Do we still need geometry for Visual Localization and Mapping?

Paul-Edouard Sarlin

50th Pattern Recognition and Computer Vision Colloquium - CVUT 2025-10-09

About me

Research scientist at Google Geo

PhD at ETH Zurich 2020-2024

More info: psarlin.com





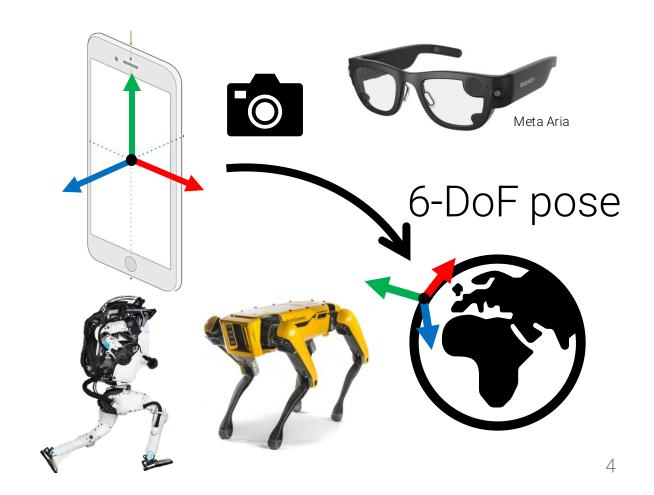


Do we still need geometry for Visual Localization and Mapping?

Visual Localization and Mapping

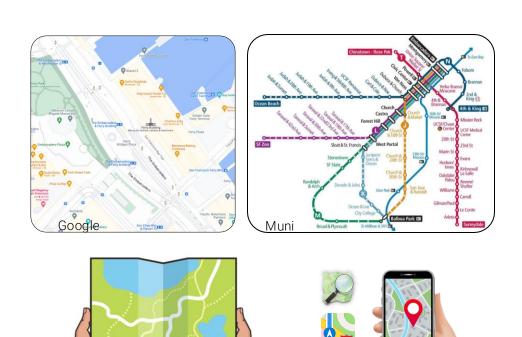
Finding where I am

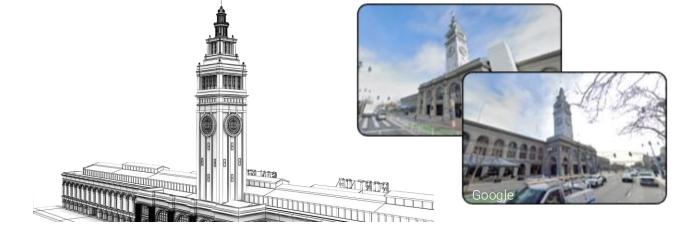




Visual Localization and Mapping

What is in the world and where









Visual Localization and Mapping

Creating maps from observations

needs the sensor pose

mapping



localization



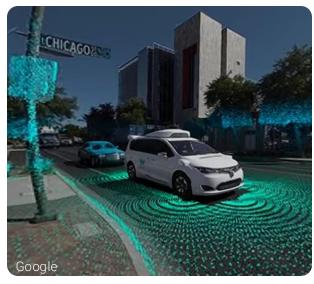




Why should we care?



