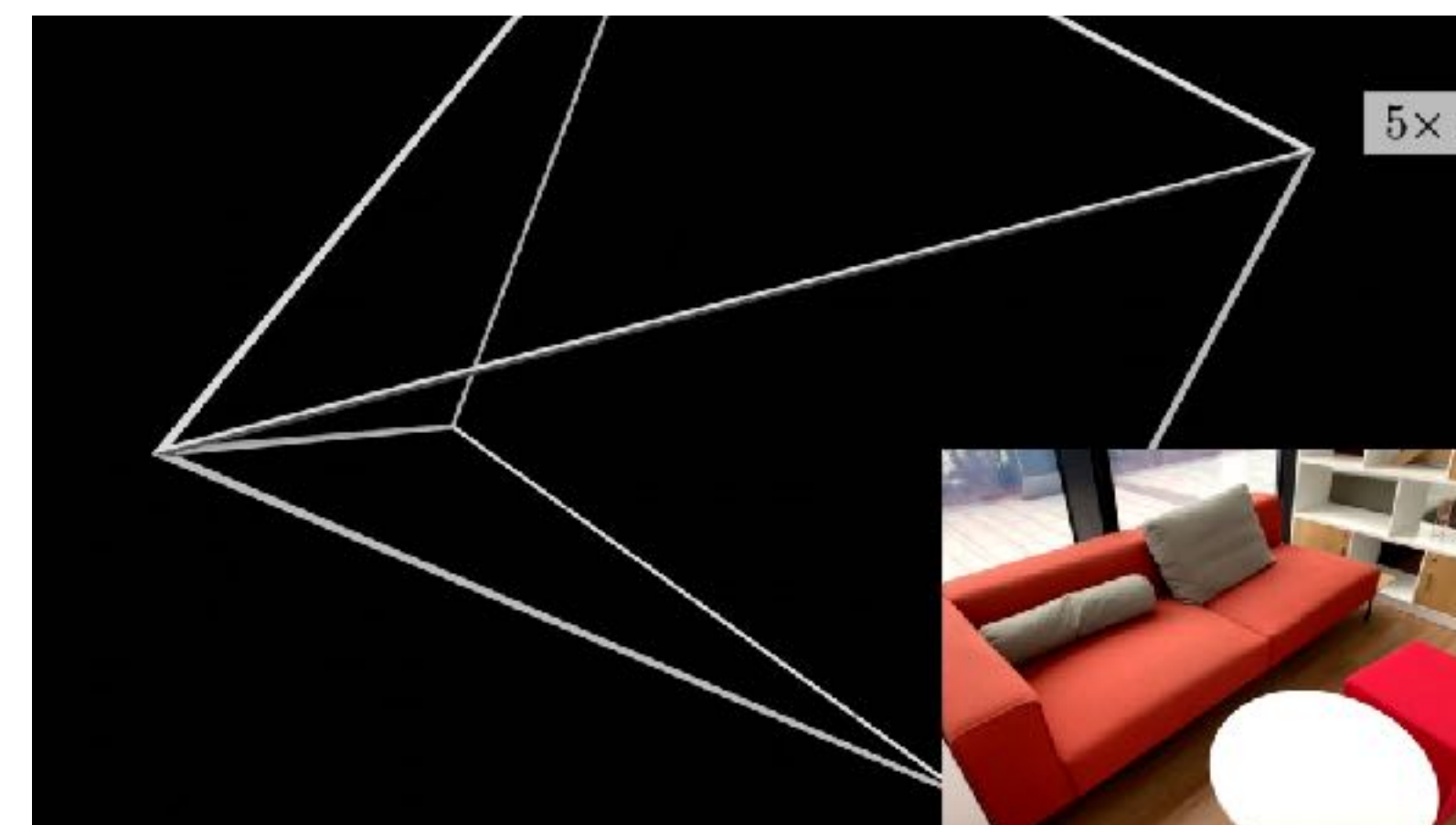
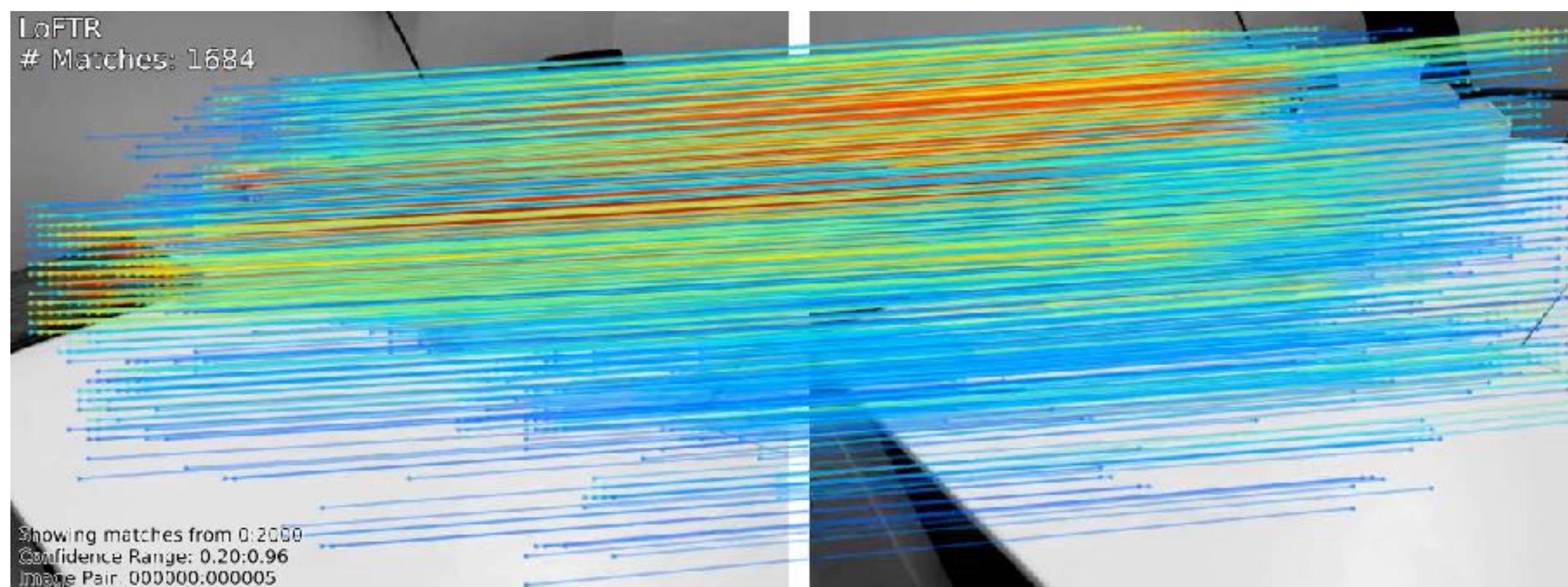




商汤
senseTime

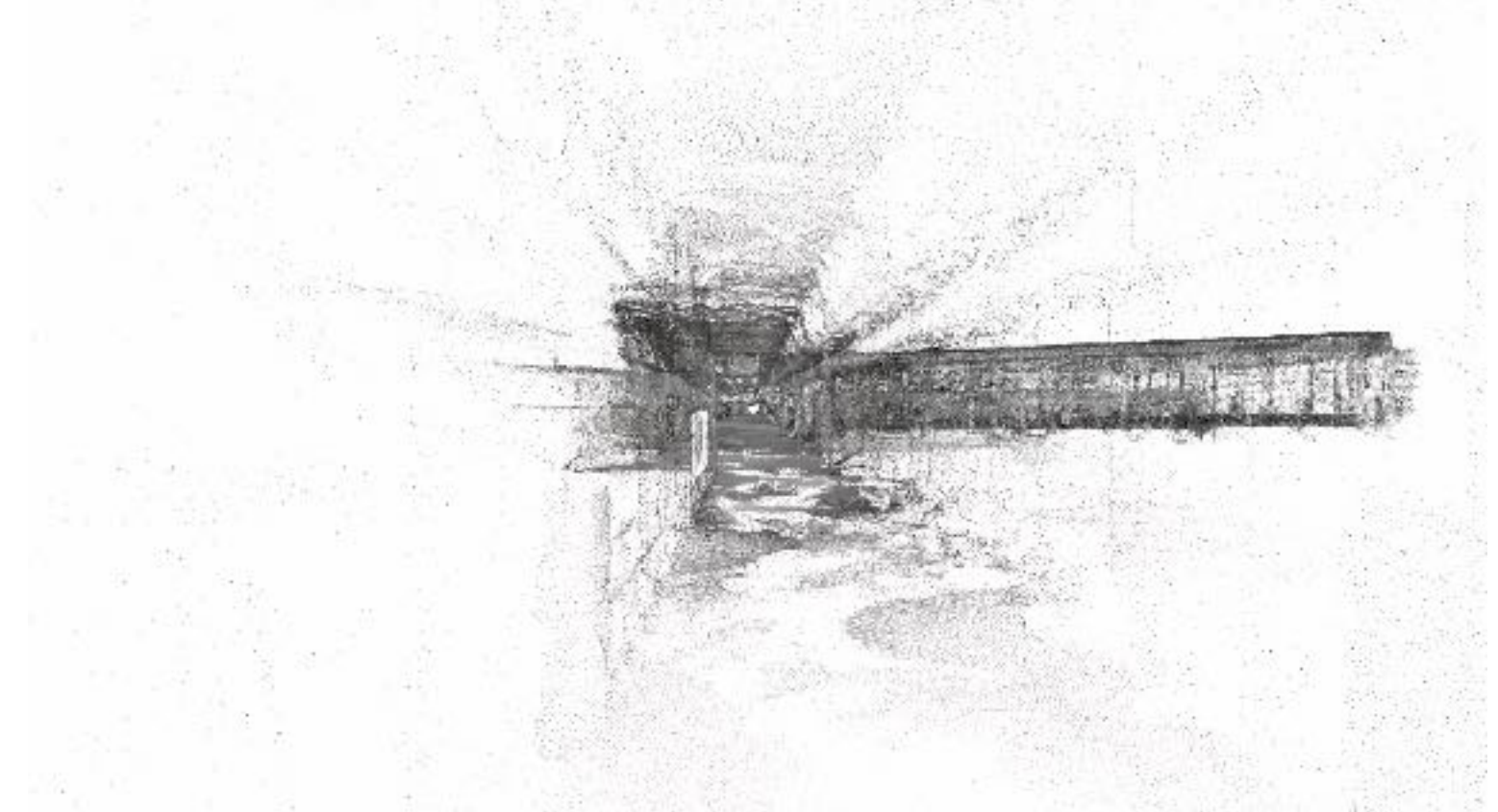
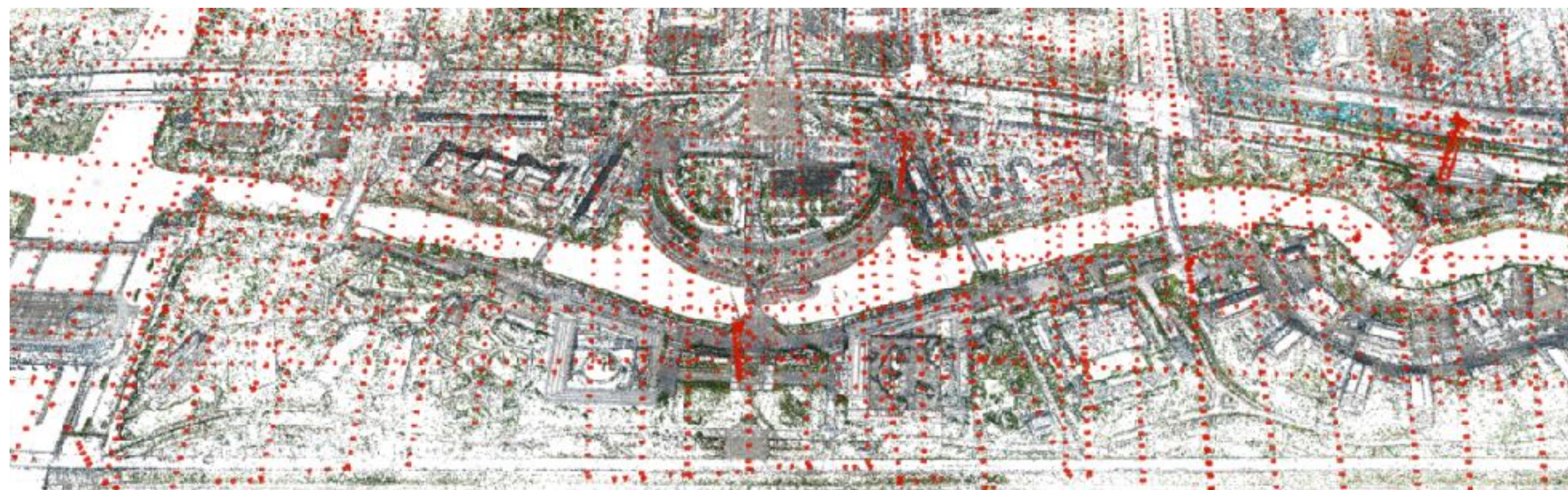
Learning-based 3D **Reconstruction** and **Localization**: What Should be Learned?

