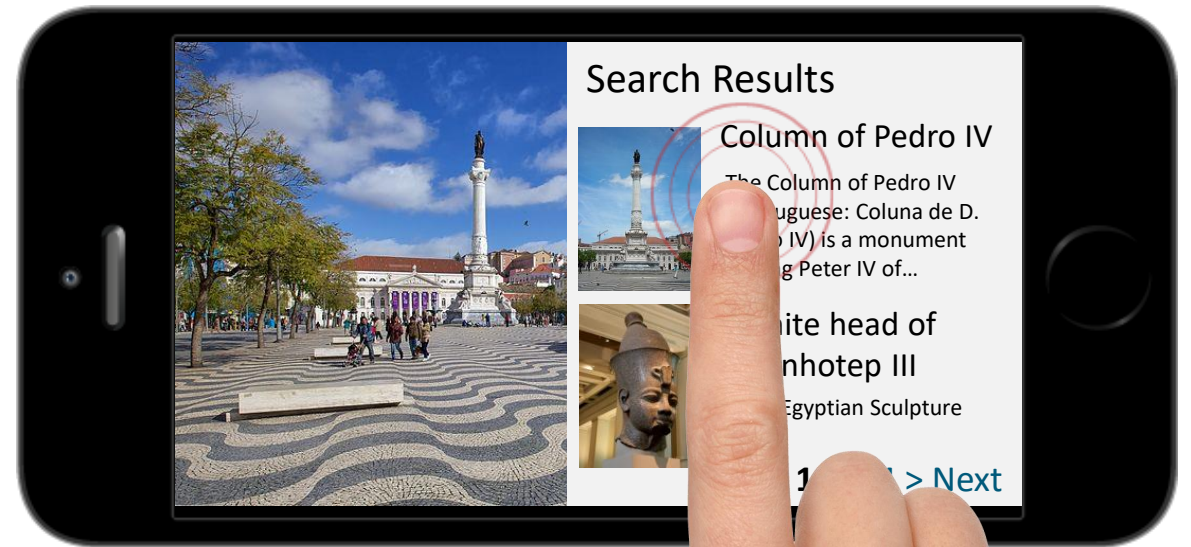


Location Only



Location and Image

Connecting real world to Knowledge Graph



Connected representations

Column of Pedro IV (Q5150196)

WIKIPEDIA
The Free Encyclopedia

Not logged in | Talk | Contributions | Create account | Log out

Article | Talk | Read | Edit | View history | Search Wikipedia

Wiki Loves Earth Italy is the photographic contest dedicated to the Italian protected areas. Take a photo, promote the territory, help Wikipedia and win!

Column of Pedro IV

From Wikipedia, the free encyclopedia

The **Column of Pedro IV** (Portuguese: Coluna de D. Pedro IV) is a monument to King Peter IV of Portugal and the Algarves, located in the centre of Rossio Square in Lisbon, Portugal. The monument was erected in 1870.

History and details [edit]

The first monument to King Peter IV, called "o galheteiro" ("the cruet-stand") by Lisboners, in a c. 1860 photograph.

King John VI had a monument to the Constitution built in 1821 on the spot the column stands today, which was rased two years later by the same king, after Infante Michael of Portugal (supported by Queen Carlota Joaquina) successfully led a counter-revolution to reinstate the absolute monarchy.

A first monument to King Peter IV was erected in 1852 with Queen Mary II (King Peter's