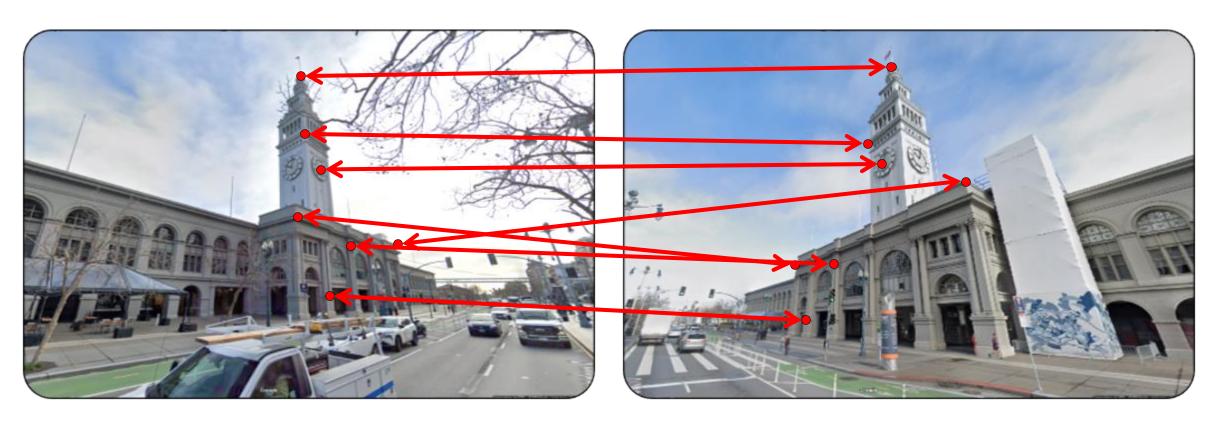






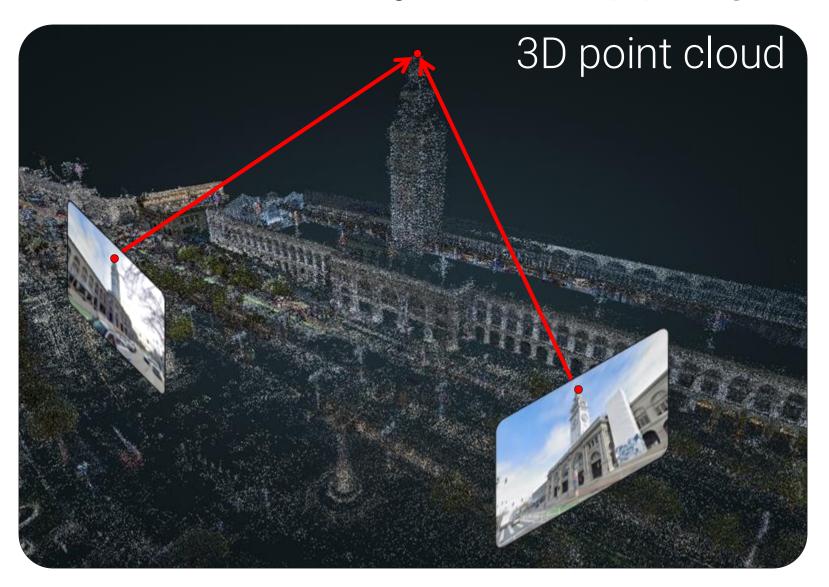


3D Geometry for Mapping



mapping images

3D Geometry for Mapping



3D Geometry for Localization

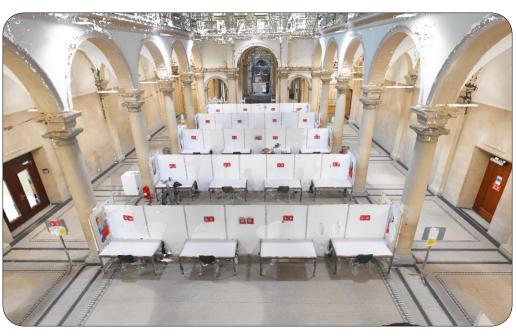


3D point cloud

query image

Real-world challenges





temporal changes







symmetries





Do we still need geometry for Visual Localization and Mapping?



We now have end-to-end learning that works

2016

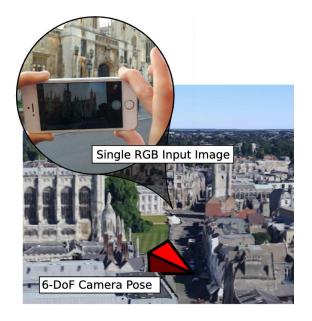


2023

2025

fast, robust, generalizes

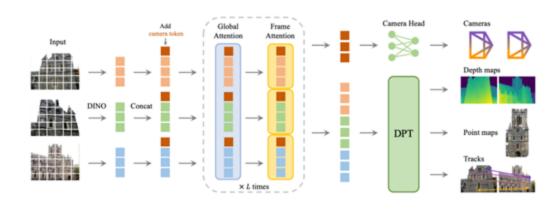
PoseNet



DUSt3R



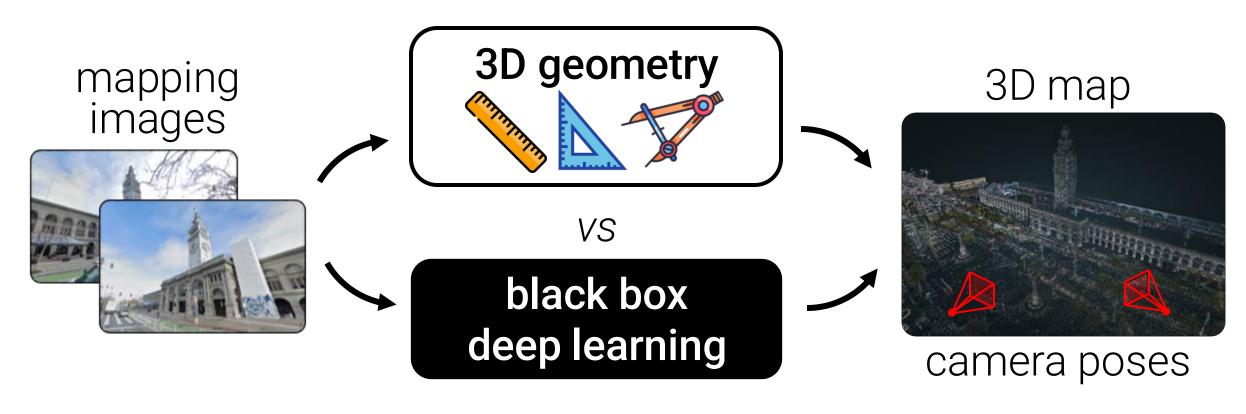
VGGT



But what did we lose?

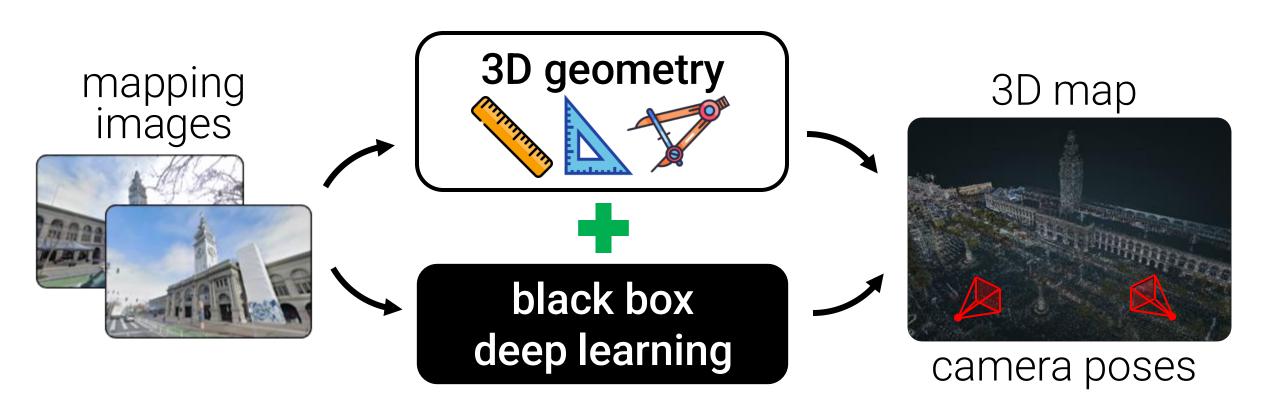
From Geometry to Deep Learning?

accurate, scalable, interpretable



none of the above! limited to <<1k views

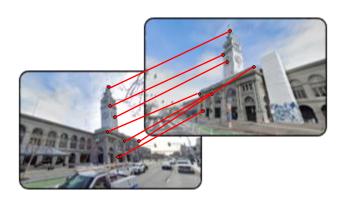
From Geometry to Deep Learning?