Jiaming Sun
SenseTime & ZJU-3DV



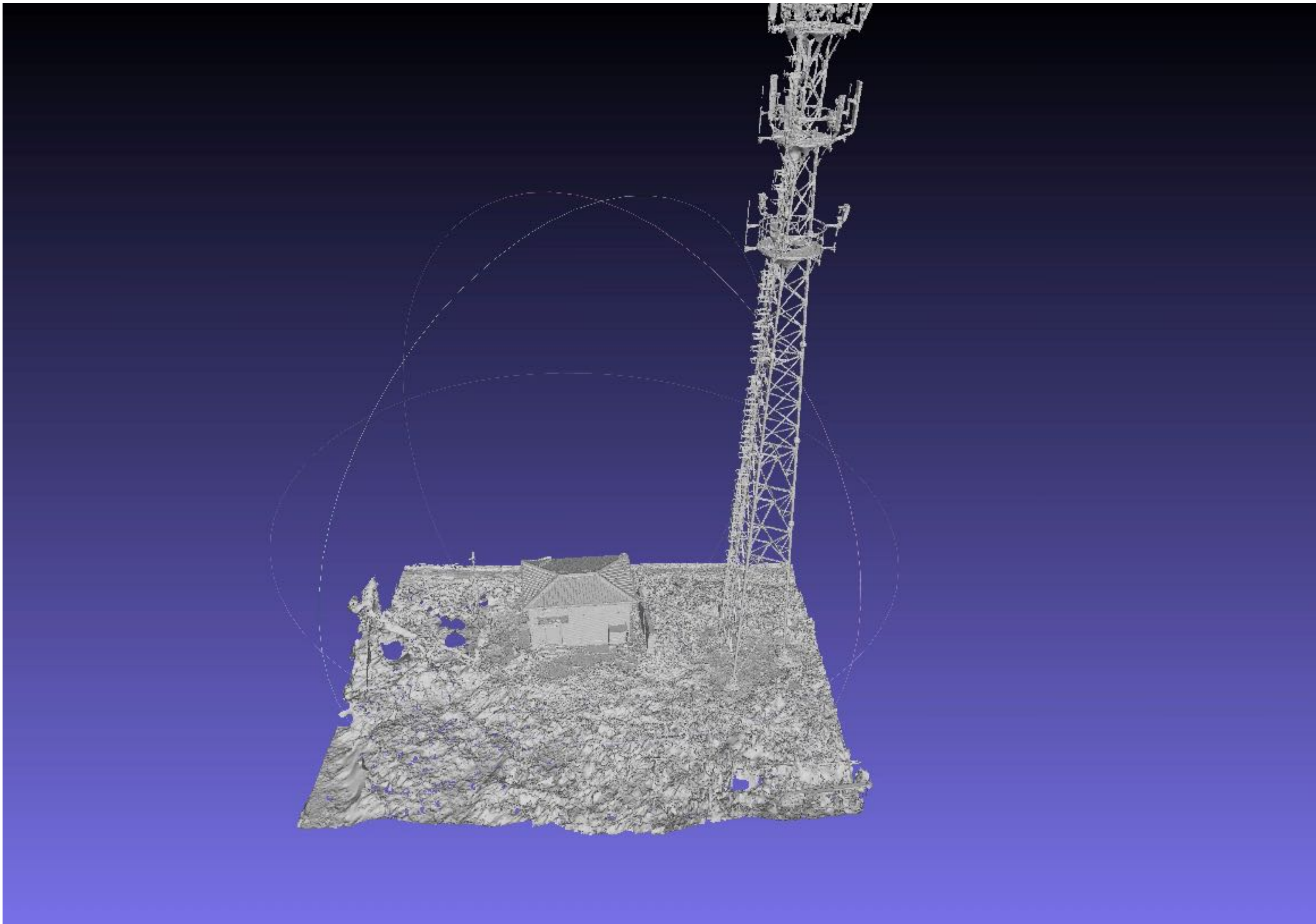
Applications at SenseTime

Localization



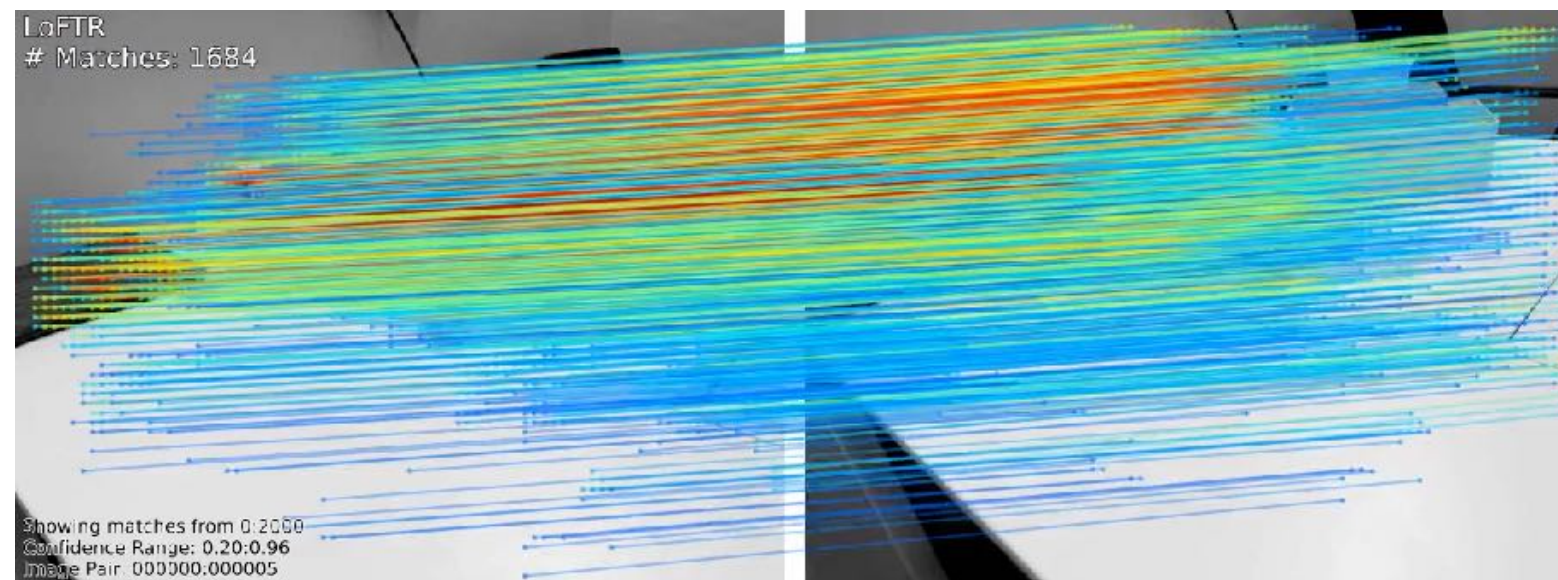
Applications at SenseTime

3D reconstruction



Today's Topic

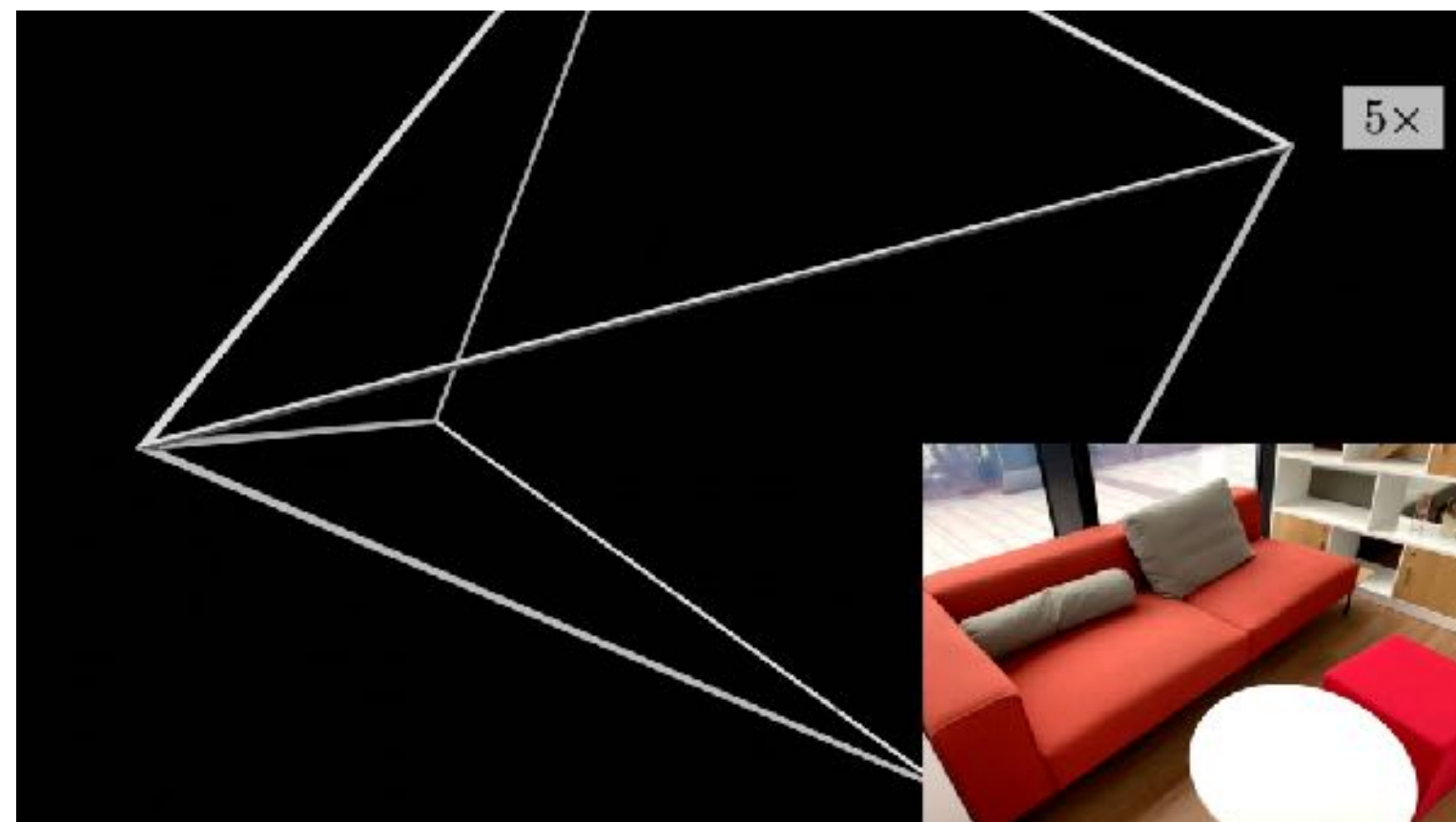
Visual localization



LoFTR: **Detector-Free** Local Feature Matching with **Transformers**

Jiaming Sun* Zehong Shen* Yu'ang Wang* Hujun Bao Xiaowei Zhou
CVPR 2021

3D Reconstruction

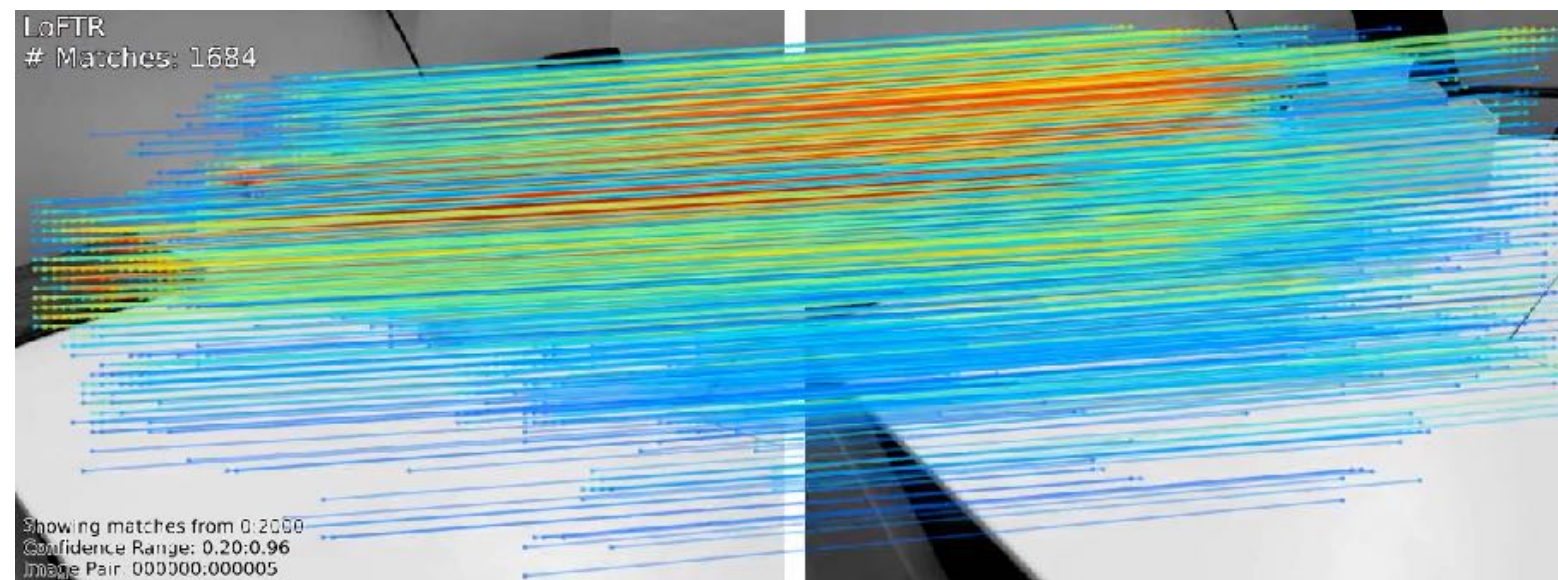


NeuralRecon: **Real-Time Coherent** 3D Reconstruction from Monocular Video

Jiaming Sun* Yiming Xie* Linghao Chen Xiaowei Zhou Hujun Bao
CVPR 2021 (Oral)

Today's Topic

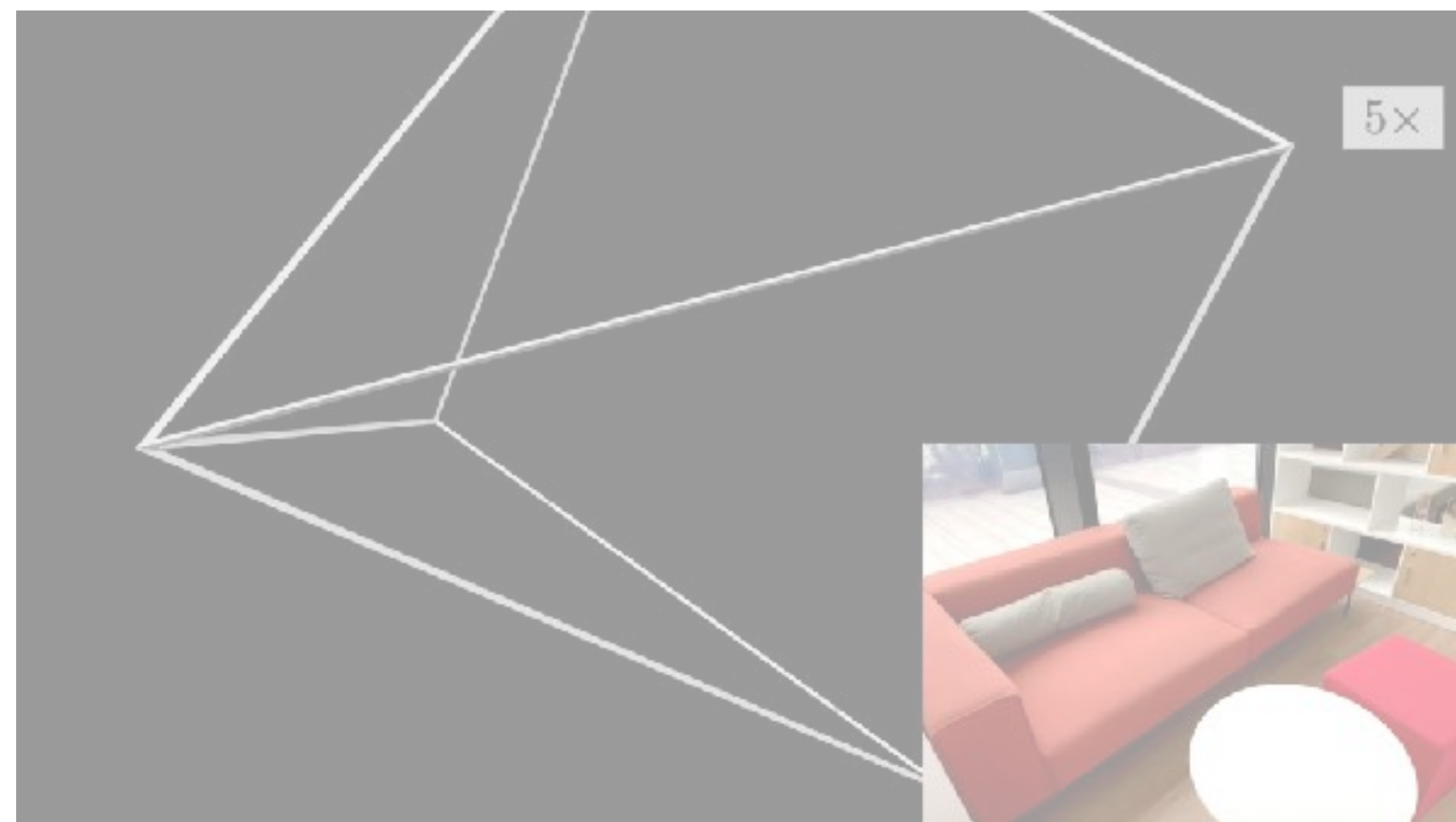
Visual localization



LoFTR: **Detector-Free** Local Feature Matching with **Transformers**

Jiaming Sun* Zehong Shen* Yu'ang Wang* Hujun Bao Xiaowei Zhou
CVPR 2021

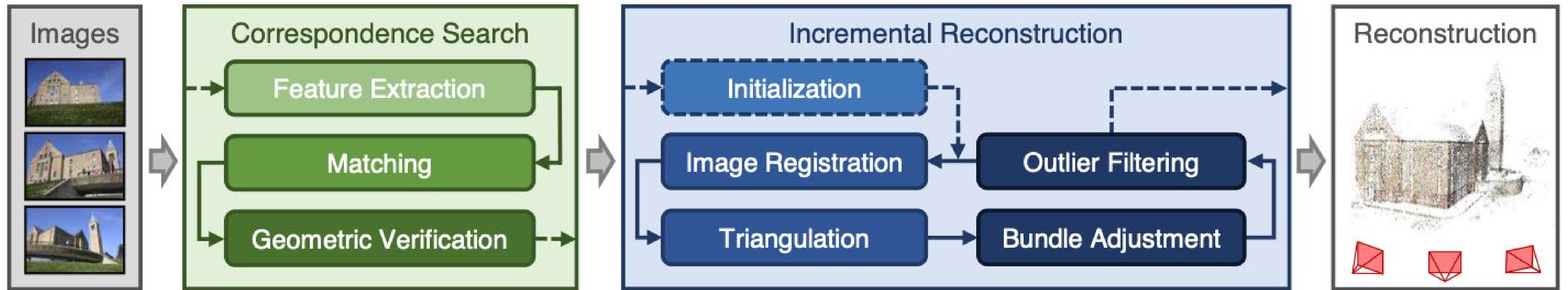
3D Reconstruction



NeuralRecon: **Real-Time Coherent** 3D Reconstruction from Monocular Video

Jiaming Sun* Yiming Xie* Linghao Chen Xiaowei Zhou Hujun Bao
CVPR 2021 (Oral)

Structure-from-Motion (SfM)



Direct Camera Pose Regression

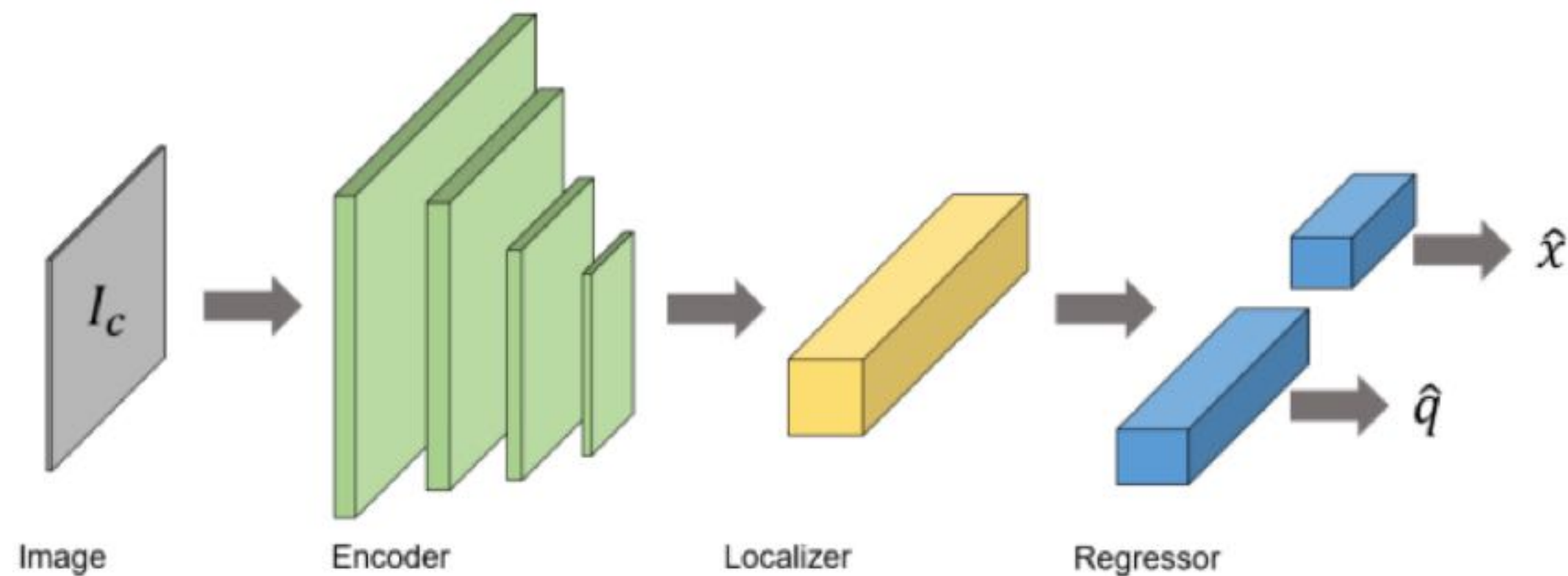
Absolute pose estimation

PoseNet: A Convolutional Network for Real-Time 6-DOF Camera Relocalization

Alex Kendall

Matthew Grimes
University of Cambridge

Roberto Cipolla



•
•
•

Relative pose estimation

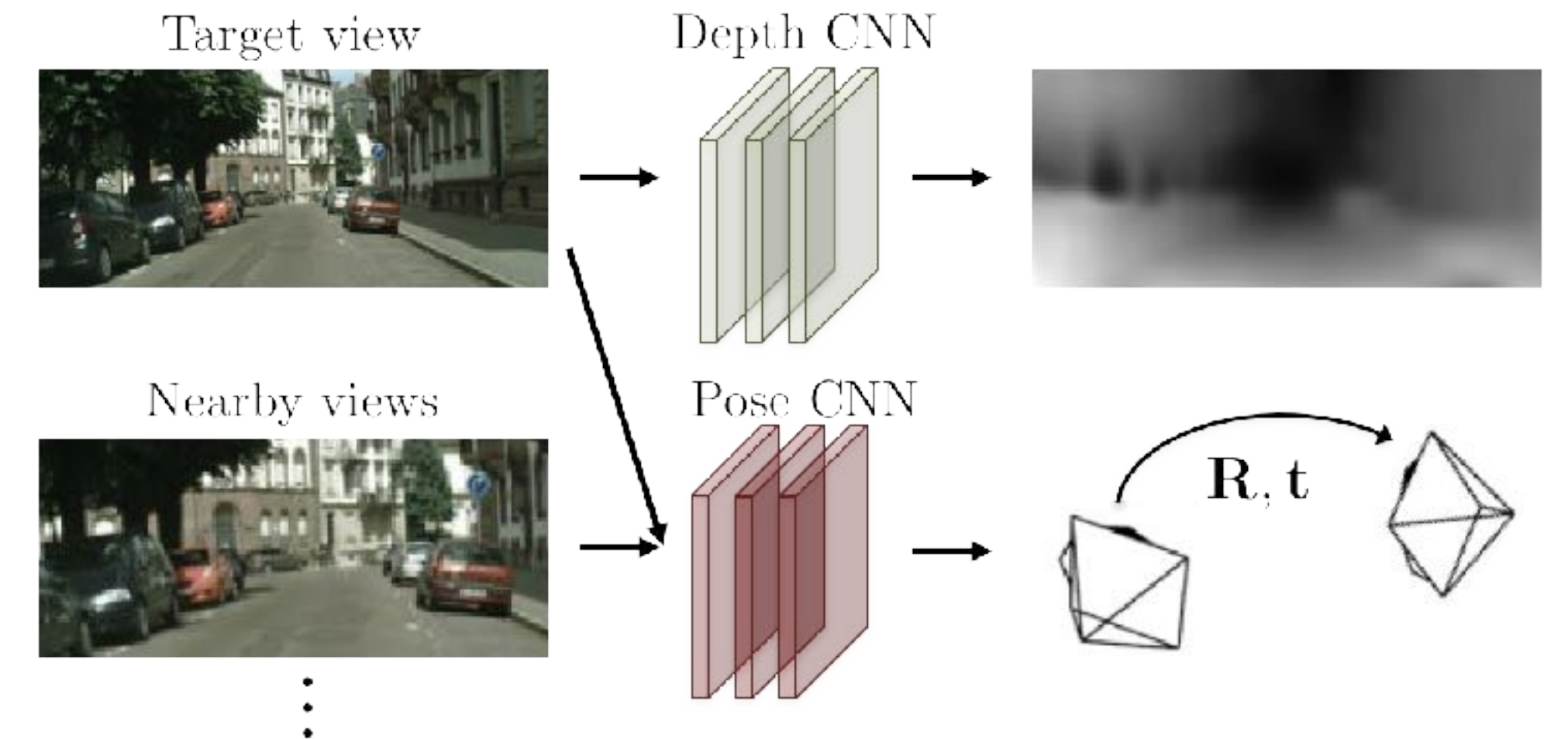
Unsupervised Learning of Depth and Ego-Motion from Video

Tinghui Zhou*
UC Berkeley

Matthew Brown
Google

Noah Snavely
Google

David G. Lowe
Google



(b) Testing: single-view depth and multi-view pose estimation.