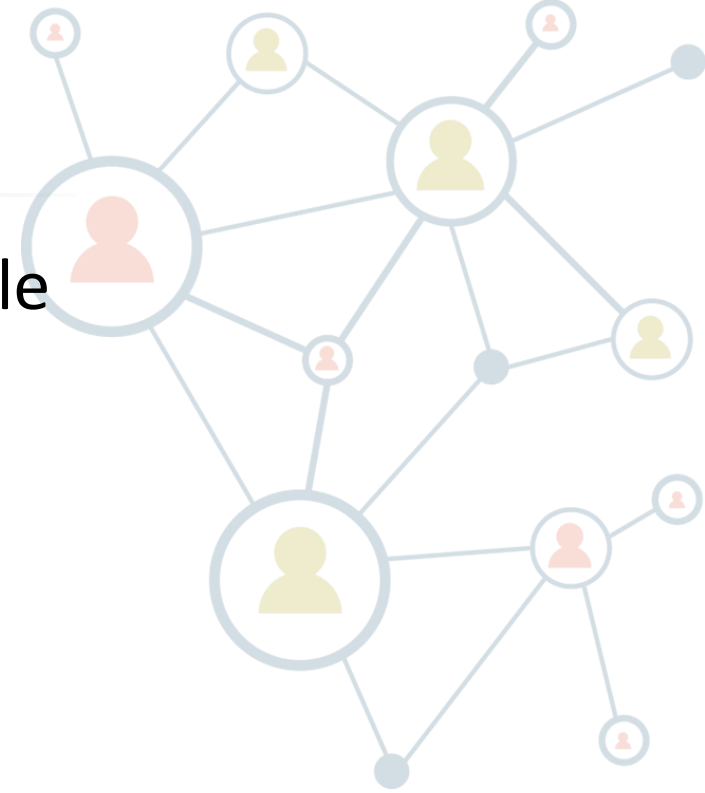
Column of Pedro IV
Coluna de D. Pedro IV

The Column of Pedro IV

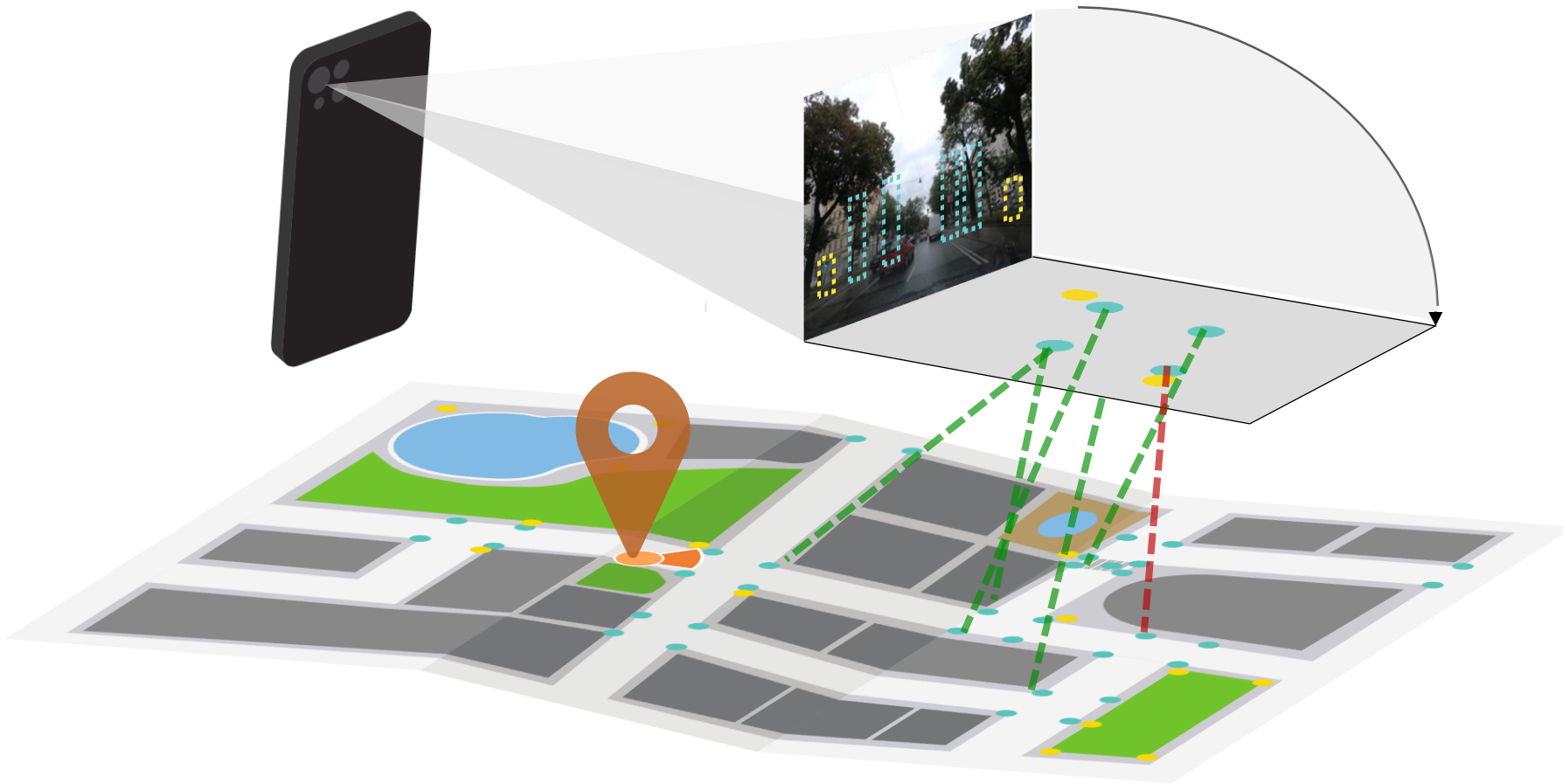
Location Rossio Square, Lisbon, Portugal