Do we have to choose one?

One step at a time

image matching





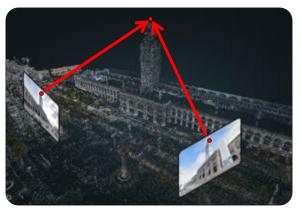
SuperGlue [CVPR'20] LightGlue [ICCV'23] pose estimation





GeoCalib [ECCV'24]

bundle adjustment



PixSfM [ICCV'21] calibration

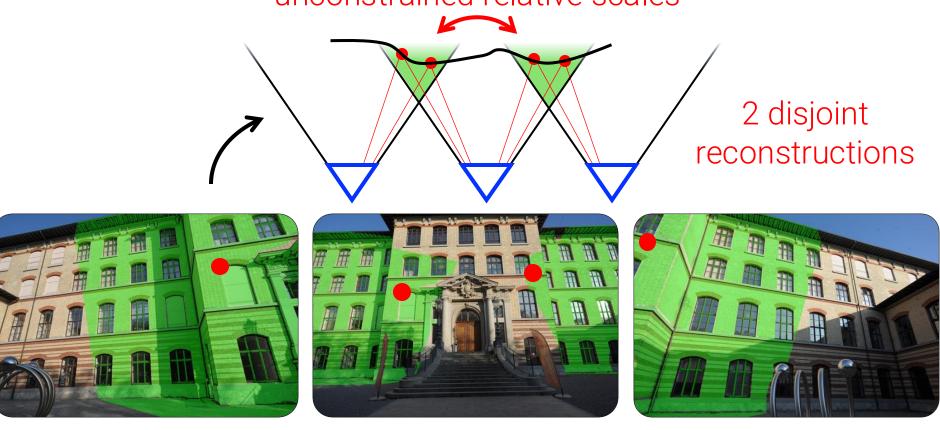


A simple recipe

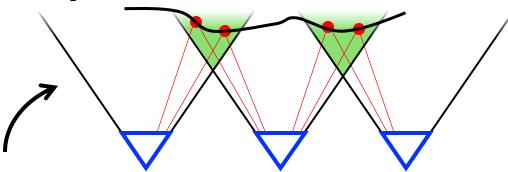
- 1. Identify **structural weaknesses** in geometry
 - unconstrained? unobservable? noisy inputs?
- 2. Think ML: what **priors** do we need? Inferred from what **data**?
- 3. What **geometry** shouldn't be learned? Camera models!
- 4. Bake it into an optimization problem with learned data terms
 - Learned priors as data terms
 - Flexible constraints
 - Predictive uncertainty