•
•
•

Direct Camera Pose Regression

Understanding the Limitations of CNN-based Absolute Camera Pose Regression

Torsten Sattler¹

Qunjie Zhou²

Marc Pollefeys^{3,4}

Laura Leal-Taixé²

¹Chalmers University of Technology

²TU Munich

³ETH Zürich

⁴Microsoft

of APR techniques. Based on this model, we show that APR approaches are **more closely related to approximate pose estimation via image retrieval** (Sec. 5) than to accurate pose estimation via 3D geometry (Sec. 4). **ii) Using**

Depth CNN



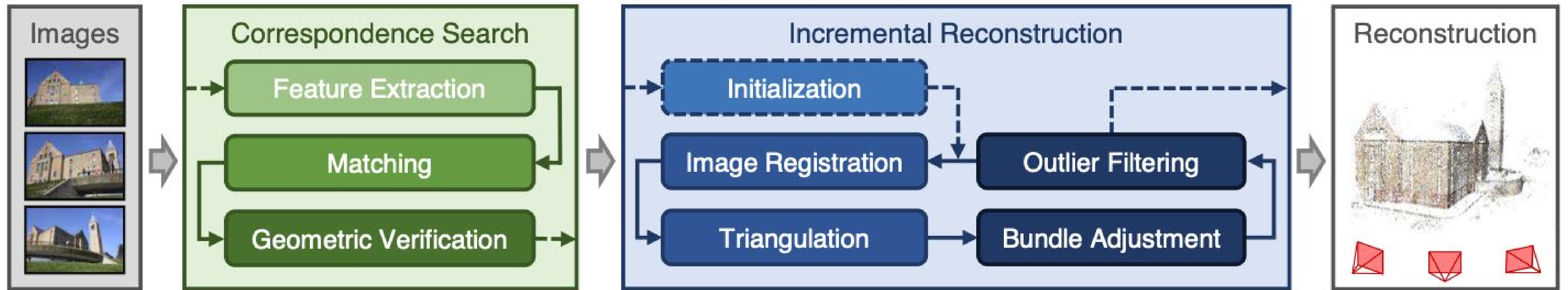
To Learn or Not to Learn: Visual Localization from Essential Matrices

Qunjie Zhou¹, Torsten Sattler², Marc Pollefeys^{3,4}, Laura Leal-Taixé¹

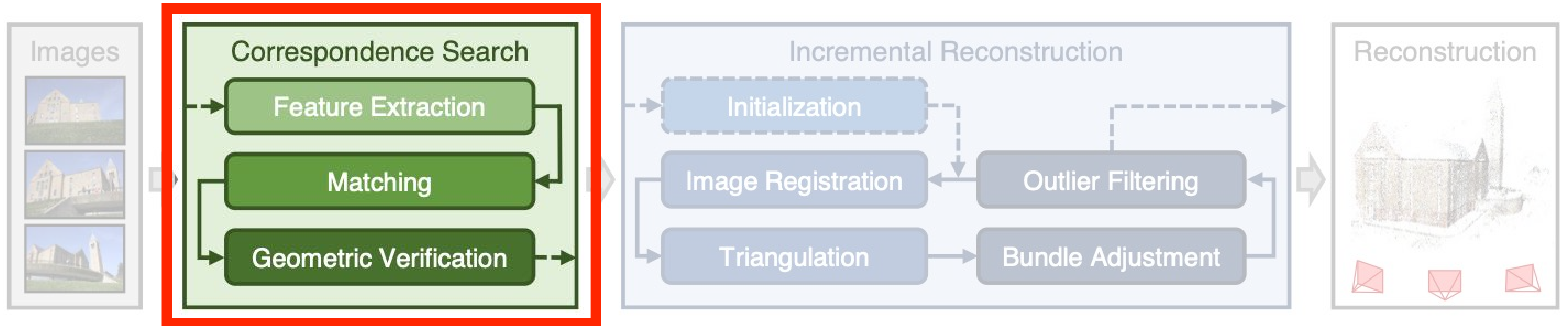
ing their results, we have shown that the **purely data-driven approach does not generalize well** and have identified the reason for this failure as the relative pose regression layers.

-
-
-

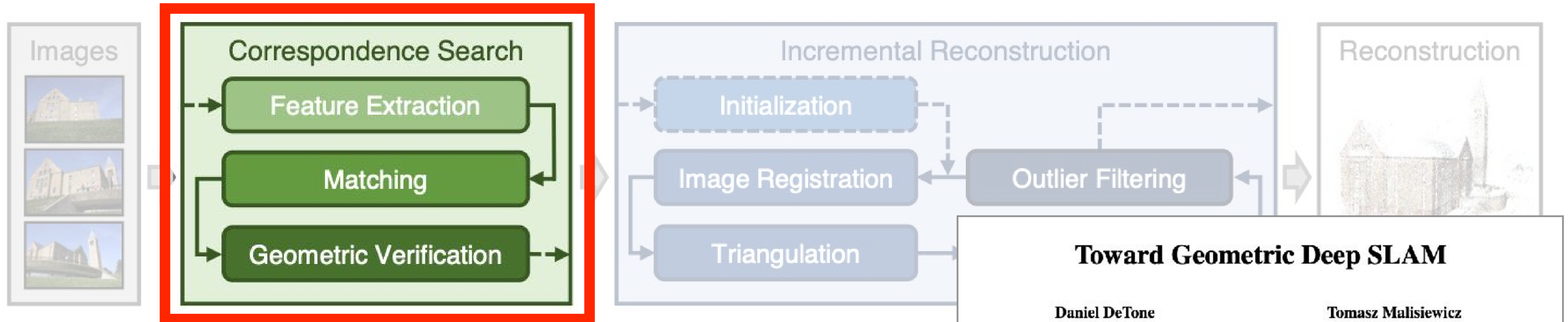
Revisiting Structure-from-Motion (SfM)



Learning to Find Correspondences



Learning to Find Correspondences



MatchNet: Unifying Feature and Metric Learning for Patch-Based Matching

Xufeng Han[†] Thomas Leung[‡] Yangqing Jia[‡] Rahul Sukthankar[‡] Alexander C. Berg[†]
[†]University of North Carolina at Chapel Hill [‡]Google Research

Toward Geometric Deep SLAM

Daniel DeTone
Magic Leap, Inc.
Sunnyvale, CA
ddetone@magicleap.com

Tomasz Malisiewicz
Magic Leap, Inc.
Sunnyvale, CA
tmalisiewicz@magicleap.com

Andrew Rabinovich
Magic Leap, Inc.
Sunnyvale, CA
arabinovich@magicleap.com

LIFT: Learned Invariant Feature Transform

Kwang Moo Yi^{*,1}, Eduard Trulls^{*,1}, Vincent Lepetit², Pascal Fua¹

SuperPoint: Self-Supervised Interest Point Detection and Description

Daniel DeTone
Magic Leap
Sunnyvale, CA
ddetone@magicleap.com

Tomasz Malisiewicz
Magic Leap
Sunnyvale, CA
tmalisiewicz@magicleap.com

Andrew Rabinovich
Magic Leap
Sunnyvale, CA
arabinovich@magicleap.com

Recently: D2-Net, R2D2, GIFT, DISK, SuperGlue...

Local Feature Applications @ SenseTime

Benchmark测试

- 自研特征检测提取模块SVCNN相比superpoint定位成功率提升~10个点
- 整体耗时10ms

B5地下停车场					
	# query	成功率	<0.5m	<1.0m	<2.0m
V1 (SIFT)	13284	71.3339%	65.53%	68.4131%	69.9262%
V2 (SuperPoint)	13284	65.0407%	58.6646%	61.8714%	63.5125%
V3 (SVCNN)	13284	74.39%	66.2376%	70.1144%	72.3050%

国博F1					
	# query	成功率	<0.5m	<1.0m	<2.0m
V1 (SIFT)	13212	30.5101%	16.3185%	20.7236%	24.3264%
V2 (SuperPoint)	13212	40.0318%	24.7427%	29.7381%	32.8716%
V3 (SVCNN)	13212	51.2413%	29.1023%	36.1792%	41.0157%

Applications



商汤SenseMARS助力成都IFS打造全国首个全场景城市综合体AR导航

此次在成都IFS覆盖区域多达46万平方米，综合定位成功率高达99%，定位精度达“厘米”级别，单次时长达“毫秒”级别，且不会出现偏移、闪烁等情况。

Local Feature Applications @ SenseTime

Benchmark测试

- 自研特征检测提取模块SVCNN相比superpoint定位成功率提升~10个百分点
- 整体耗时10ms

OK, this is all pretty cool and exciting.
What's next?

	# query	成功率	<0.5m	<1.0m	<2.0m
V1 (SIFT)	13284	71.3339%	65.53%	61.43%	69.9722%
V2 (SuperPoint)	13284	65.0407%	58.6646%	61.8714%	63.5125%
V3 (SVCNN)	13284	74.39%	66.2376%	70.1144%	72.3050%

国博F1					
	# query	成功率	<0.5m	<1.0m	<2.0m
V1 (SIFT)	13212	30.5101%	16.3185%	20.7236%	24.3264%
V2 (SuperPoint)	13212	40.0318%	24.7427%	29.7381%	32.8716%
V3 (SVCNN)	13212	51.2413%	29.1023%	36.1792%	41.0157%