You are here!

Finding position and orientation on a 2D map from a single image: The Flatlandia localization problem and dataset



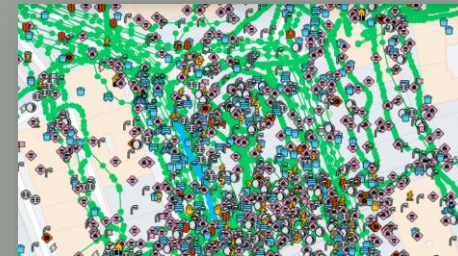
3DoF visual localization from object detections



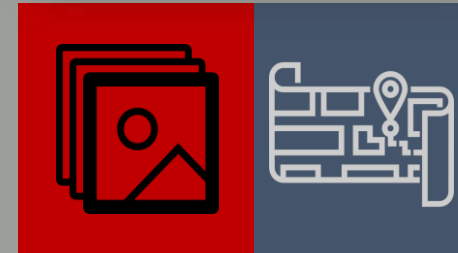
Why 2D maps?



Most static objects do not change appearance in time



We have large-scale object maps



Maps are Concise to store in contrast to Street-level images

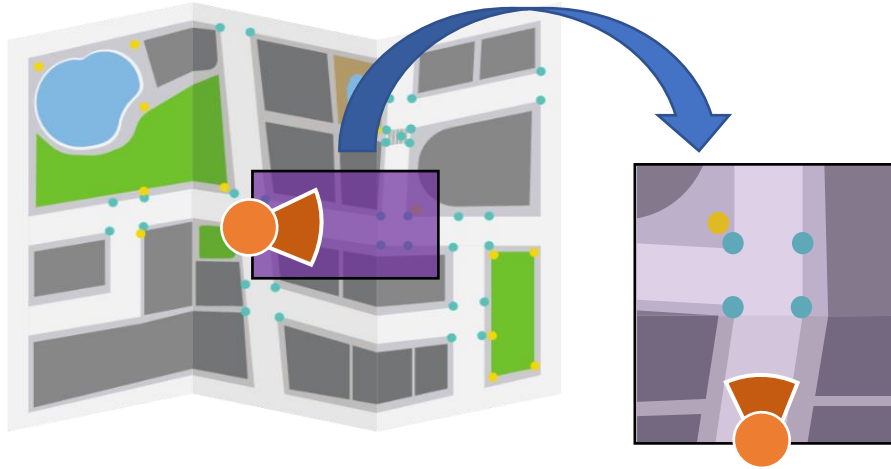


No privacy concerns over stored or transmitted data

Local Map Generation (phone view)