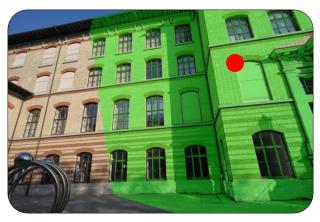
Example: three-view overlap

unconstrained relative scales



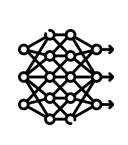
Example: three-view overlap

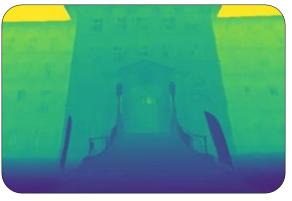


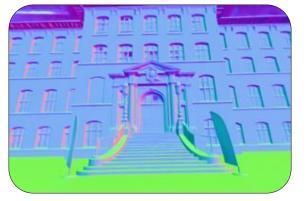






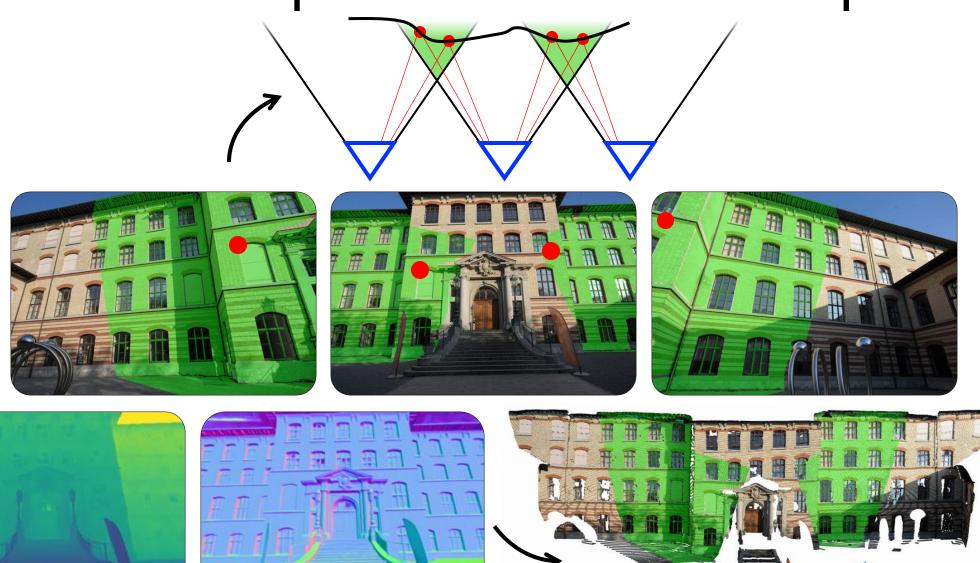






Monocular surface priors depth & normals

Example: three-view overlap







MP-SfM Monocular Surface Priors for Robust Structure-from-Motion

CVPR 2025



Zador Pataki¹



Paul-Edouard Sarlin²



Johannes Schönberger^{1,3}

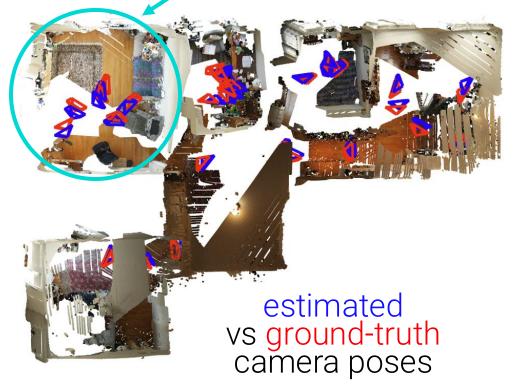


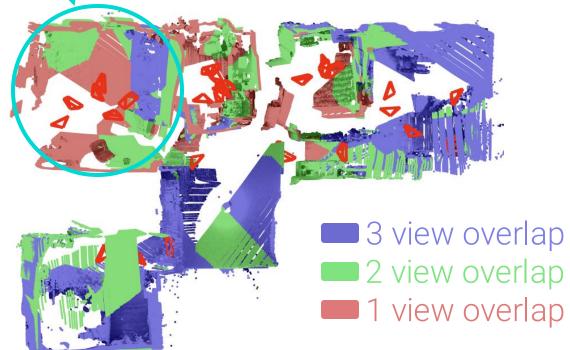
Marc Pollefeys^{1,3}

github.com/cvg/mpsfm

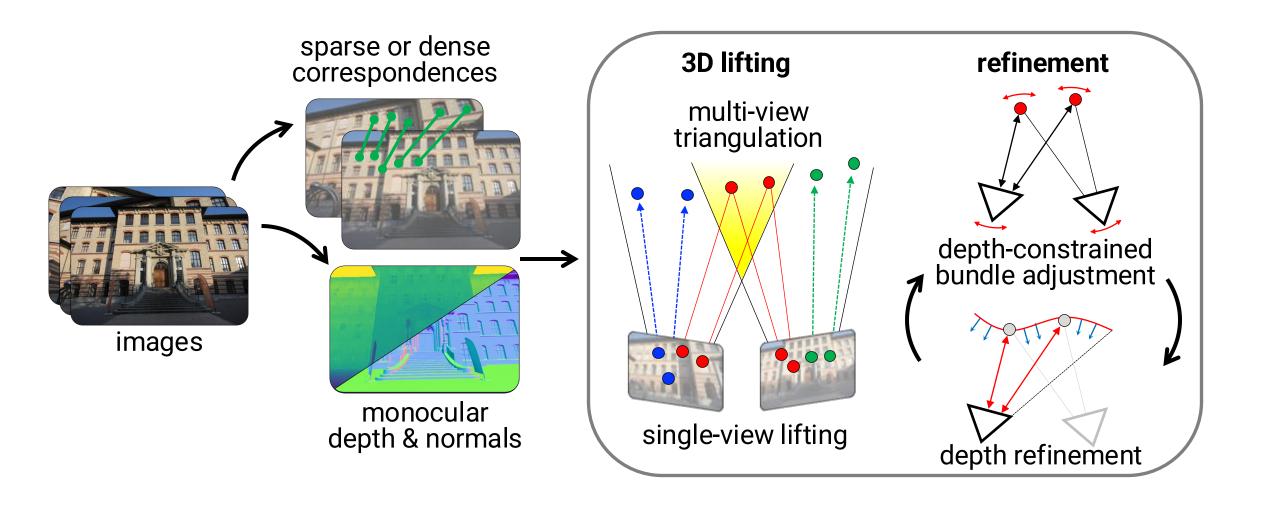


Common scenario for non-expert users

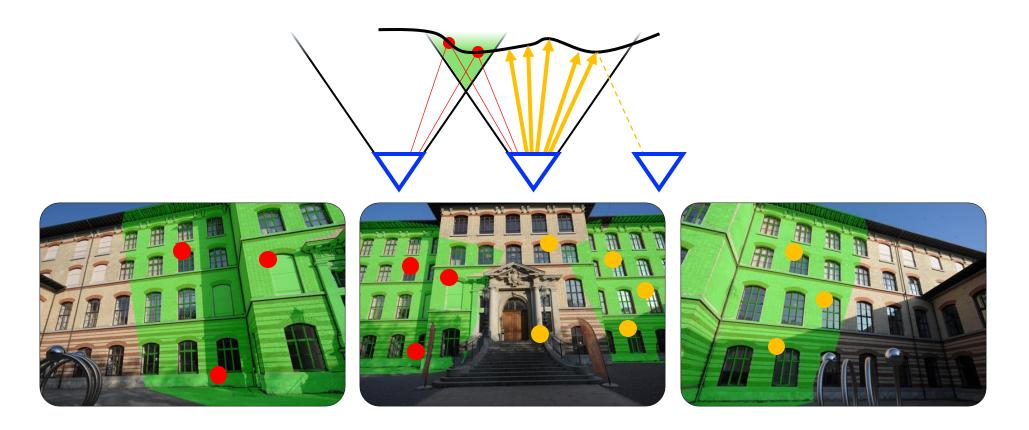




How it works



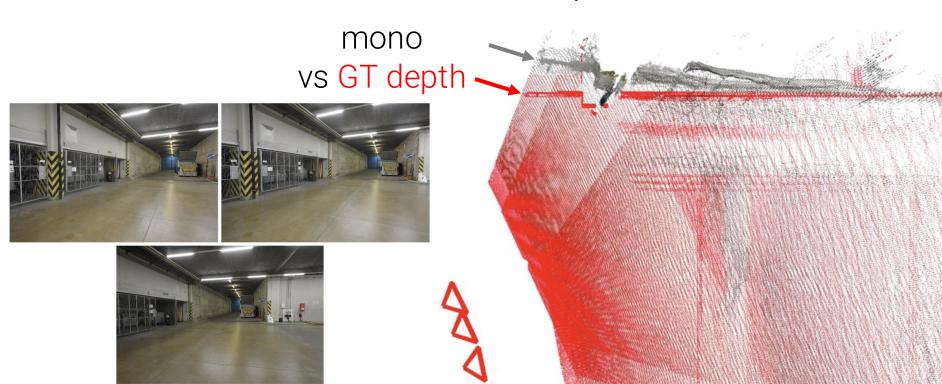
Single-view lifting with depth

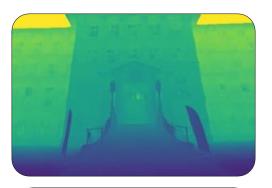


Complement regular multi-view 3D points with single-view tracks from depth

Bundle adjustment with priors

- Naïve: constrain the depth of 3D points
- But: monocular depth is rarely accurate!
- Surface normals are easier to predict