Applications

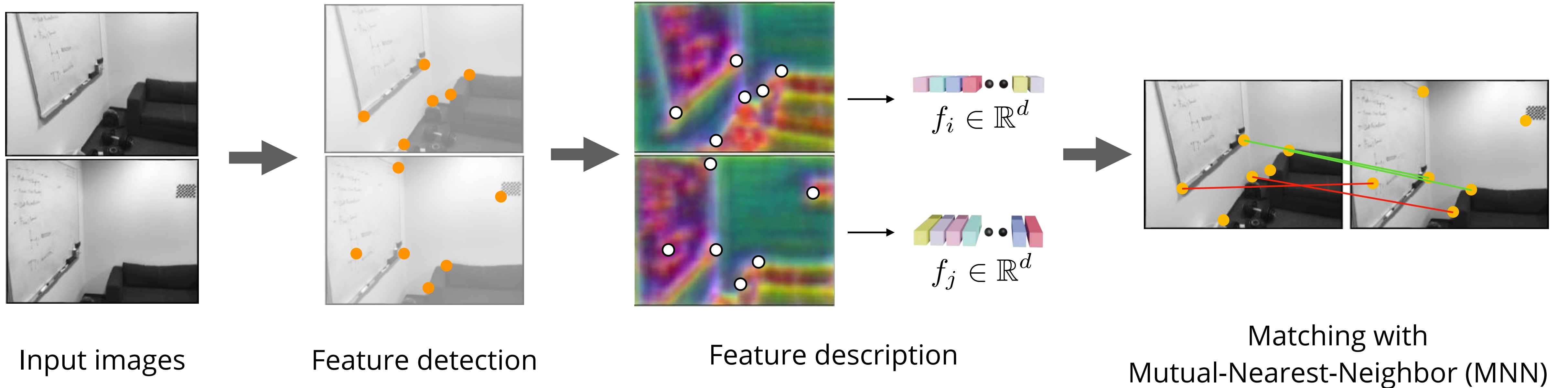


商汤SenseMARS助力成都IFS打造全国首个全场景城市综合体AR导航

此次在成都IFS覆盖区域多达46万平方米，综合定位成功率高达99%，定位精度达“厘米”级别，单次时长达“毫秒”级别，且不会出现偏移、闪烁等情况。

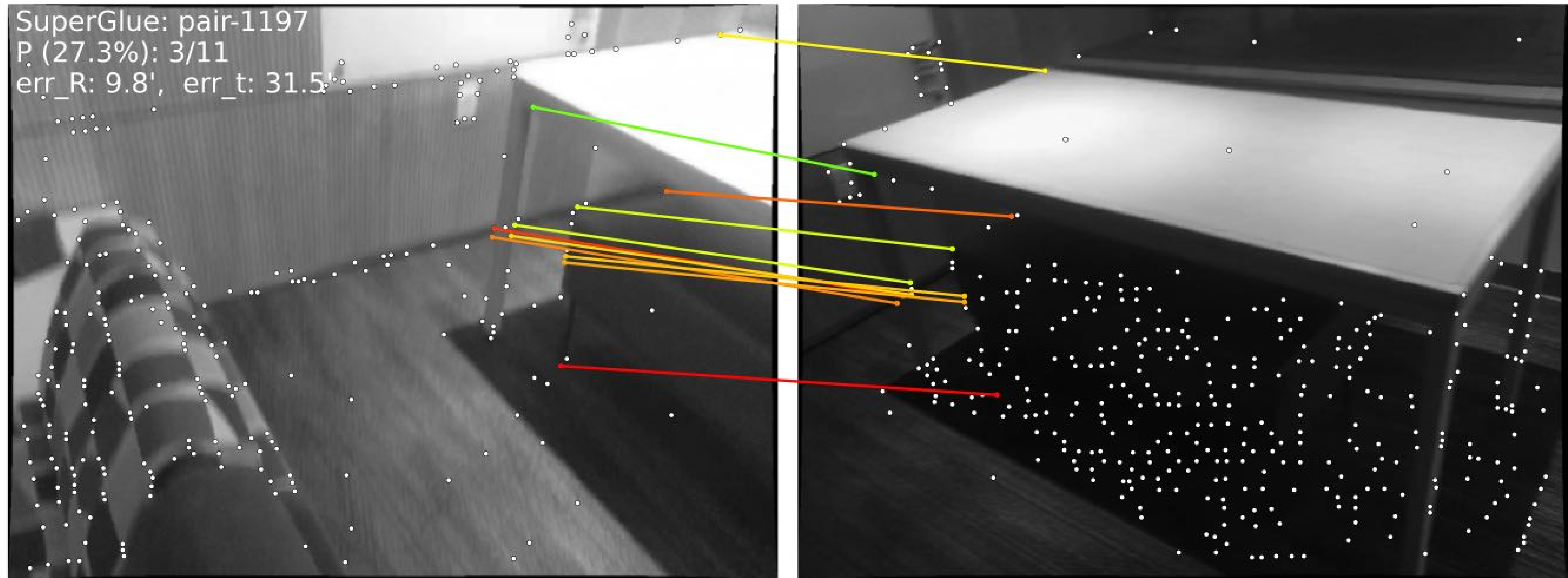
Introduction to LoFTR

The standard pipeline for local feature matching



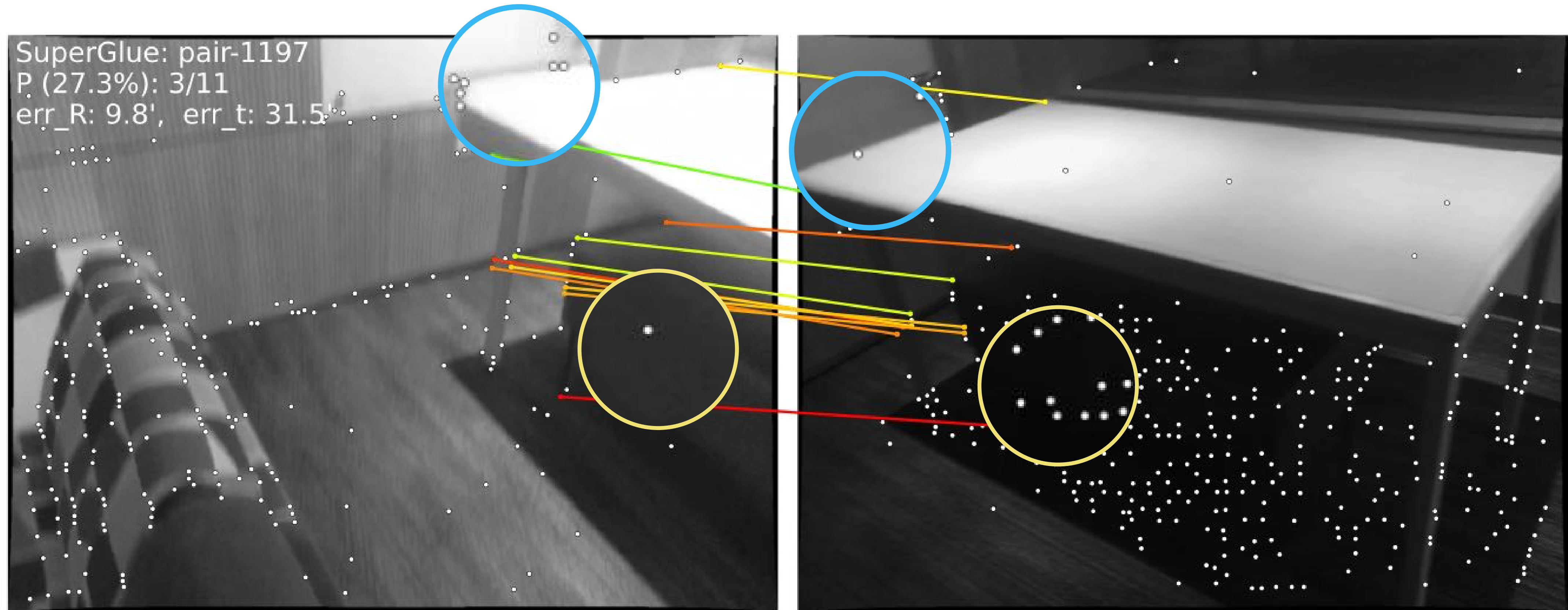
Introduction to LoFTR

Challenge 1: Detector repeatability on two (or more) images



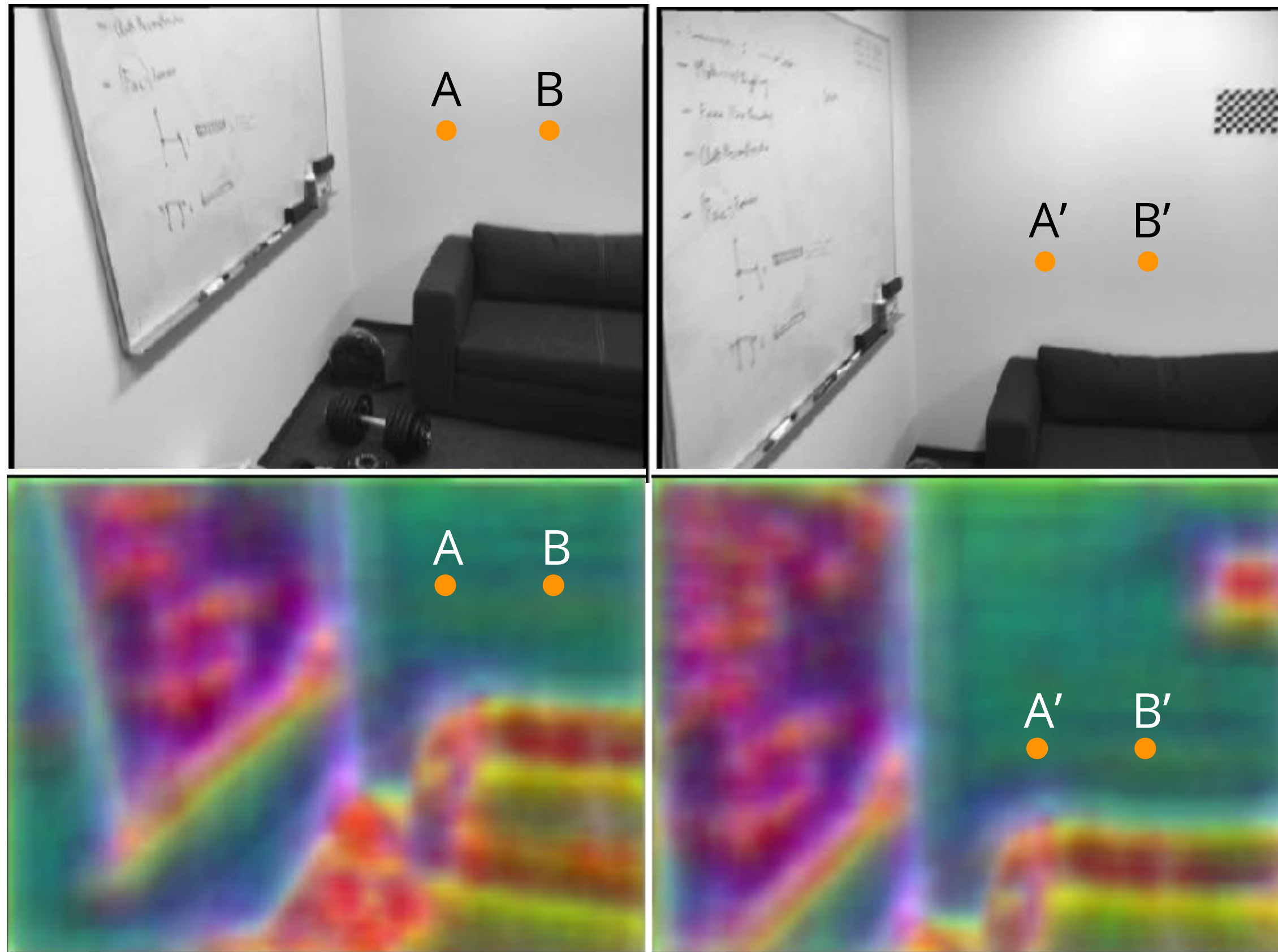
Introduction to LoFTR

Challenge 1: Detector repeatability on two (or more) images

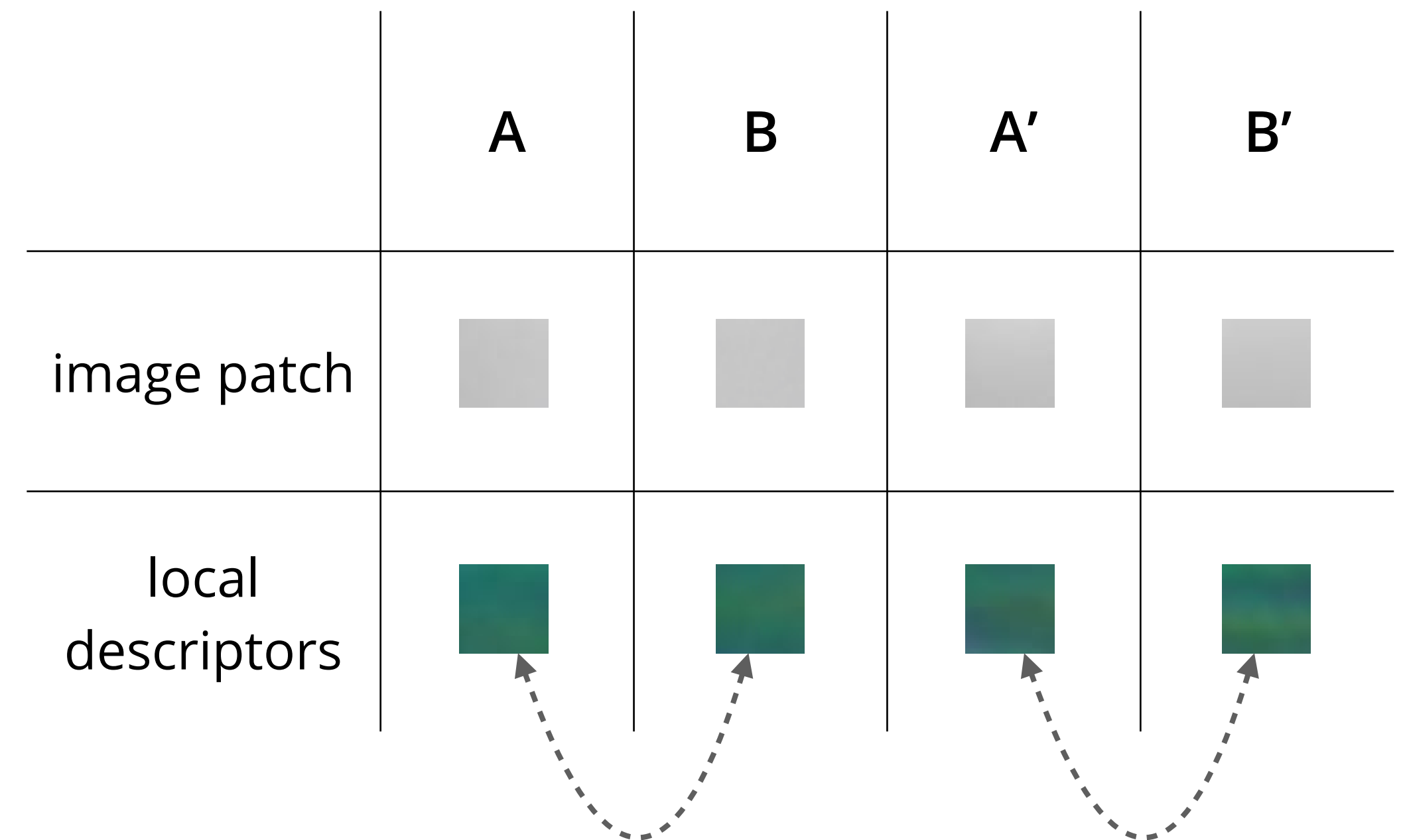


Introduction to LoFTR

Challenge 2: Local descriptors are not position-dependent



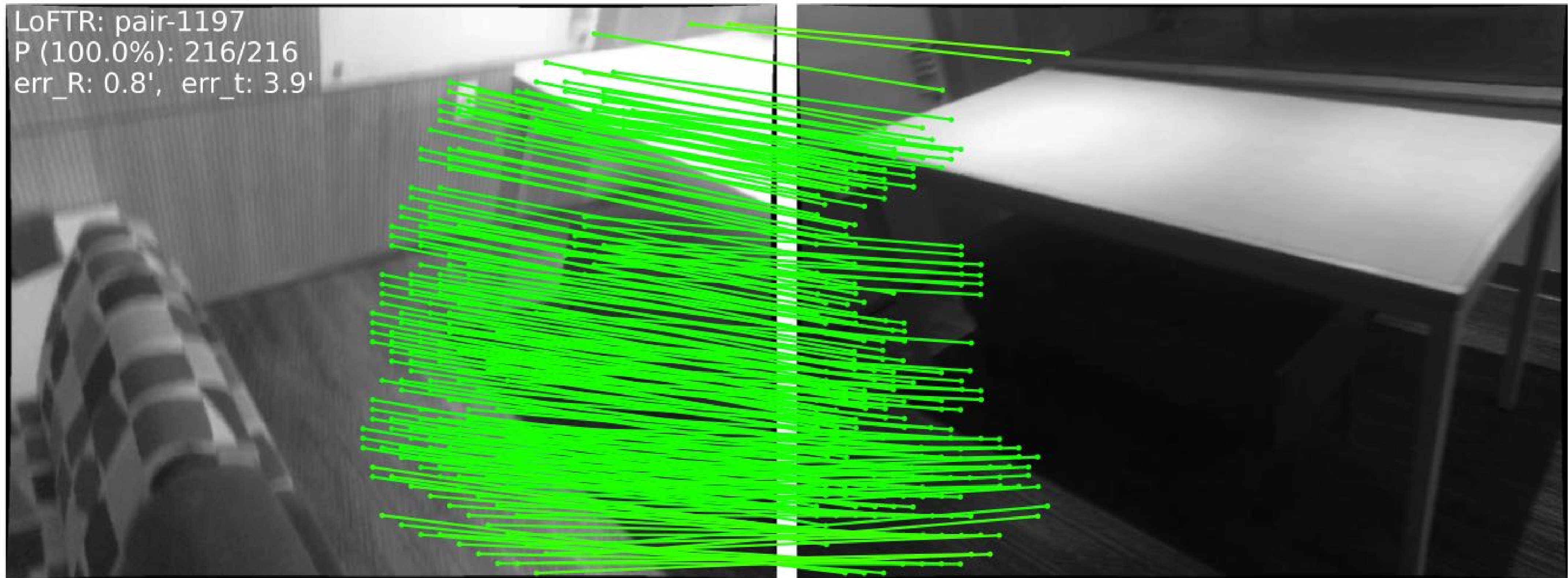
Local features visualized with PCA



Expected to be different
(position-dependent)

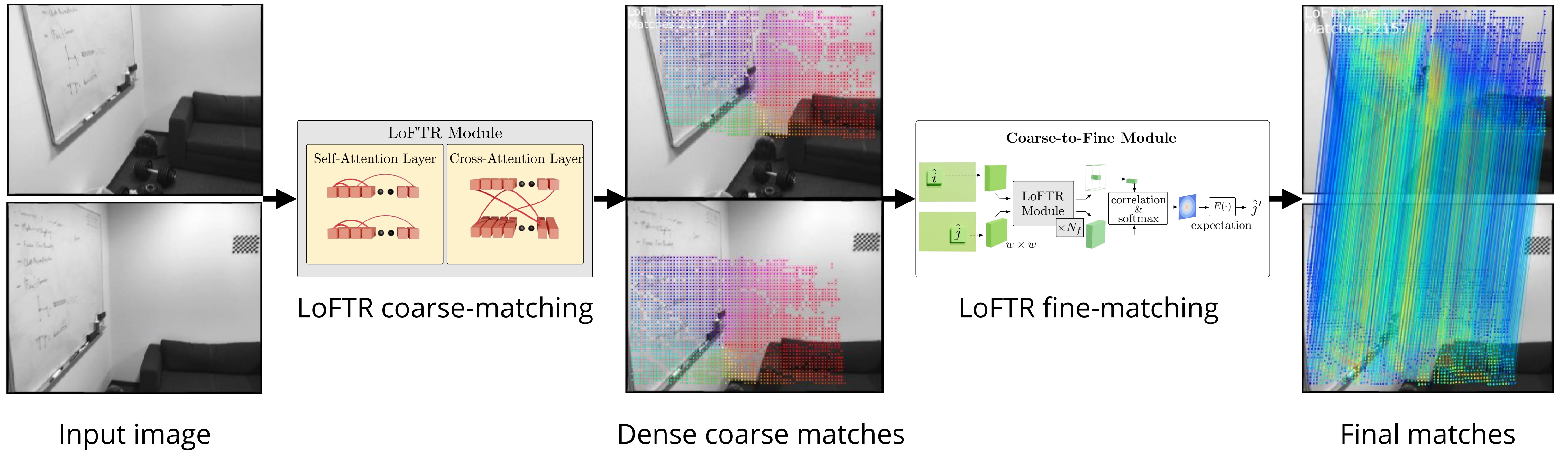
Introduction to LoFTR

Our solution: LoFTR



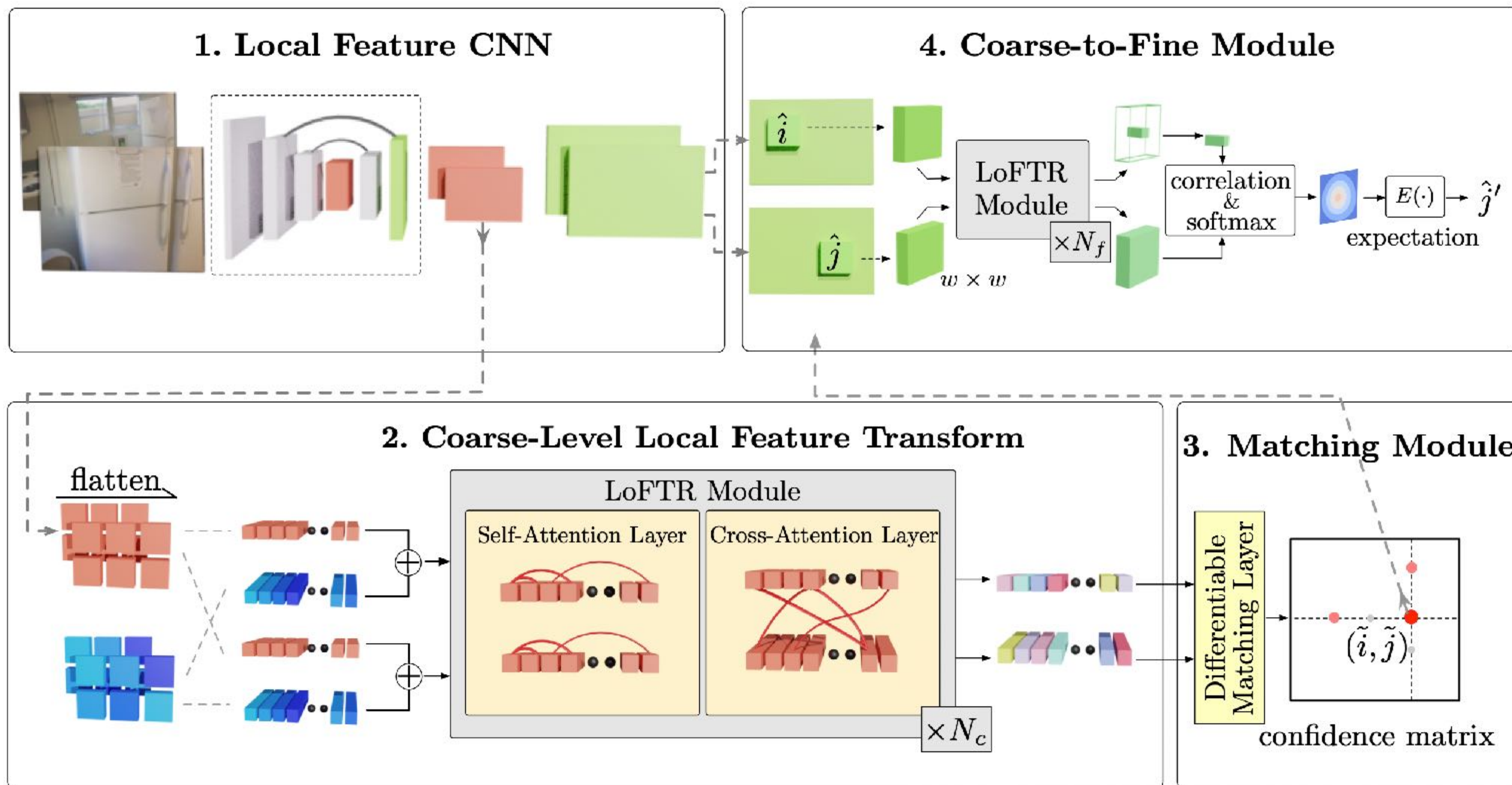
LoFTR

Pipeline overview



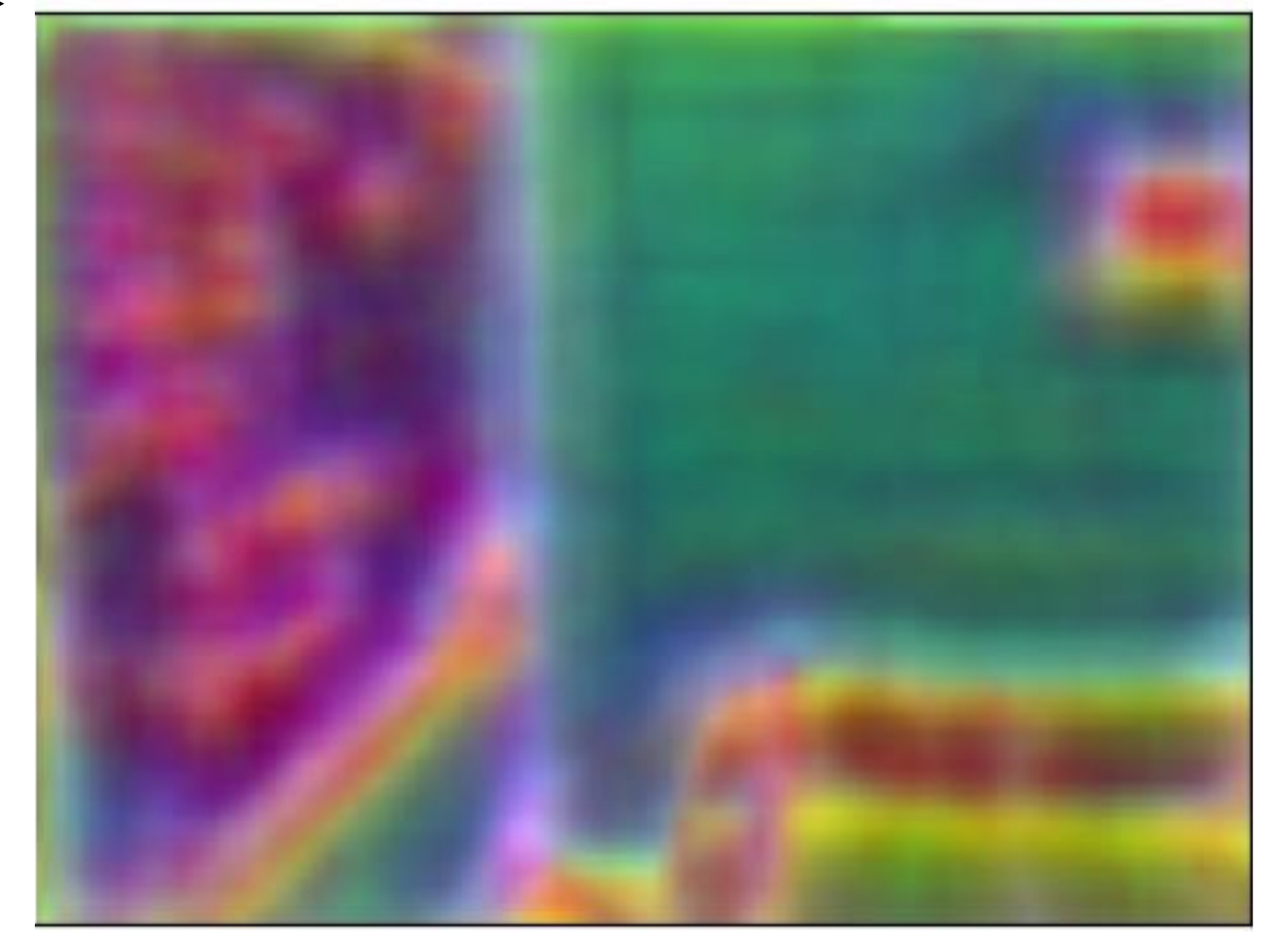
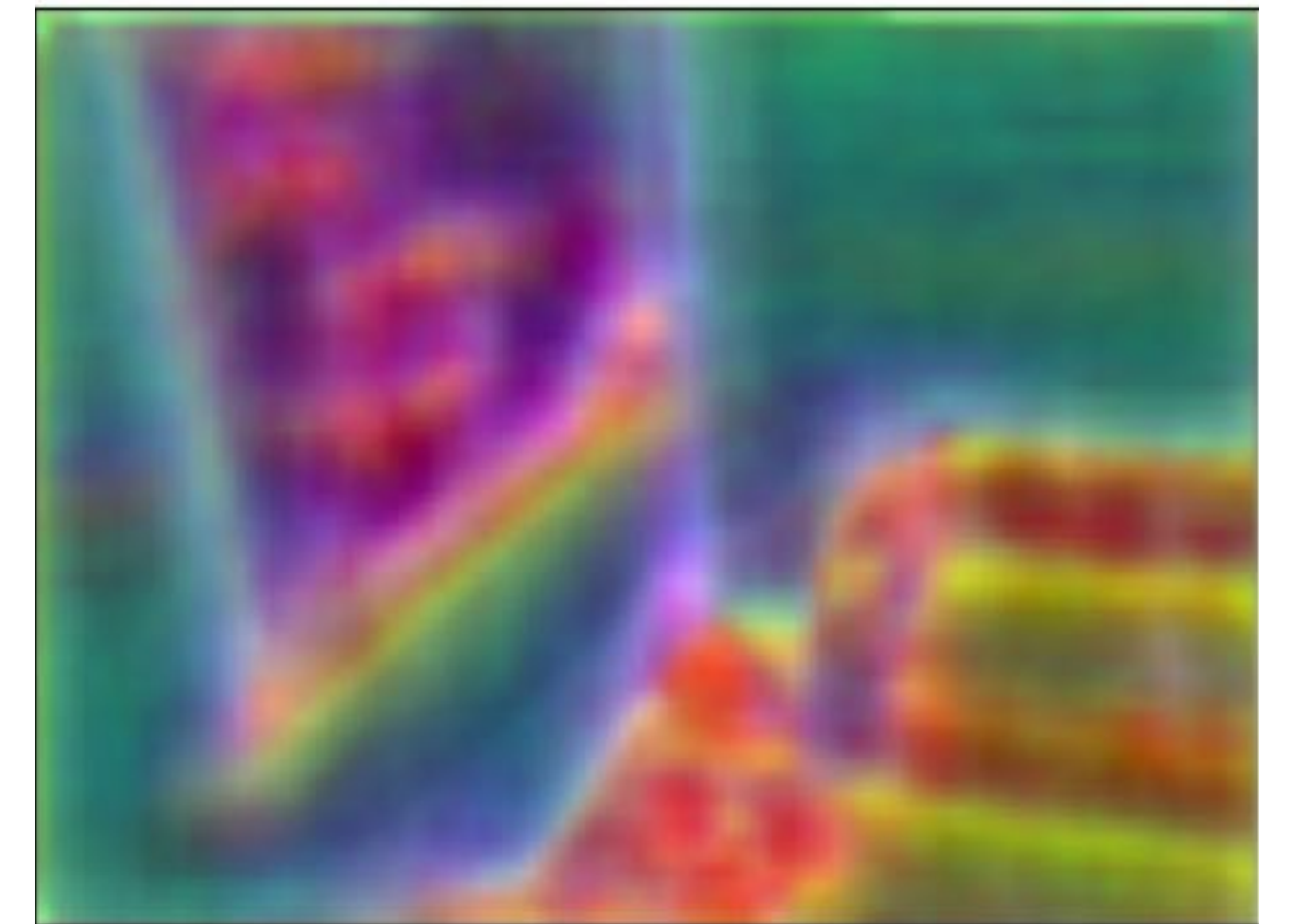
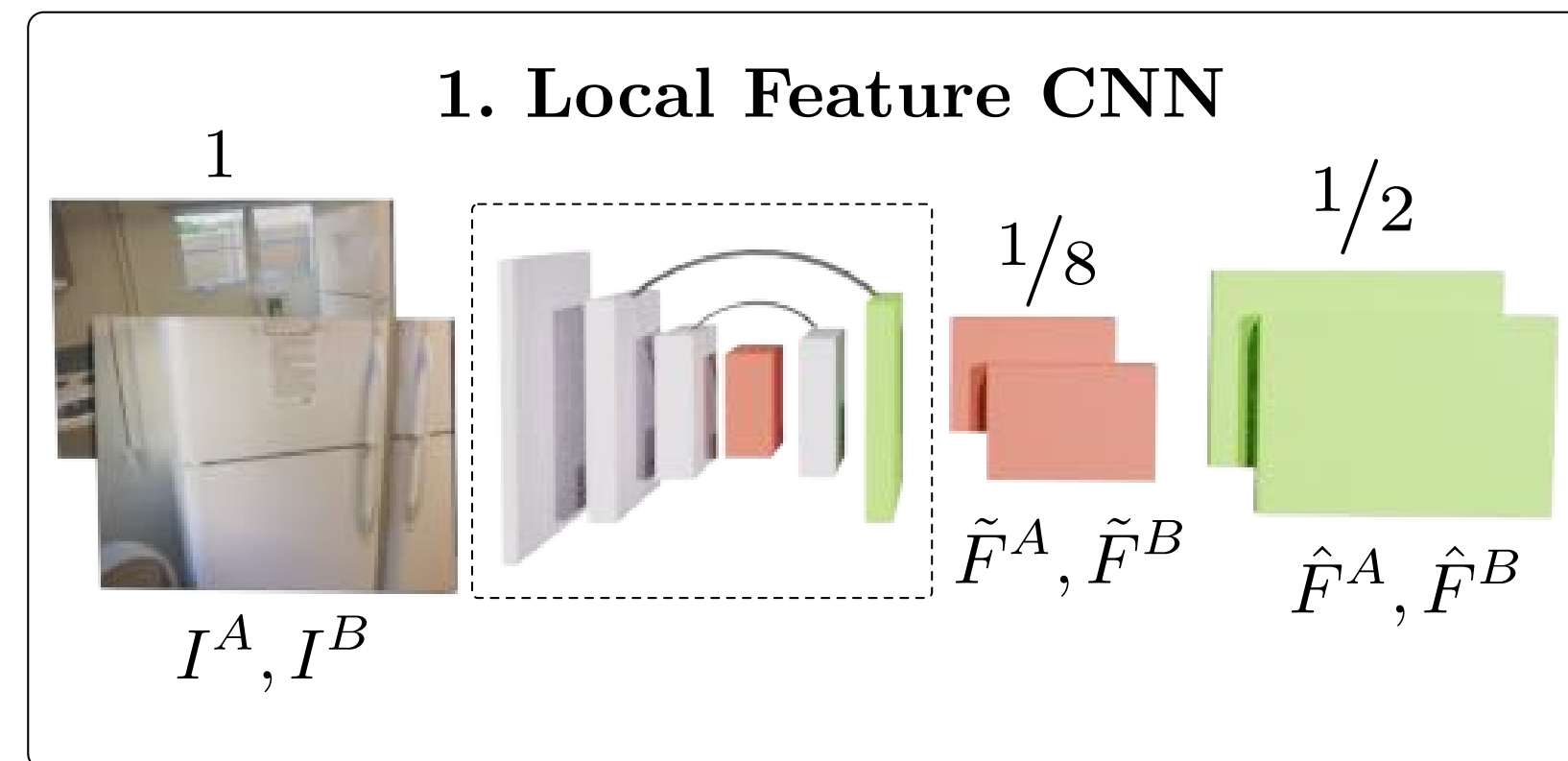
LoFTR