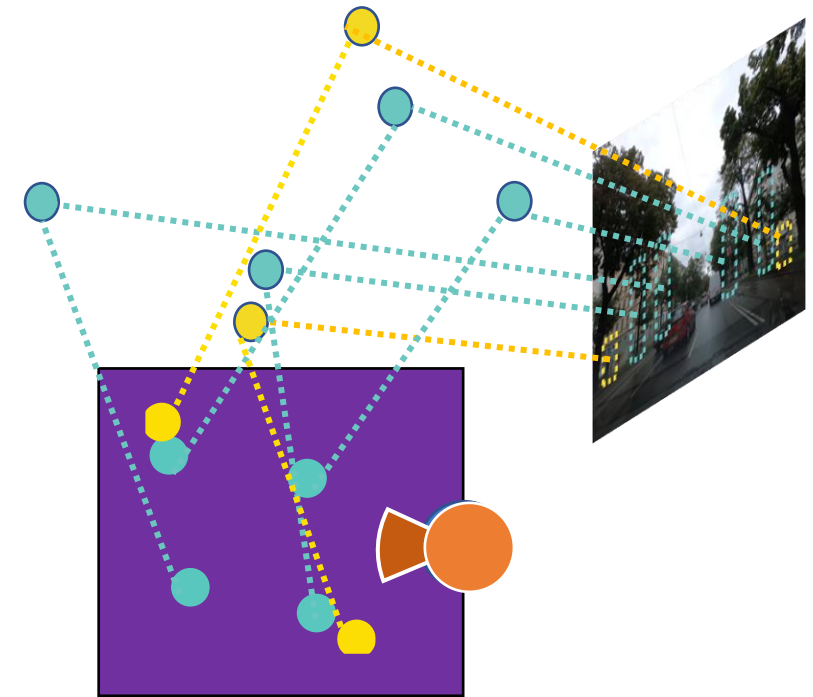


GT-based



Reference map object locations, as seen by the camera

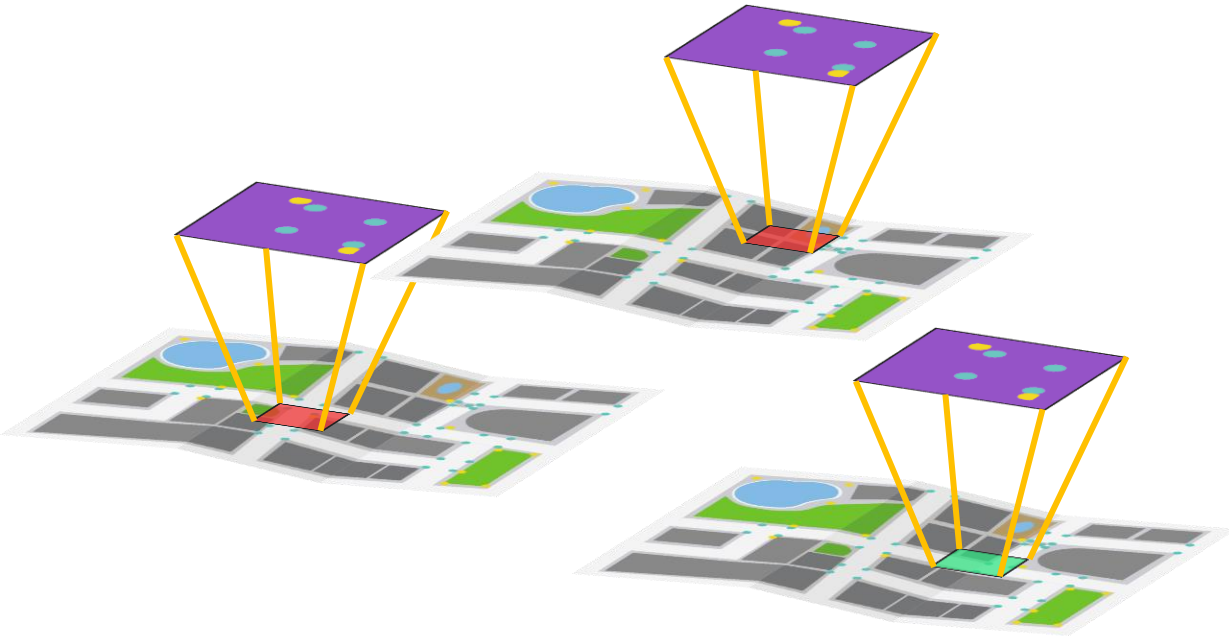
Depth-based



Detections projected using monocular depth estimation

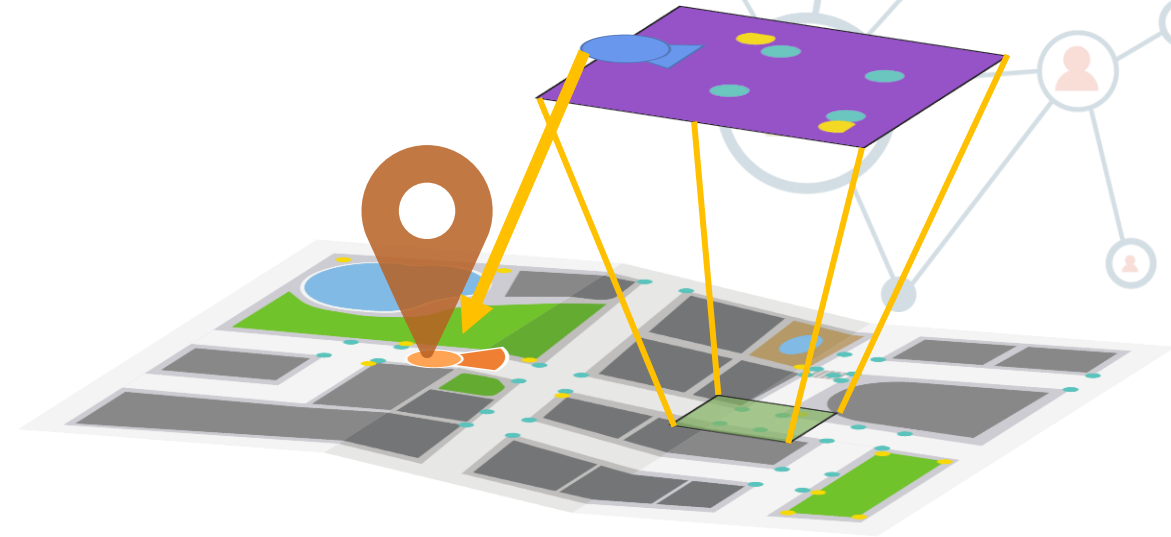
3DoF visual localization from object detections

Coarse Map Localization



Where on the reference map can I find the objects arranged like the ones I am seeing?

Fine-grained 3DoF Localization



Given a candidate reference map region, where is the camera on the reference map?