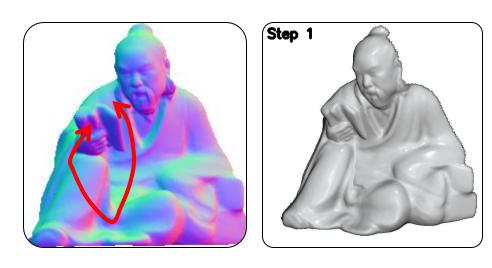




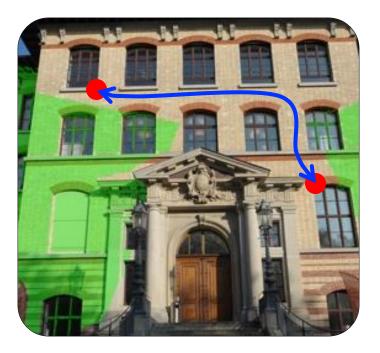


Depth refinement as normal integration

- Constrain 3D points that are on the same surface
- Optimize a dense depth map with NLS
- Estimate the depth discontinuity with IRLS

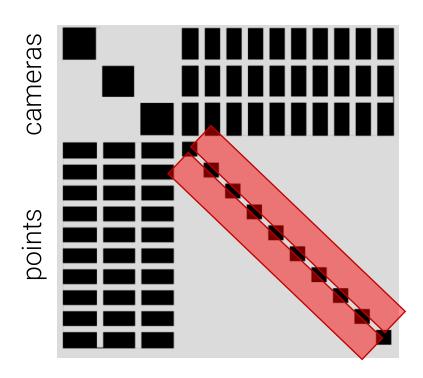






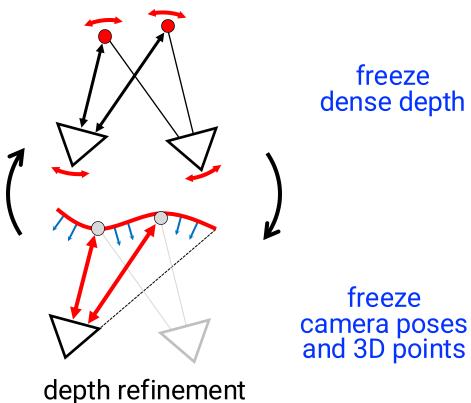
Depth refinement as normal integration

point-to-point constraints break the Schur Complement

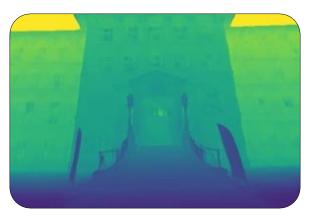


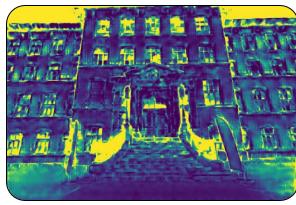
Solution: alternate optimization

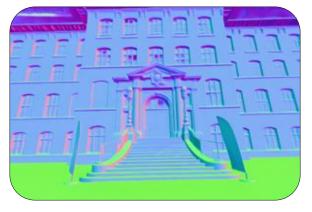
depth-constrained bundle adjustment



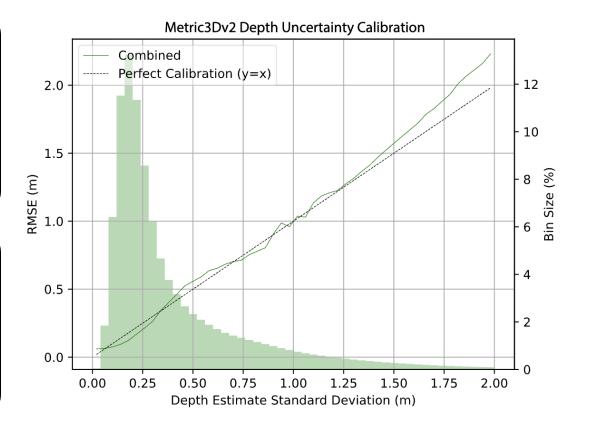
Joint optimization requires calibrated uncertainties