Architecture



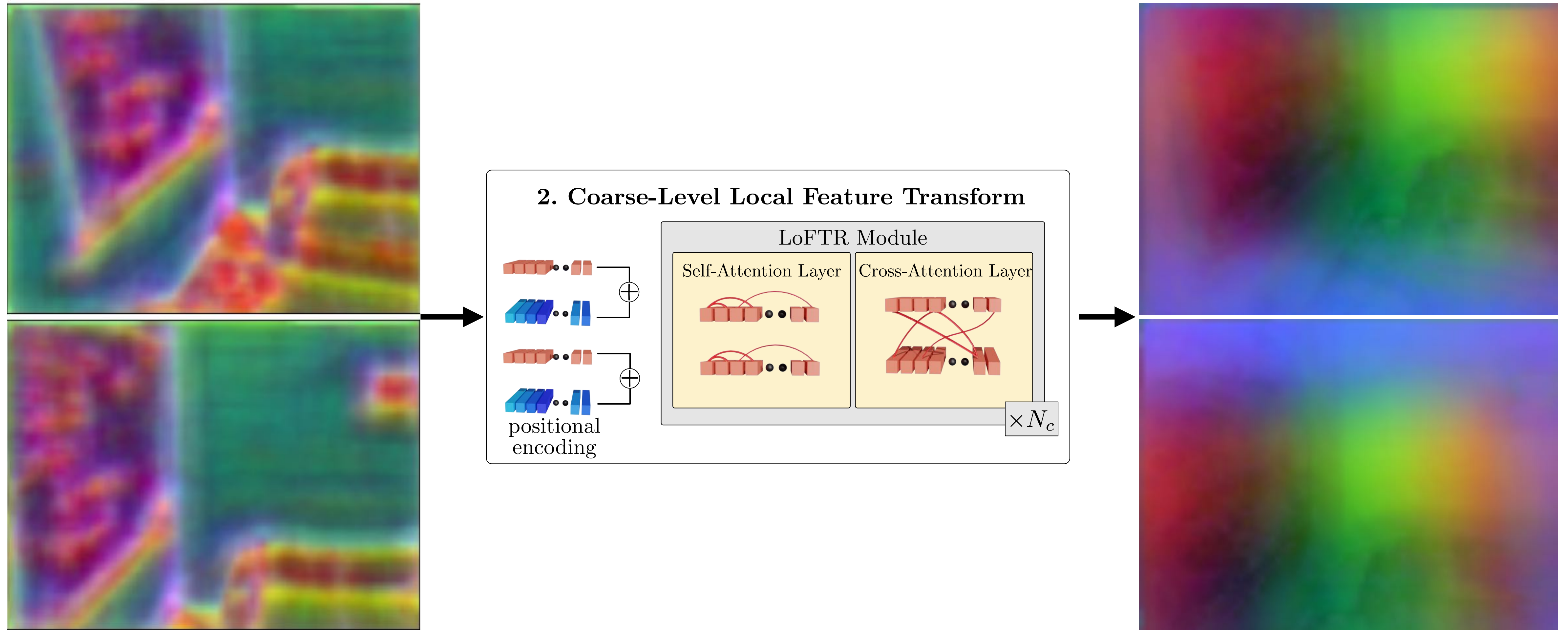
LoFTR

Local feature CNN



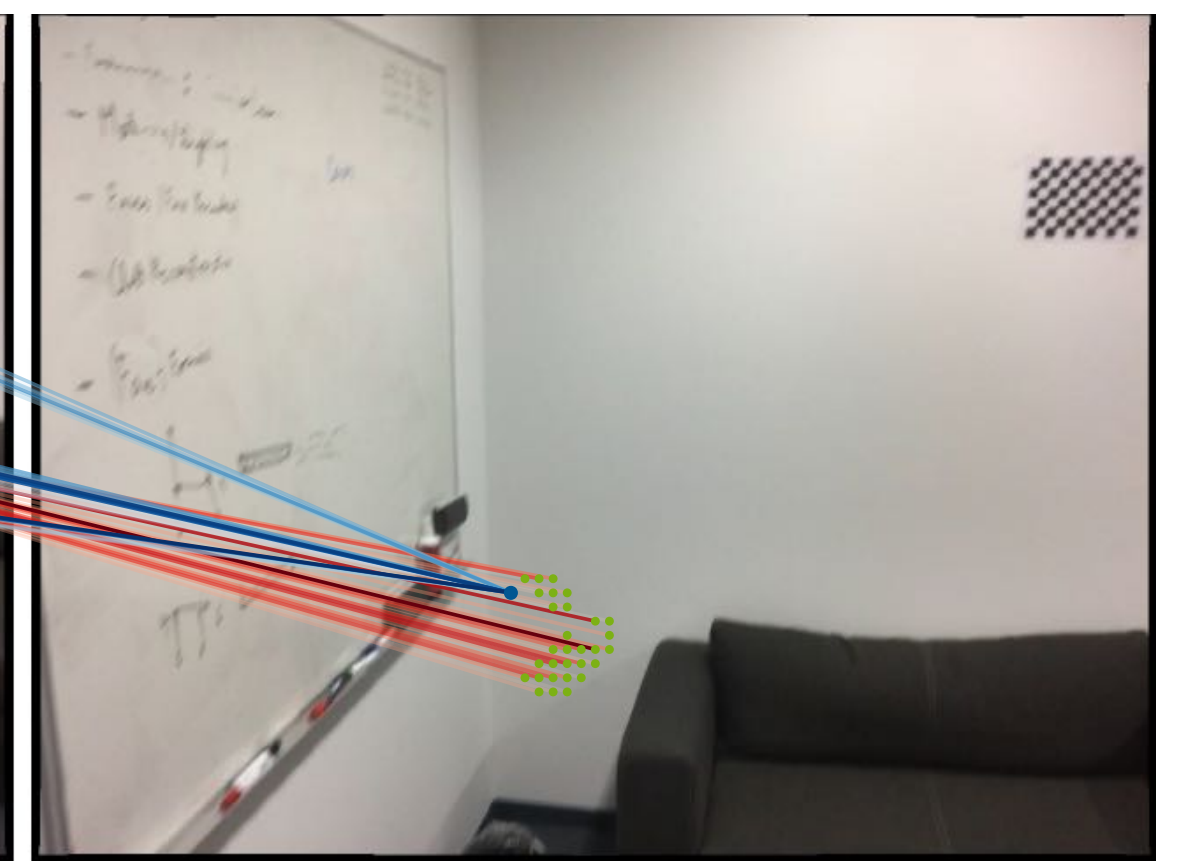
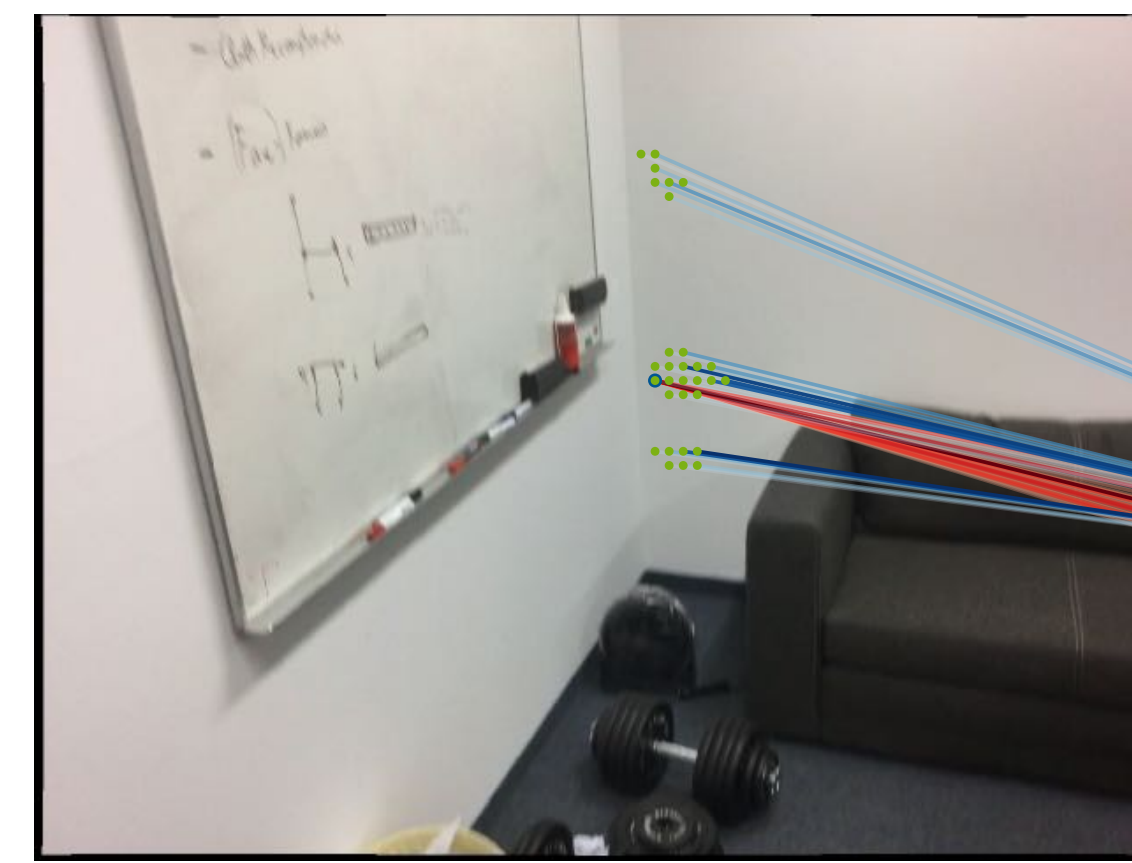
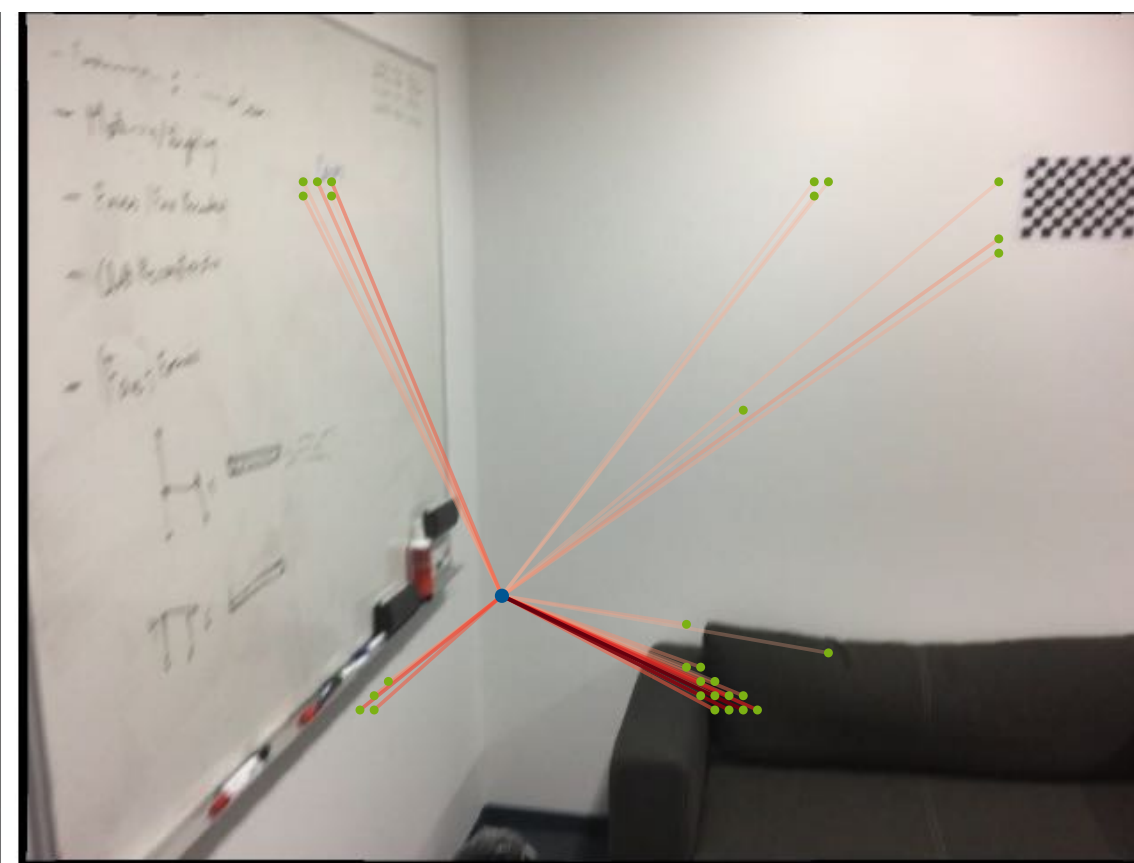
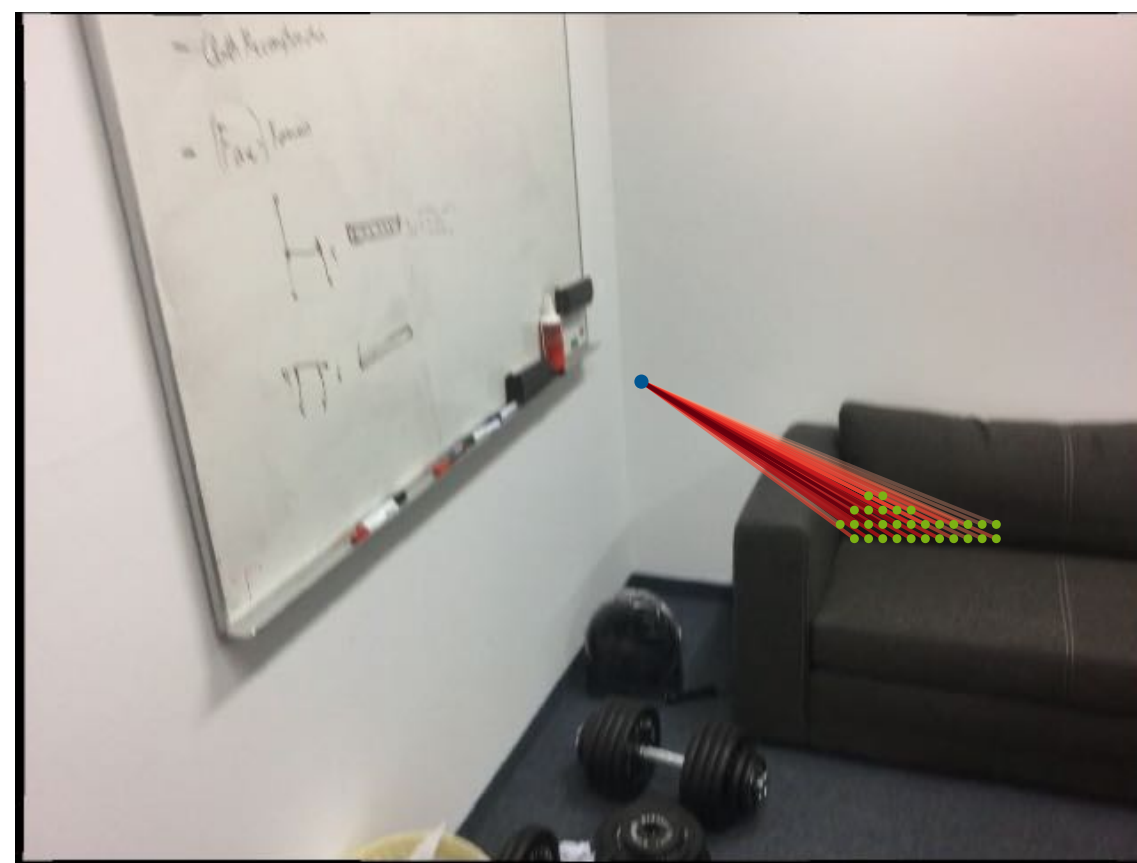
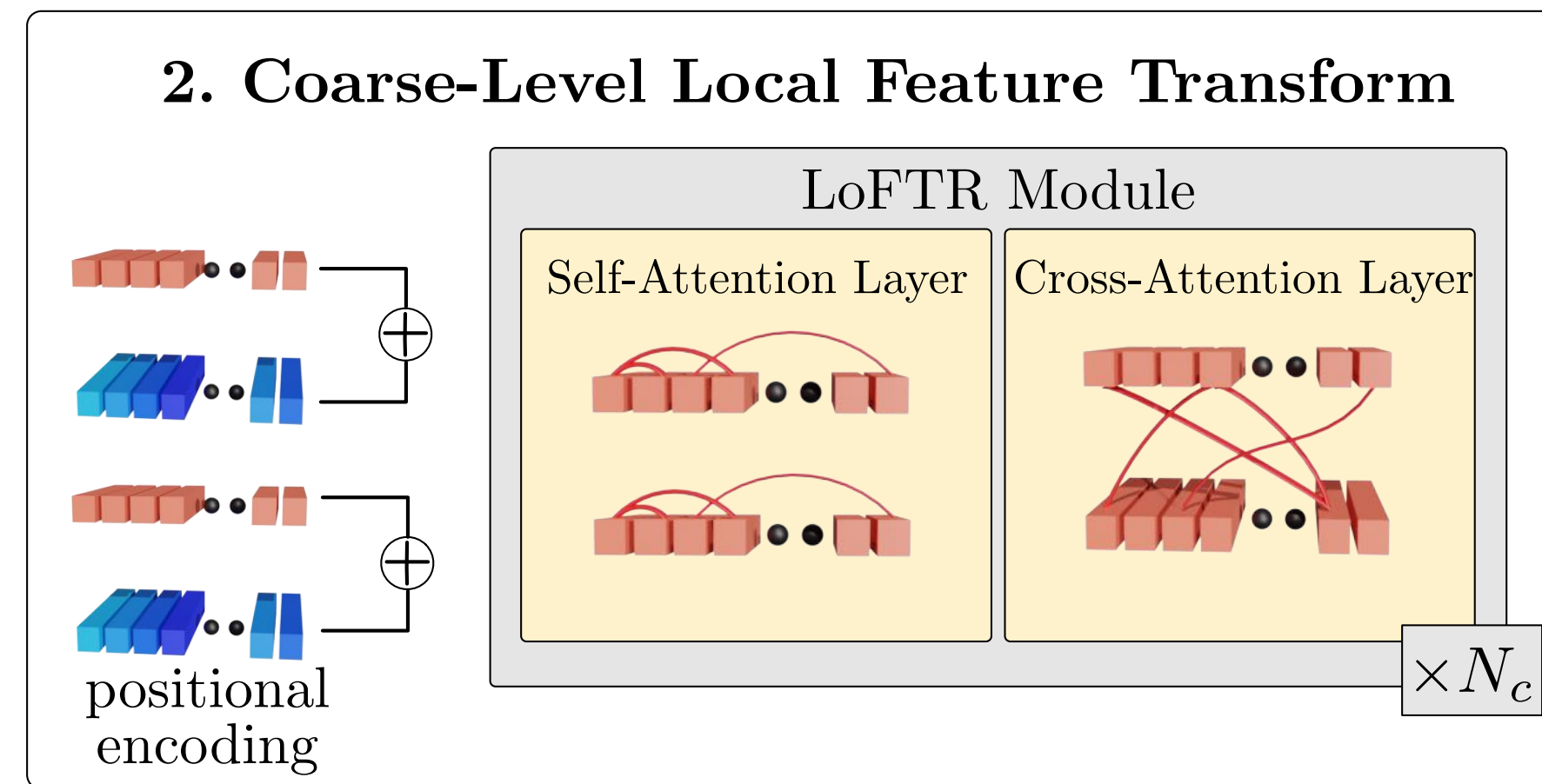
LoFTR