Coarse Map Localization - Results



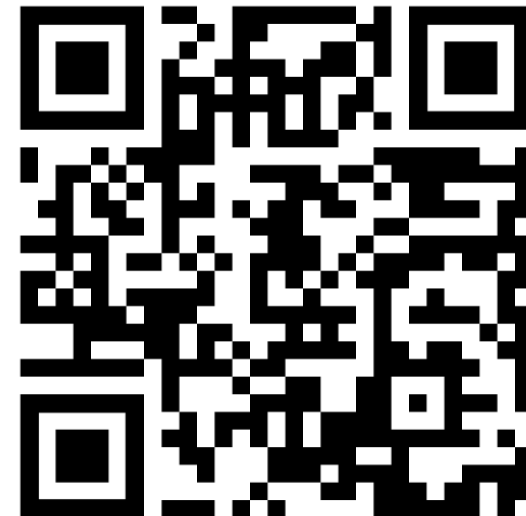
GT Local Maps			
Model	Precision	Recall	Success
Similarity	0.151	0.586	0.717
Triplet	0.161	0.559	0.693
Triplet + Similarity	0.206	0.617	0.804

More info here:

Paper: <https://arxiv.org/abs/2304.06373>

Github: <https://github.com/IIT-PAVIS/Flatlandia>

Depth Local Maps			
Model	Precision	Recall	Success
Similarity	0.149	0.594	0.720
Triplet	0.160	0.570	0.740
Triplet + Similarity	0.184	0.630	0.808



Thank you!