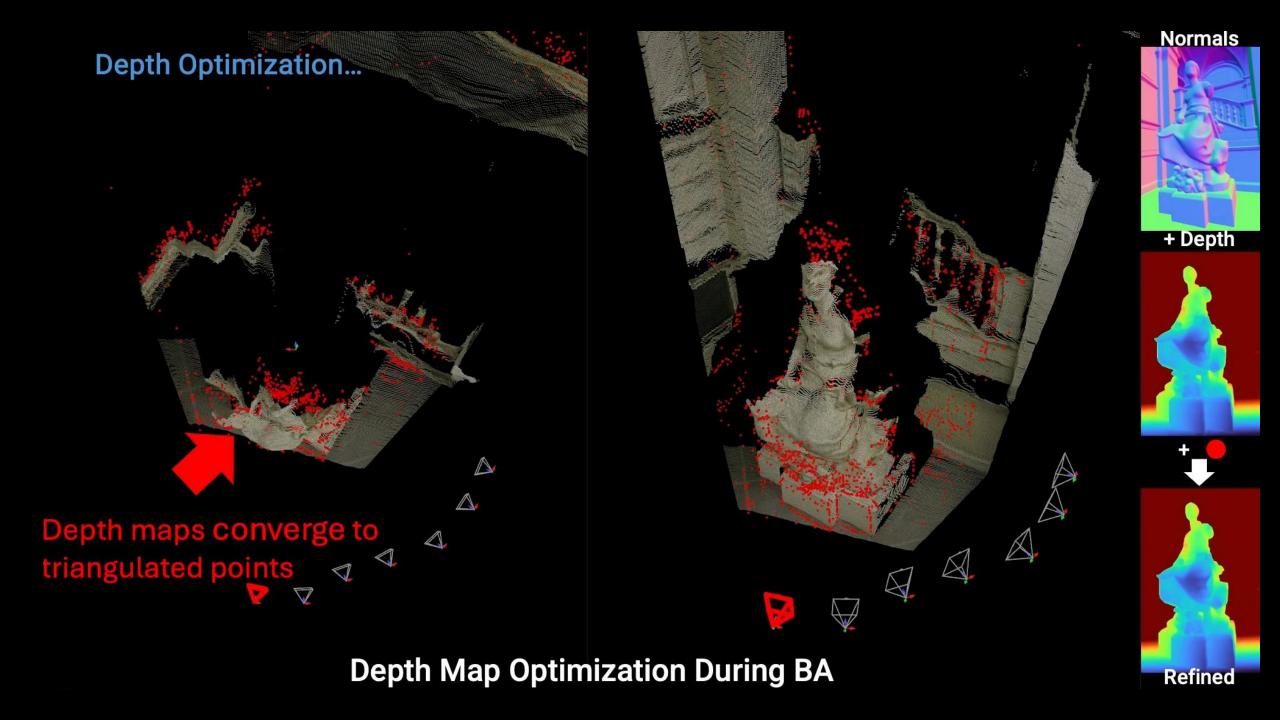


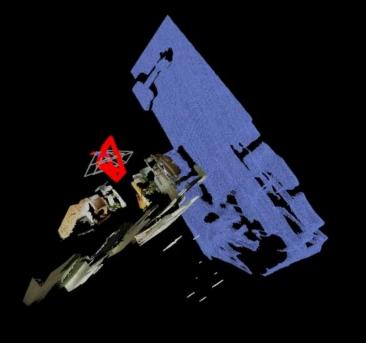












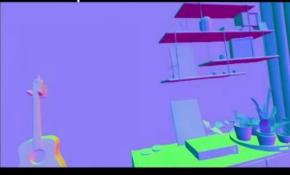
Latest Image



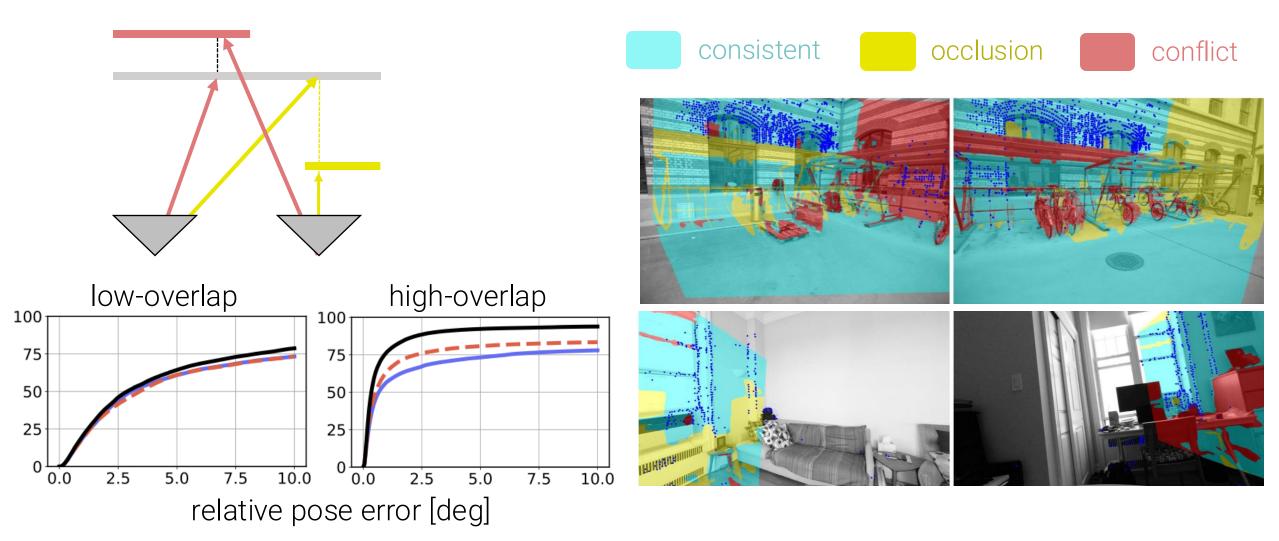
Input Depth



Input Normals

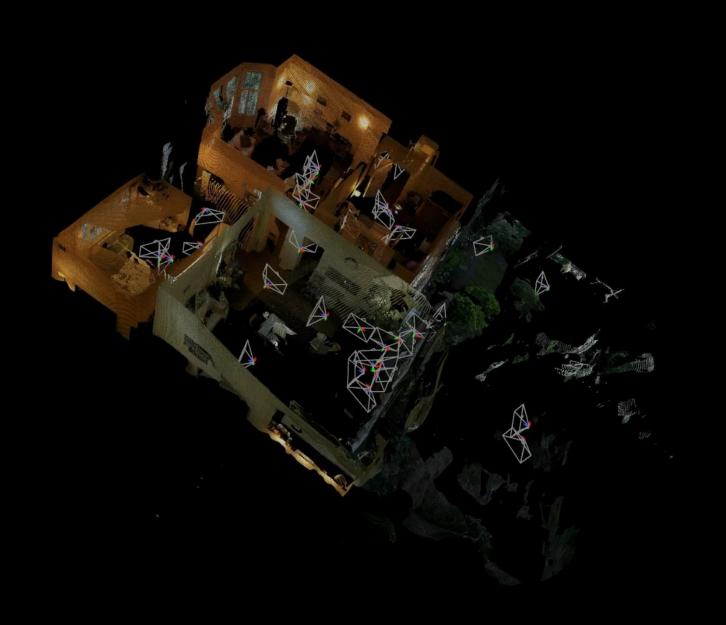


Symmetry detection with depth consistency



vanilla / with Doppelganger / with depth

MP-SfM (ours)