Coarse-level LoFTR module



LoFTR

Coarse-level LoFTR module

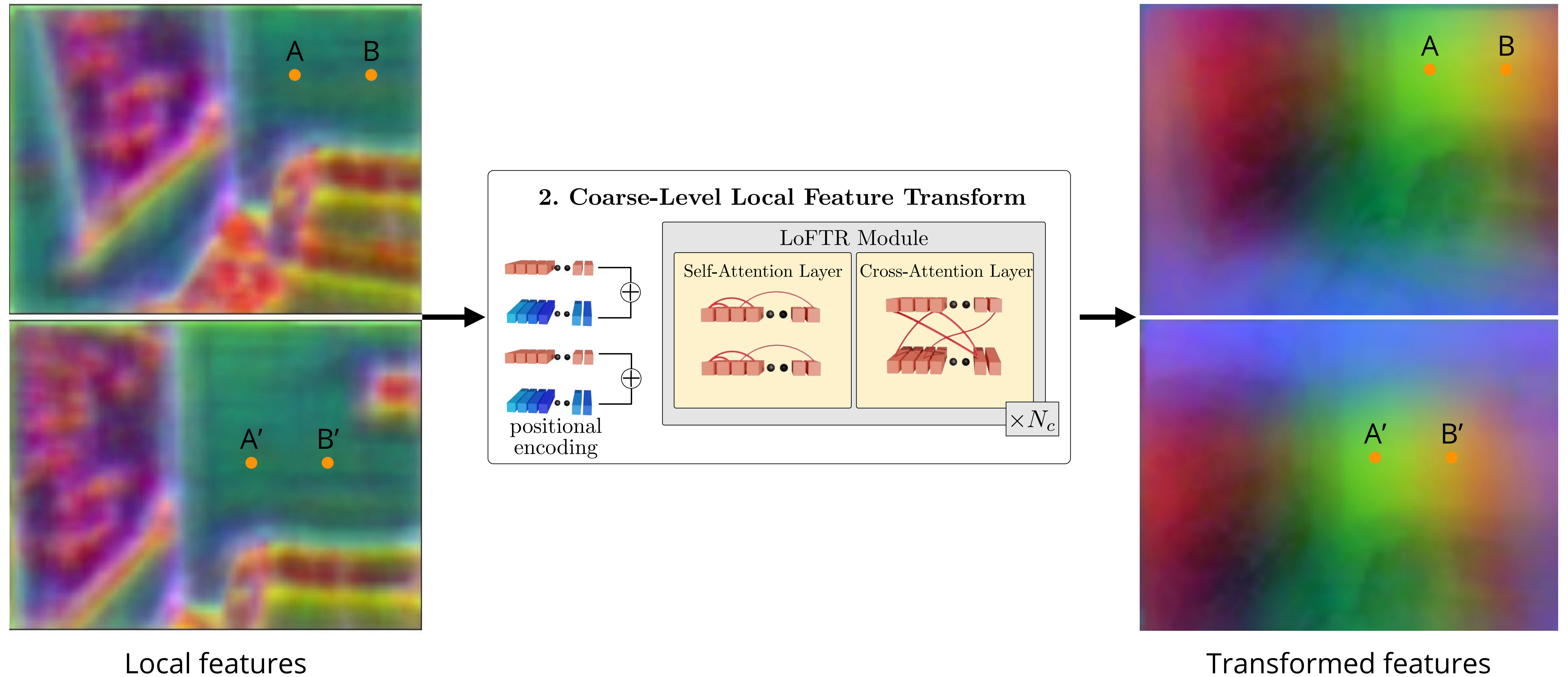


Attention weight visualization of **self-attention**

Attention weight visualization of **cross-attention**

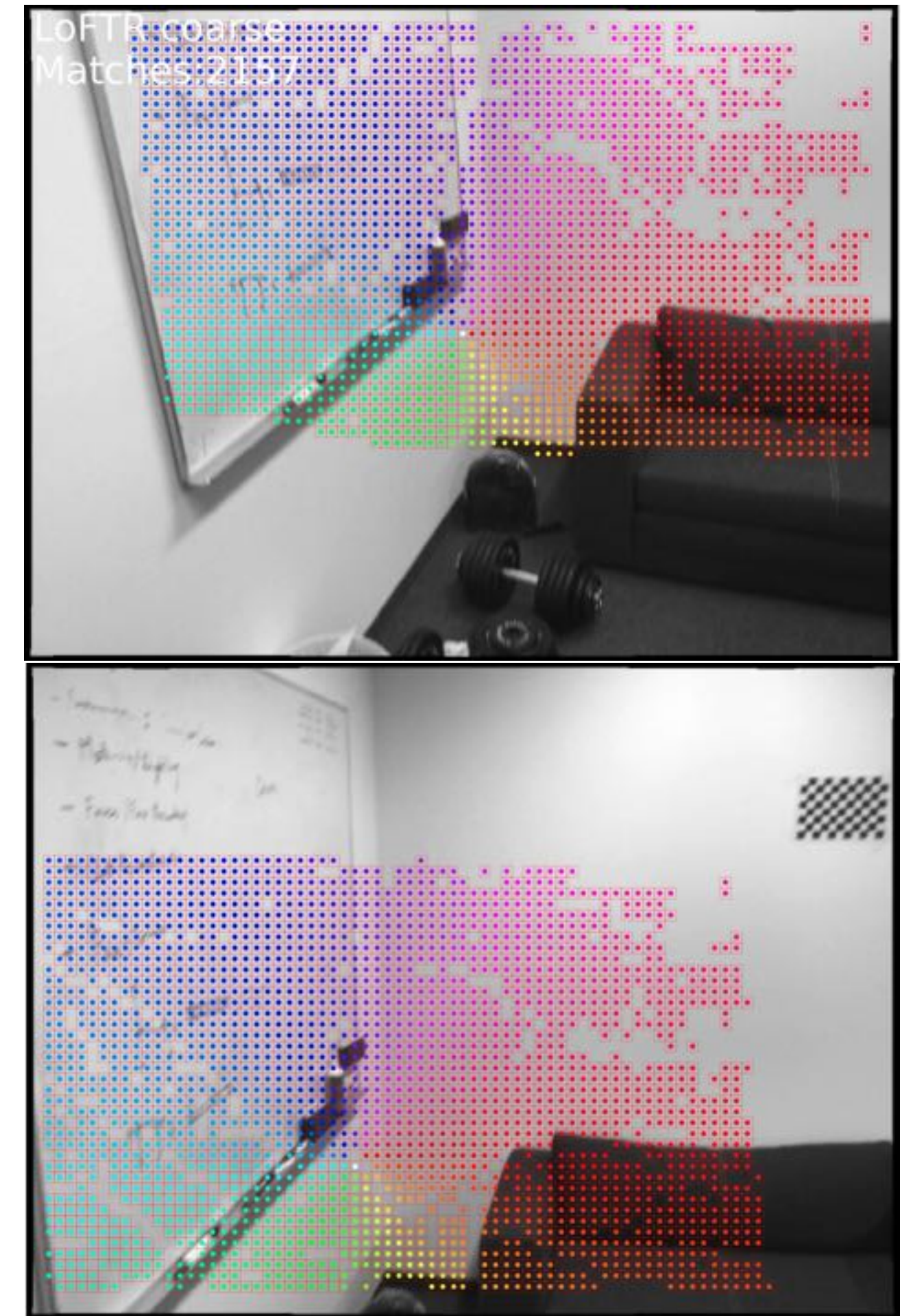
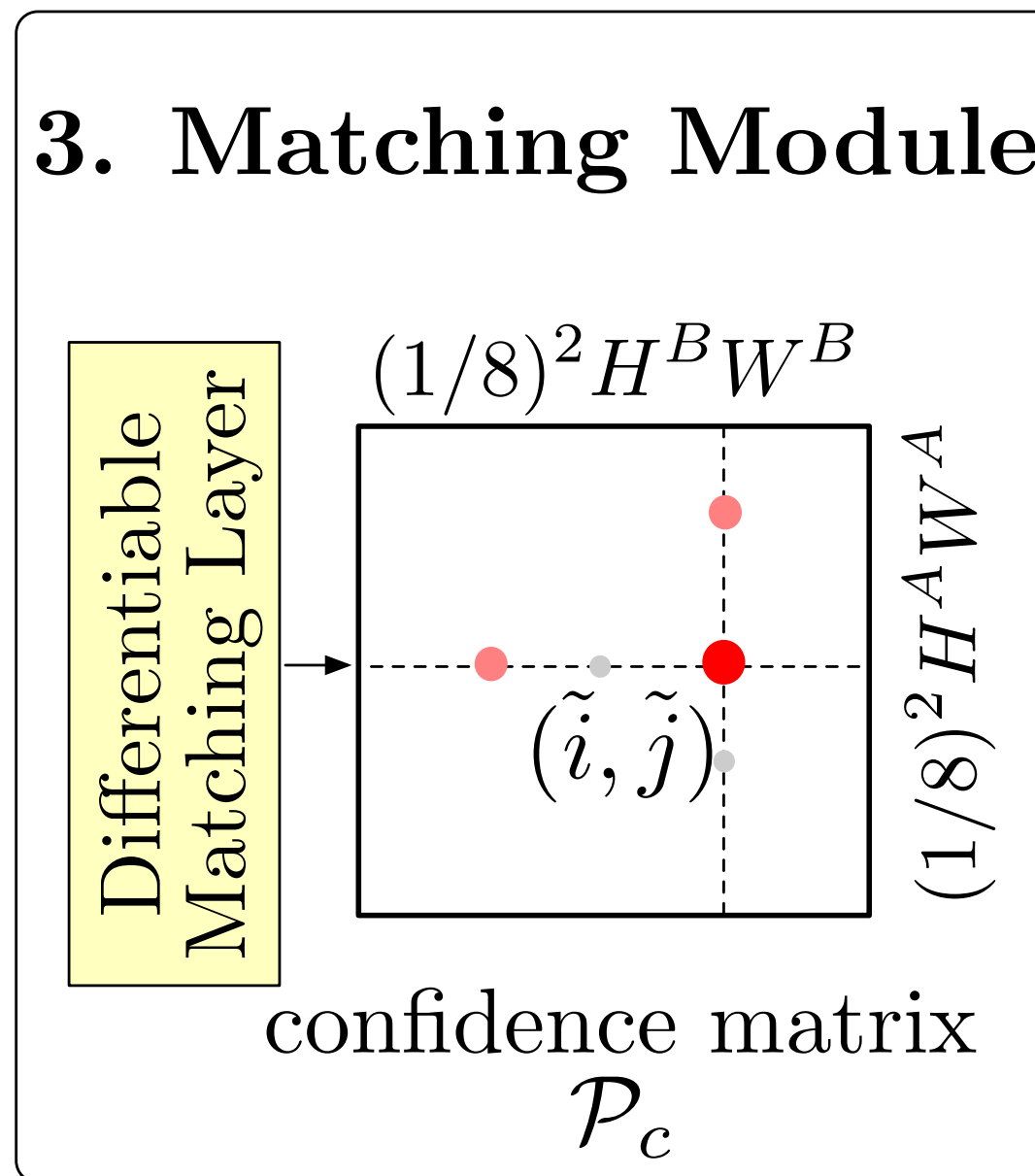
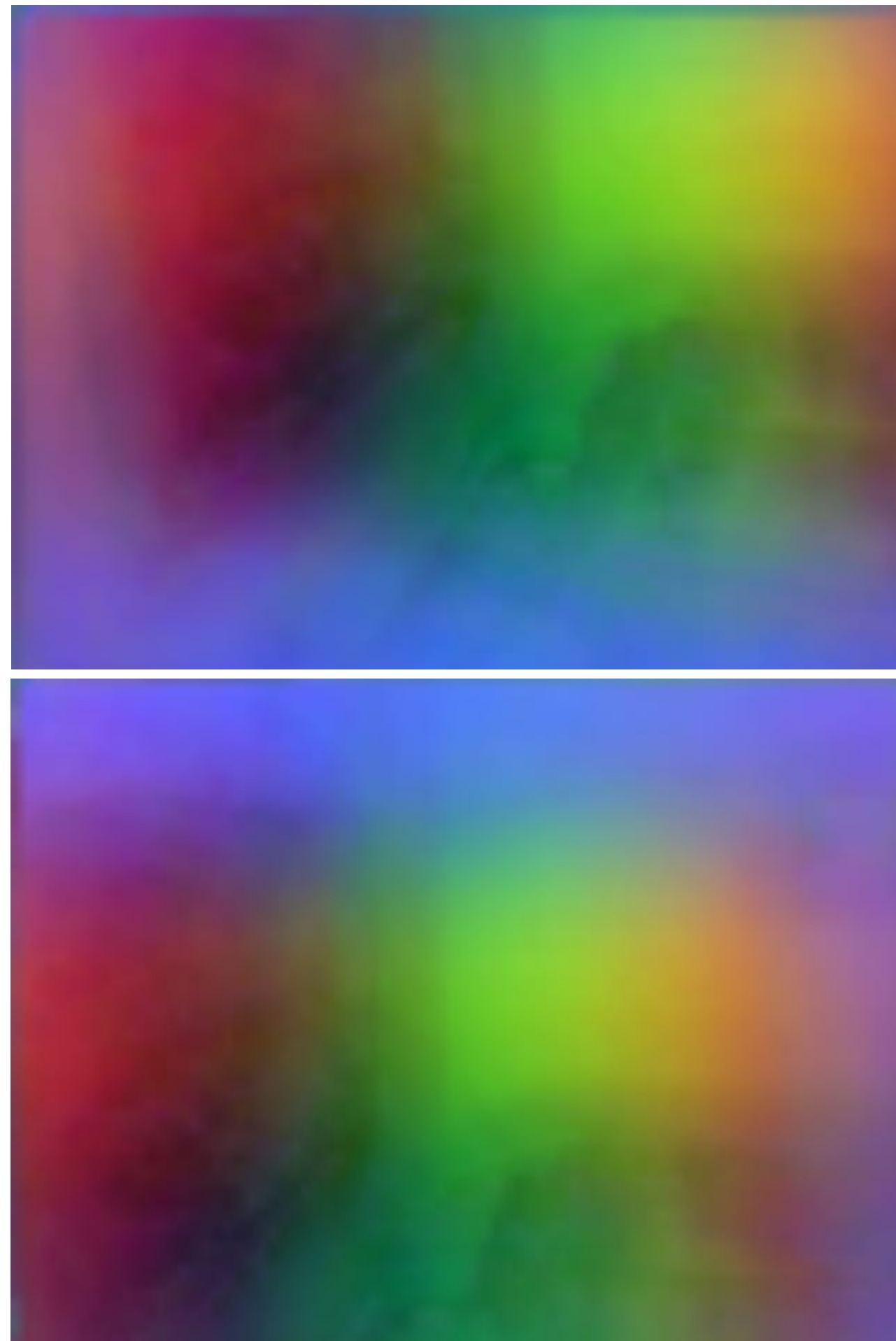
LoFTR

Coarse-level LoFTR module



LoFTR

Differentiable matching layer

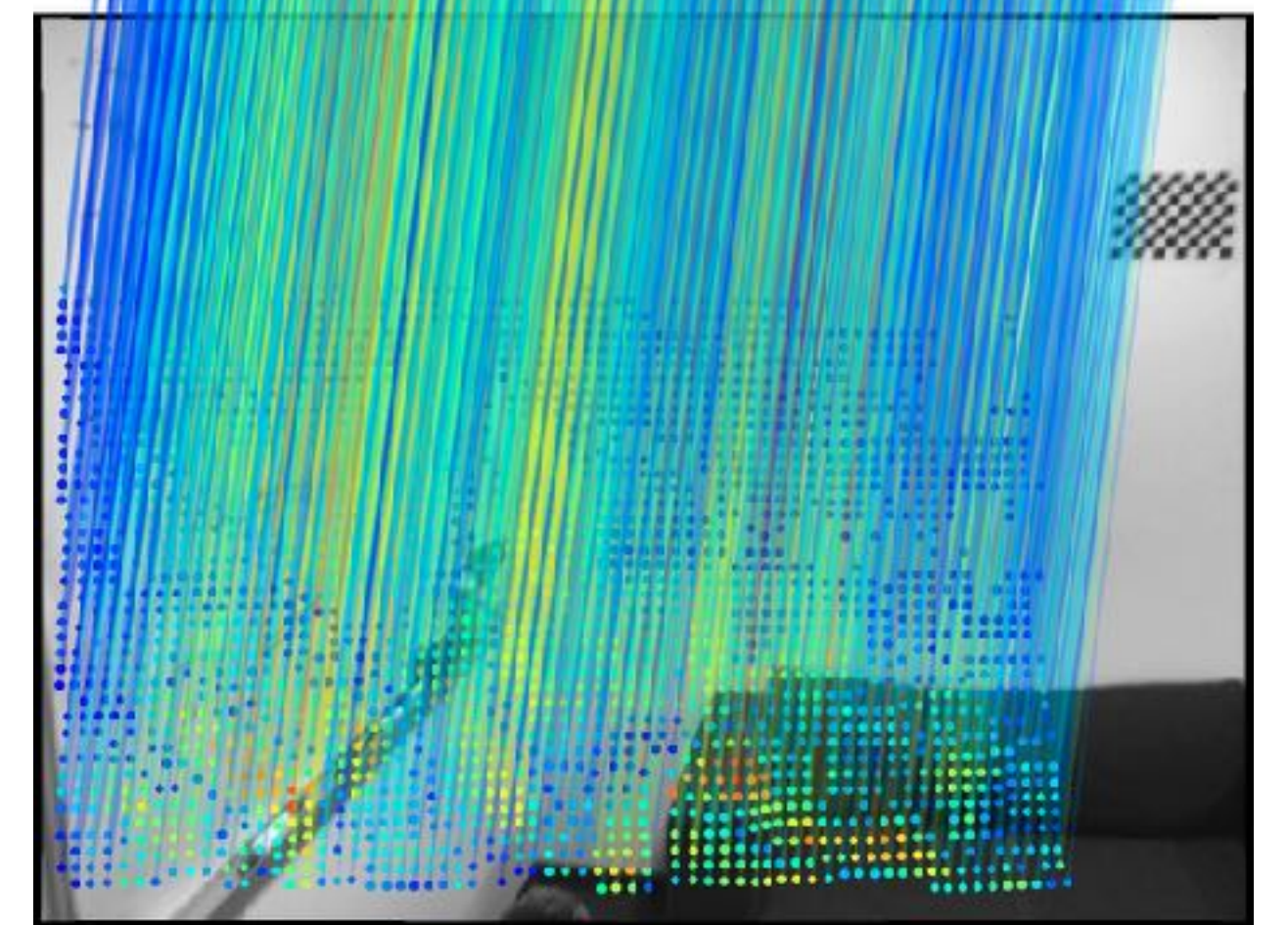
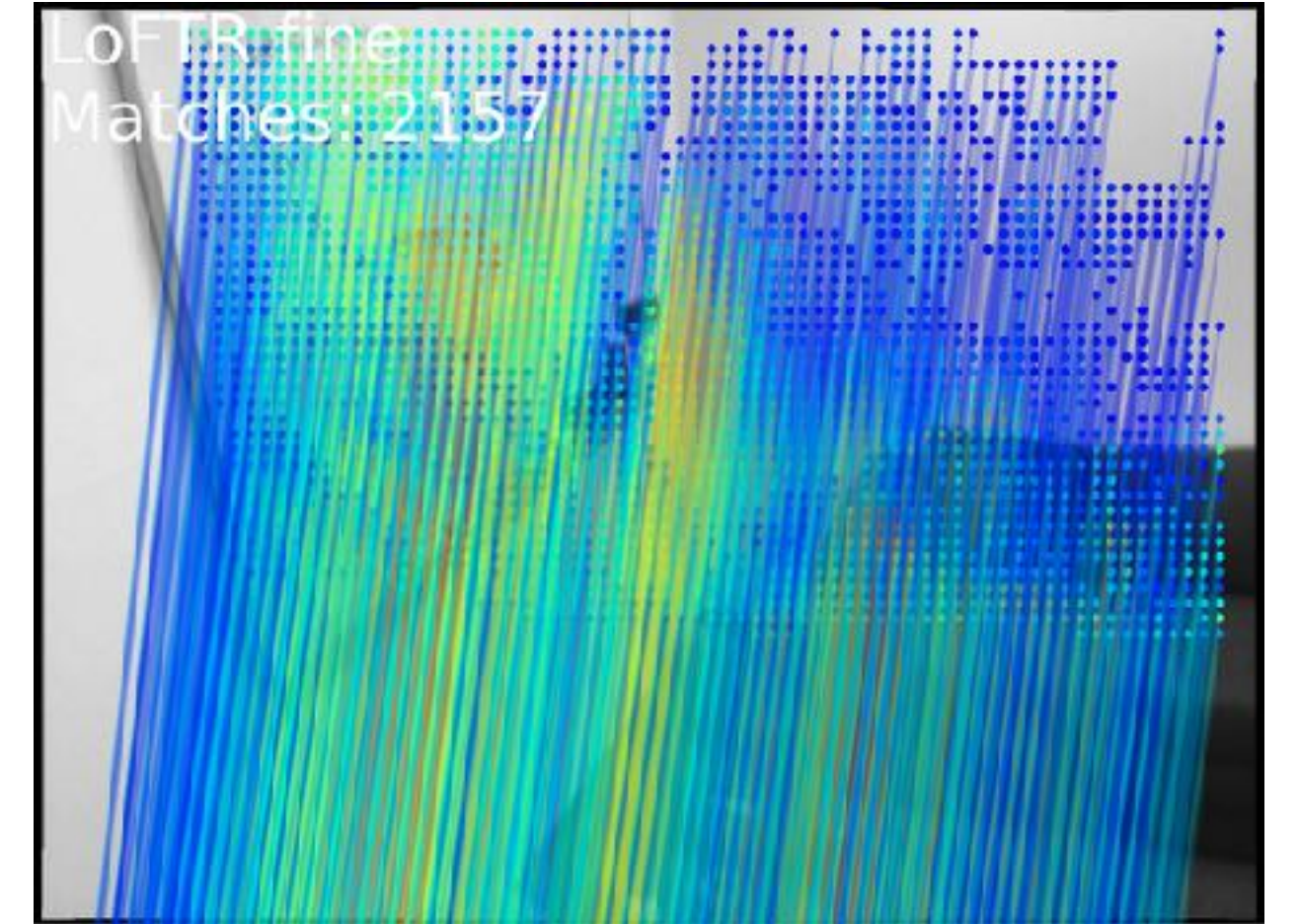
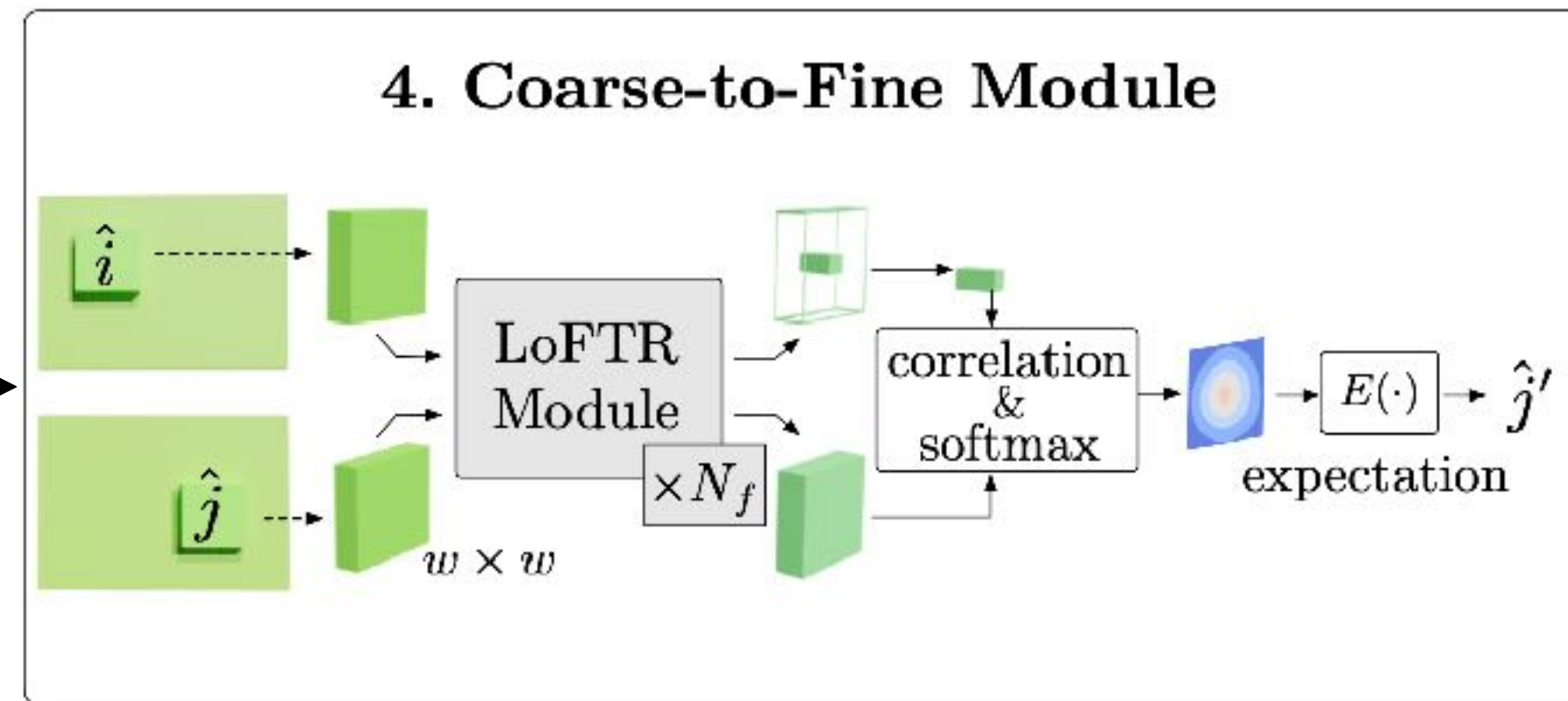
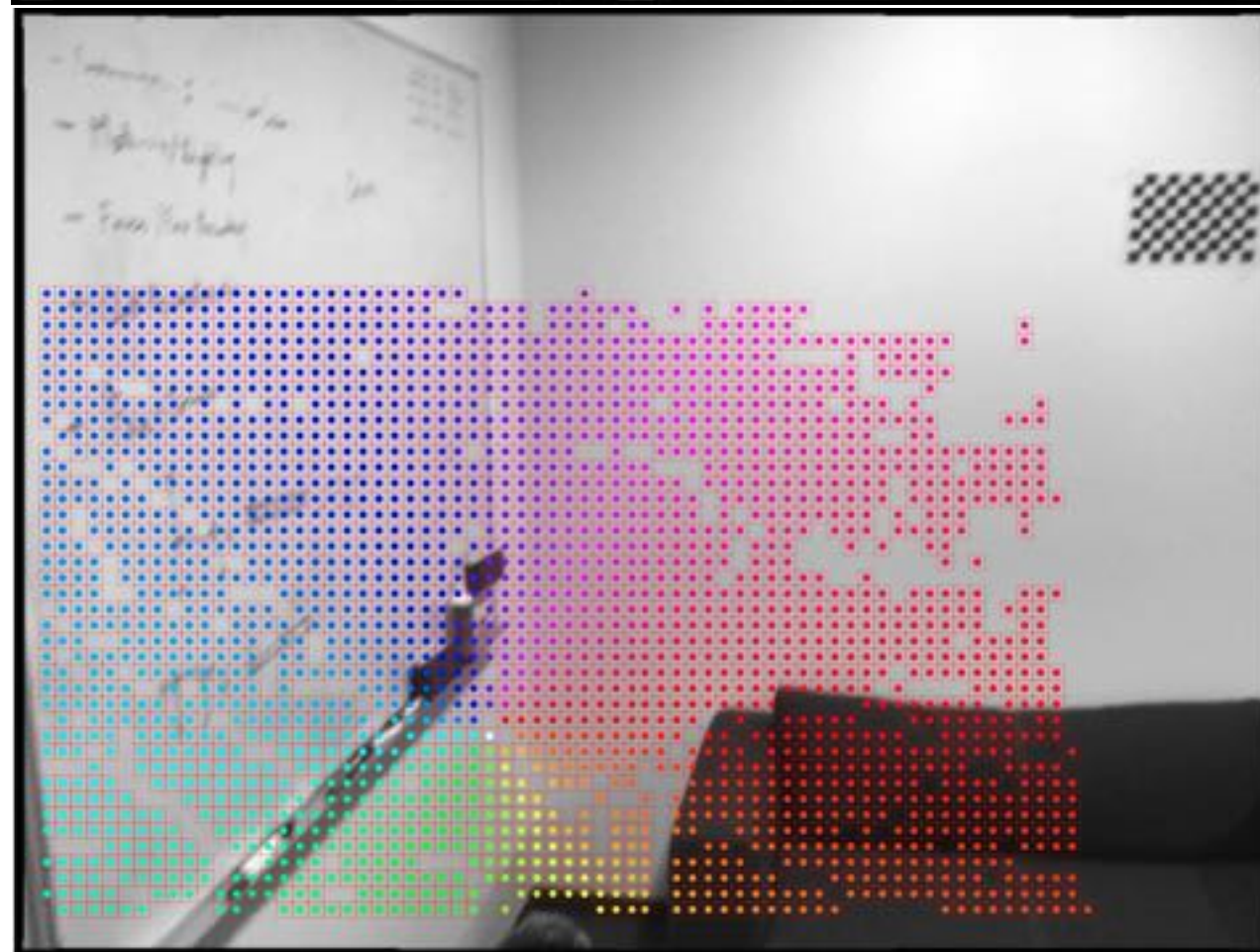


$$\mathcal{P}_c(i, j) = \text{softmax}(\mathcal{S}(i, \cdot))_j \cdot \text{softmax}(\mathcal{S}(\cdot, j))_i$$

Dual-softmax

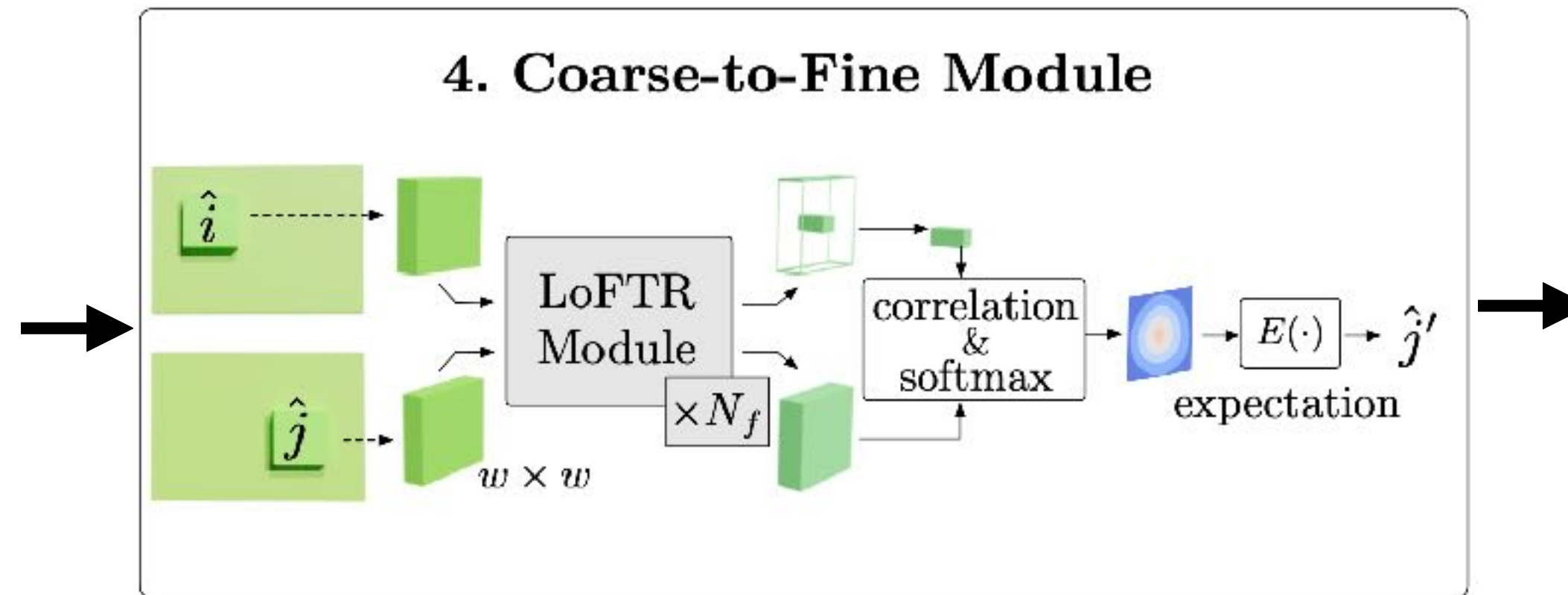
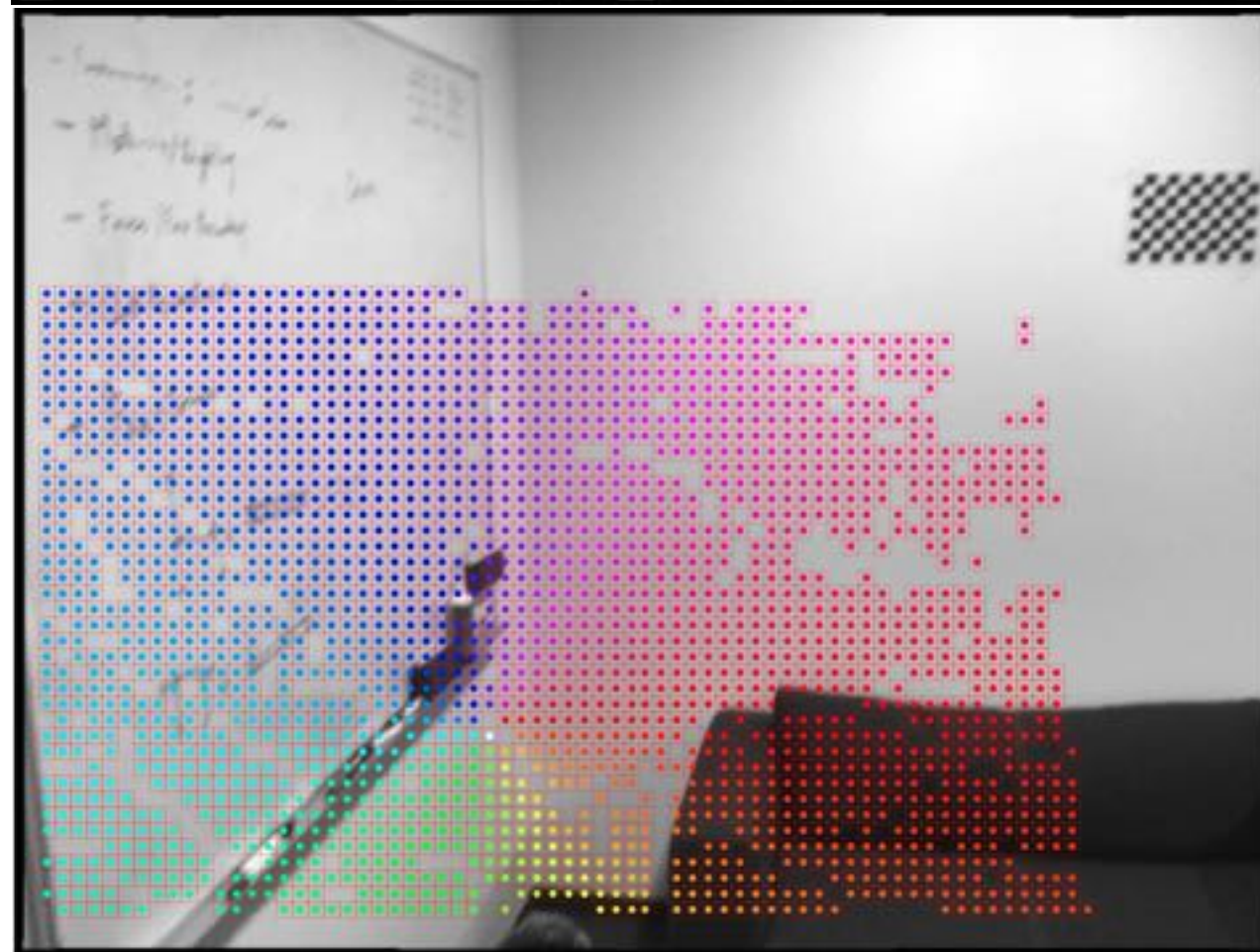
LoFTR