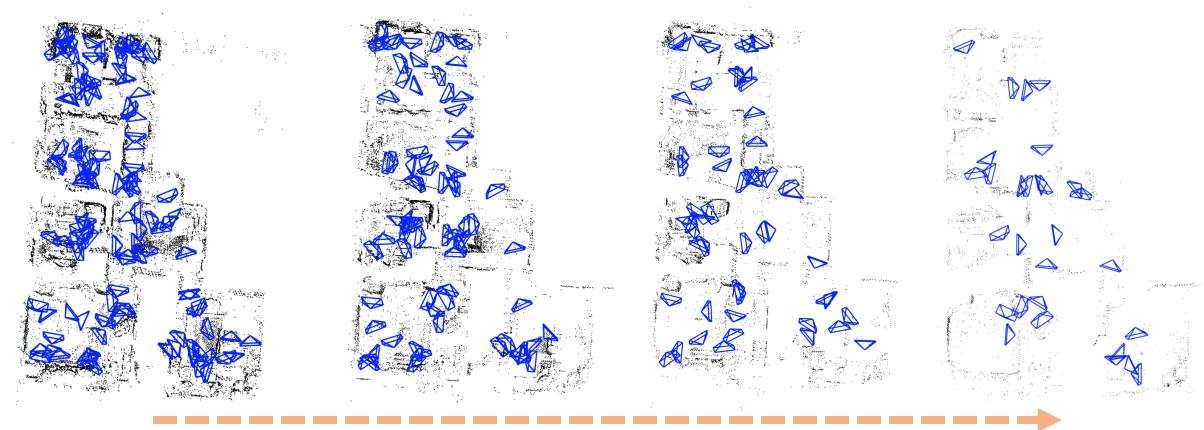
So... do we still need geometry?

YES!

- Two/multi-view models learn strong 3D priors
- But we already have this from off-the-shelf models
- Use geometry to glue them together!
- Open question: can we learn depth/normals end-to-end in SfM?
- Grand vision: a self-tuning system, reinforced with geometry "Self-supervised VGGT" with hard geometric constraints

Beware of artificial evaluations!

- 1. Take nice, slow videos with a lot of parallax
- 2. Then extract images at regular interval



Much of real data is very different

- Example: egoentric data from wearable devices
- Opportunistic capture, byproduct of the user's motion









Challenges & opportunities

- Unconstrained motion, moving environments
- Appearance: low light / exposure changes
- Long up-time → time-varying calibration

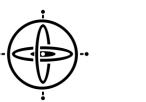
But

- Non-image sensors: IMU, GPS, microphones, WiFi/BT
- High-frequency (>60 FPS)
- Dense coverage (crowd-sourced)













Bluetooth



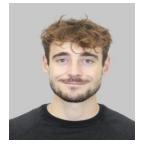
WiFi

Benchmarking Egocentric Visual-Inertial SLAM at City Scale



















Anusha Krishnan^{1*} Shaohui Liu^{1*} Paul-Edouard Sarlin^{2*} Oscar Gentilhomme¹ David Caruso³ Maurizio Monge³ Richard Newcombe³ Jakob Engel³ Marc Pollefeys^{1,4} ¹ETH Zurich ²Google ³Meta Reality Labs Research ⁴Microsoft Spatial AI Lab



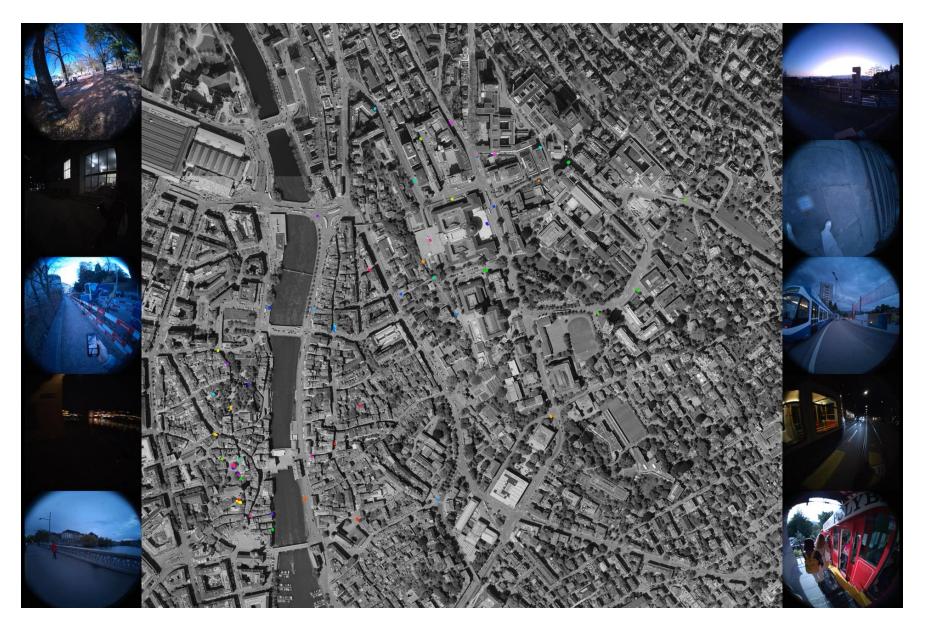






lamaria.ethz.ch

The LaMAria dataset





70km 22 hours 10-45min each

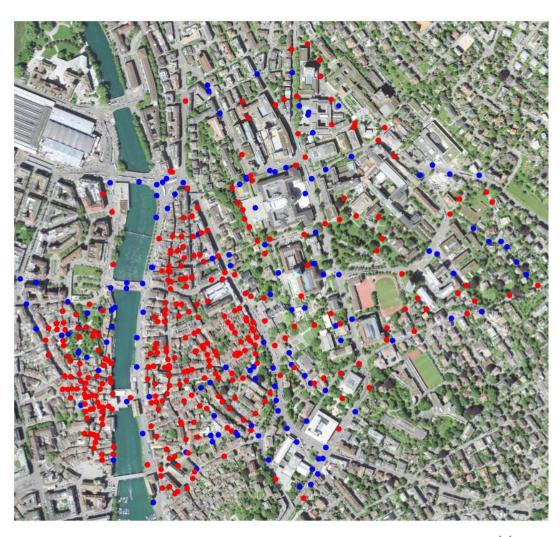
Good for multi-sensor SLAM and monocular SfM

The LaMAria dataset

Control points with cm accuracy Measure 0.1% metric scale drift







Monocular SLAM only works on easier data

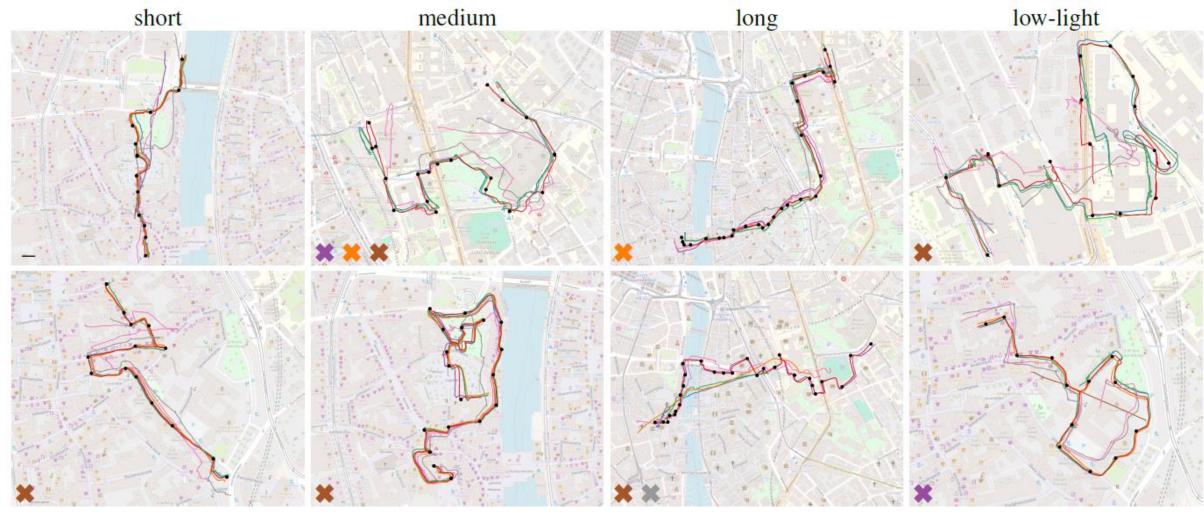


increasingly natural motion patterns

handheld

egocentric

Open-source systems are far behind industry



OpenVINS • OpenVINS+Maplab • ORB-SLAM3 • OKVIS2 •

Kimera VIO • DPVO • DPV-SLAM • Aria's SLAM •

Open-source systems are far behind industry

- monocular
- monocular inertial
- multi-camera inertial

method	causal	short			medium			long			challenge - low-light			challenge - moving platform		
			CP@1m↑	R@5m↑	score ↑	CP@1m↑	`R@5m↑	score ↑	CP@1m↑	R@5m↑	score ↑	CP@1m1	`R@5m↑	score ↑	CP@1m↑	R@5m↑
DPVO DPV-SLAM	√ ✓	9.4 7.2	1.7 1.4	21.3 14.5	5.2 5.2	1.0 1.4	10.8 10.0	1.2 0.4	0.0	1.9 0.6	3.4 1.9	0.2 0.4	7.5 3.5	2.4 1.7	0.1 0.0	
Kimera VIO ORB-SLAM3 OpenVINS OpenVINS + Maplab	✓ x ✓ x ✓ x	6.3 28.3 18.1 22.9	2.9 13.4 4.4 8.1	12.6 67.1 45.7 50.8	6.6 20.3 10.9 13.1	1.7 4.4 2.3 4.1	15.1 57.0 27.9 29.0	6.3 14.2 4.7 5.8	1.7 2.3 0.5 1.3	14.3 40.6 12.3 13.3	4.2 6.2 7.9 9.6	2.7 0.6 2.4 2.9	6.4 12.5 17.6 19.3	7.1 15.7 2.4 3.7	1.6 4.1 0.6 1.2	- - - -
OpenVINS OpenVINS + Maplab OKVIS2	✓ x ✓	22.2 26.0 24.3	6.2 9.5 12.0	57.9 61.1 54.8	17.8 21.3 13.6	5.7 7.3 6.8	46.1 50.6 28.2	10.6 12.6 2.3	1.7 1.9 2.5	25.8 30.3 1.7	16.9 16.5 15.4	6.2 4.6 5.3	38.2 37.9 38.6	11.5 13.0 4.1	2.4 3.0 2.8	- - -
Aria's SLAM	Х	90.7	99.2	_	78.5	87.4	_	70.8	75.9	_	84.2	91.6	-	53.6	51.2	_

Aria's MPS SLAM (close-sourced) is far ahead of all current academic solutions.

Moving the goalpost for pose Transformers

- How to maintain scale consistency over million of frames?
- Efficient joint inference over multiple sequences?
- How to use sensors like IMU in e2e models?
- Ubiquitous symmetry indoors
- GT is orders of magnitude more accurate
- Plenty of extra sequences for SSL



Conclusion

- Yes, e2e models are **impressive**!
- But e2e is not a requirement to leverage 3D priors
- And e2e is sometimes a **liability** e.g. sensor fusion is harder
- Open question

How to retain the benefits of geometry while learning even higher-level priors e2e? accurate, scalable, interpretable, flexible w.r.t. sensors

Thank you!

psarlin.com