Coarse-to-fine module



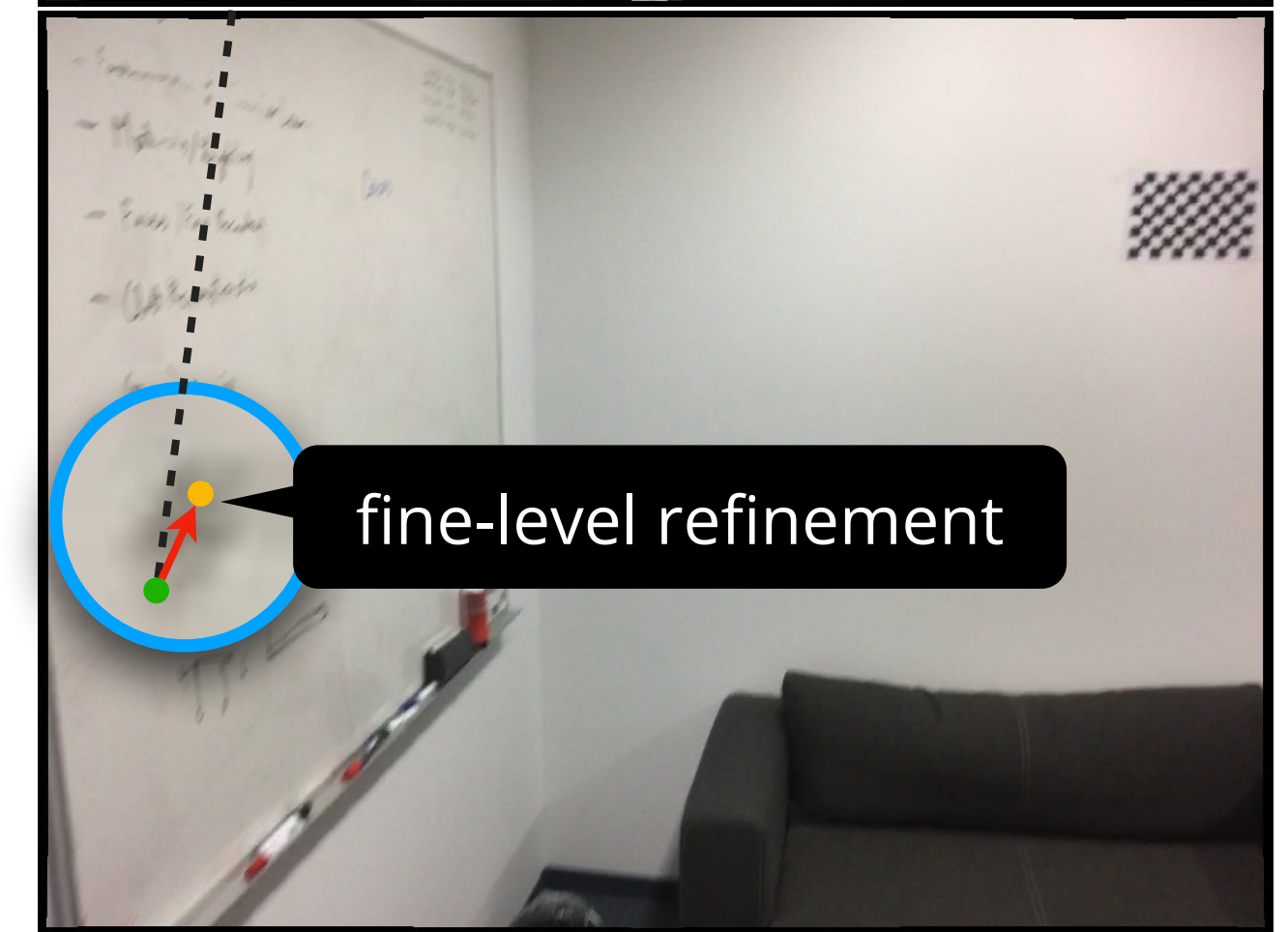
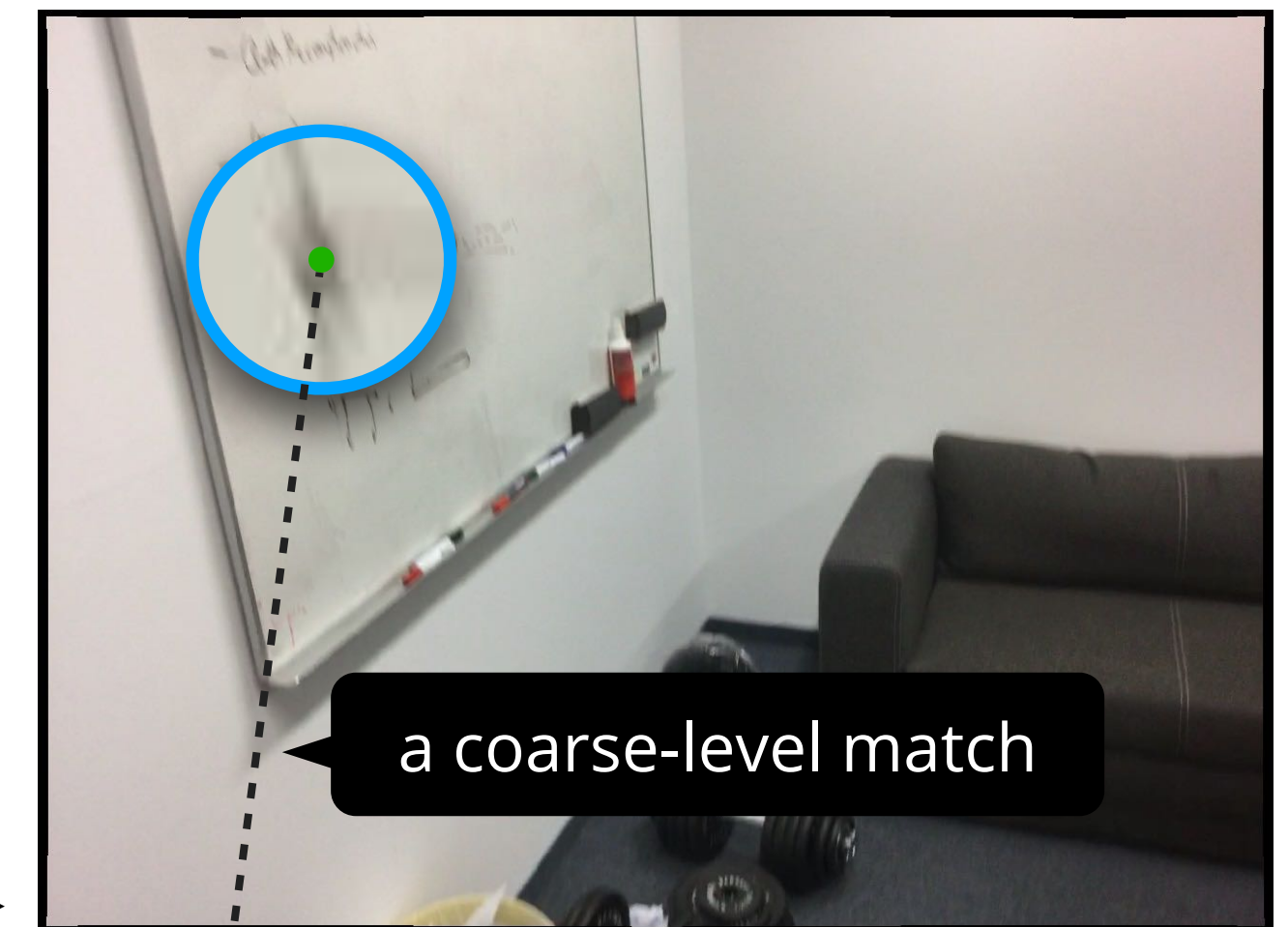
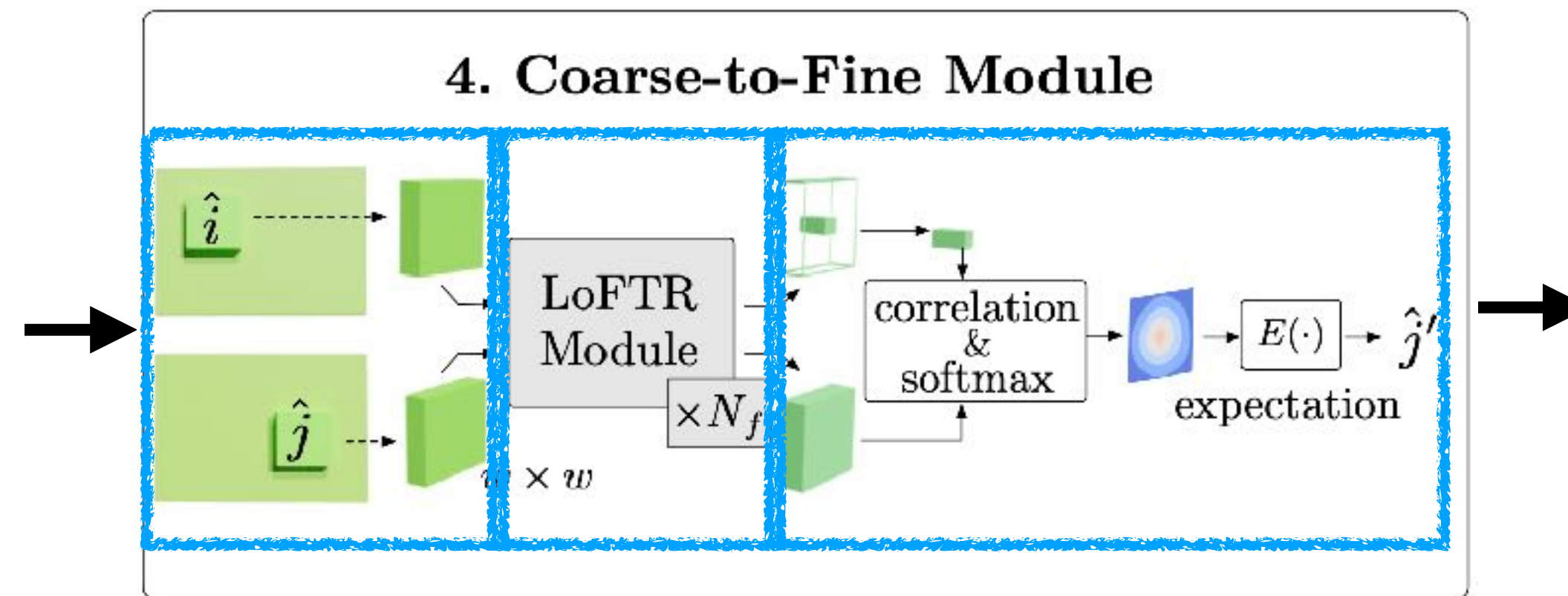
LoFTR

Coarse-to-fine module



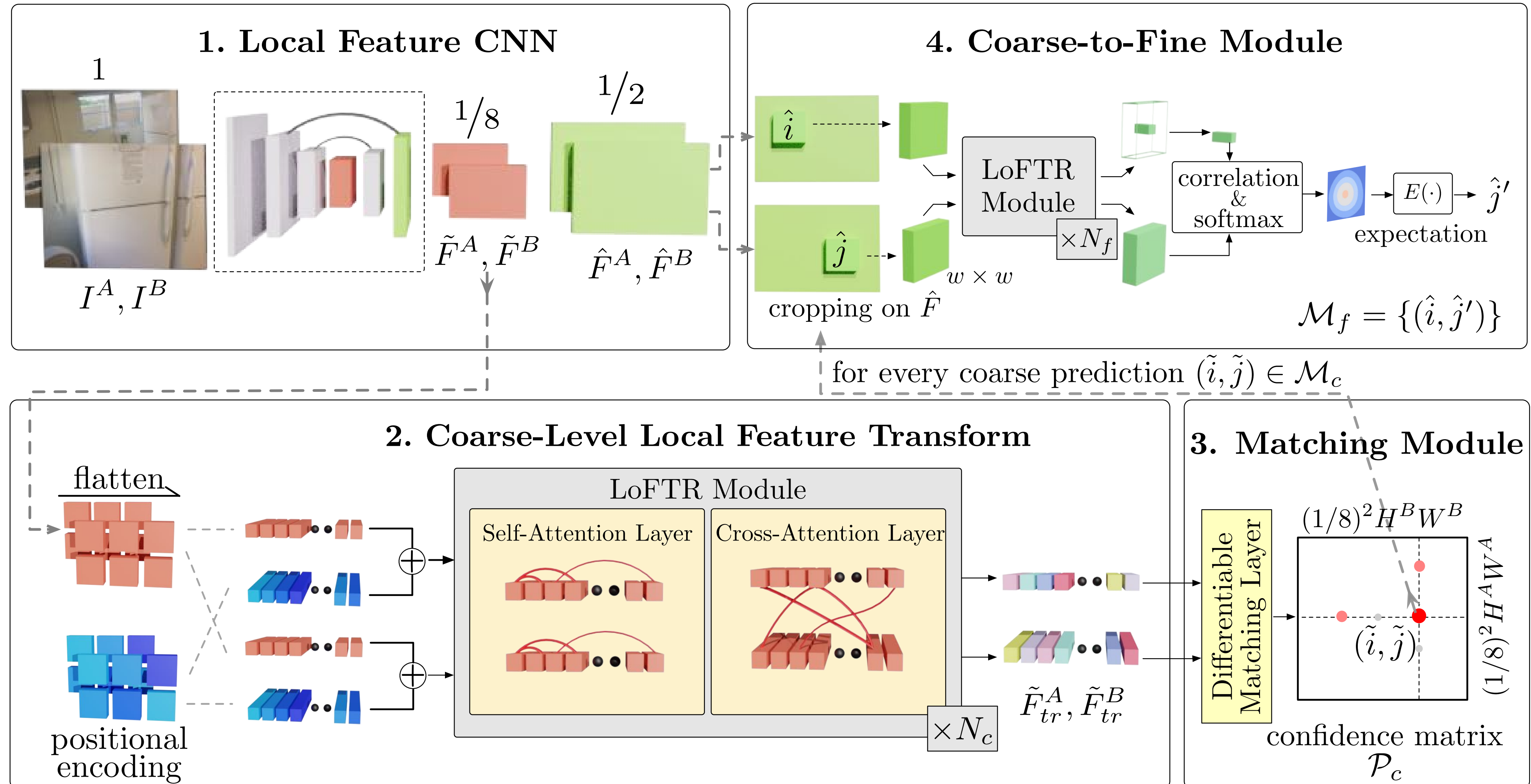
LoFTR

Coarse-to-fine module



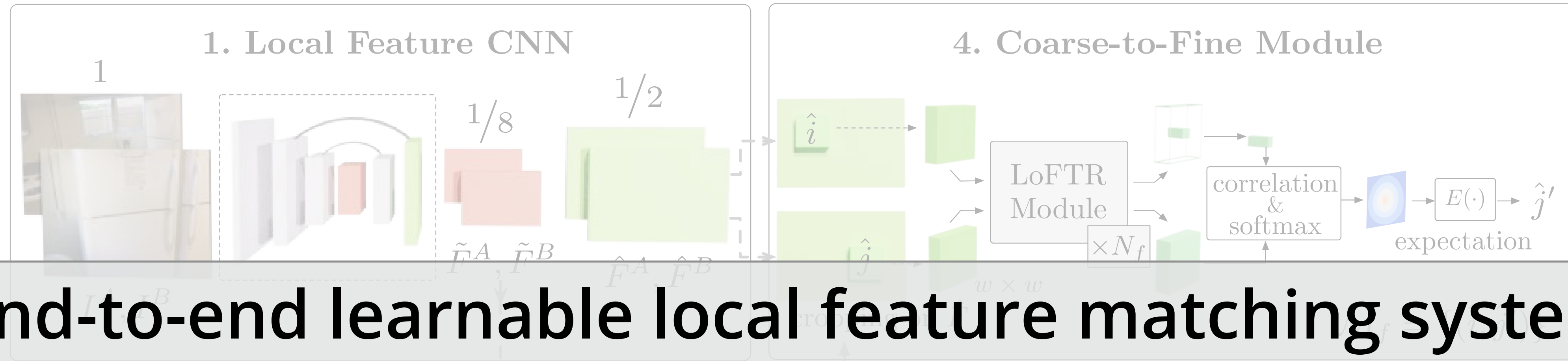
LoFTR

Architecture

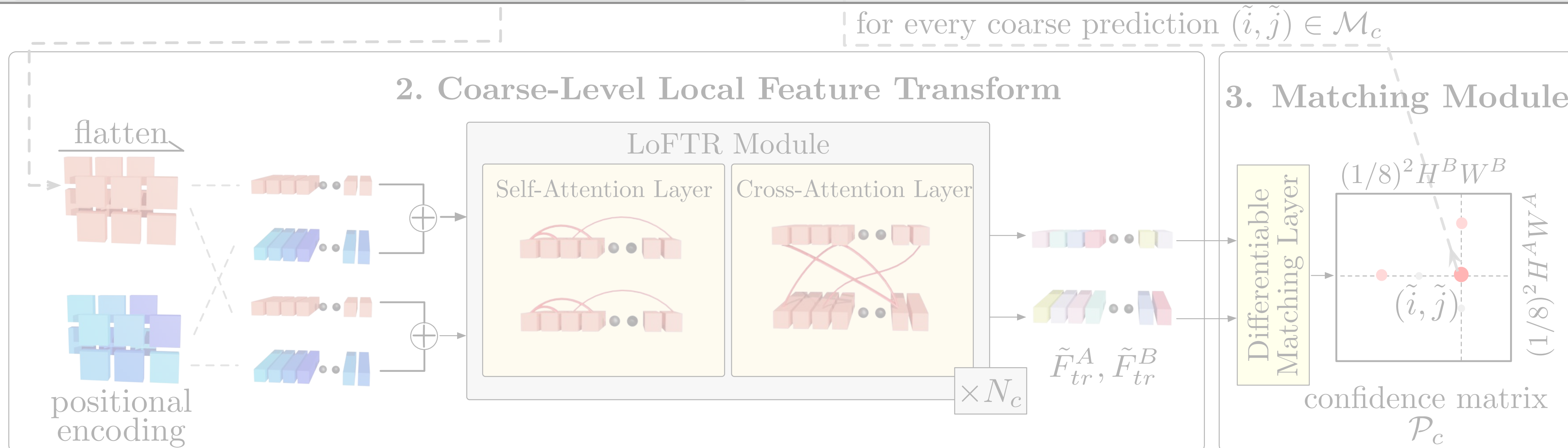


LoFTR

Architecture



End-to-end learnable local feature matching system

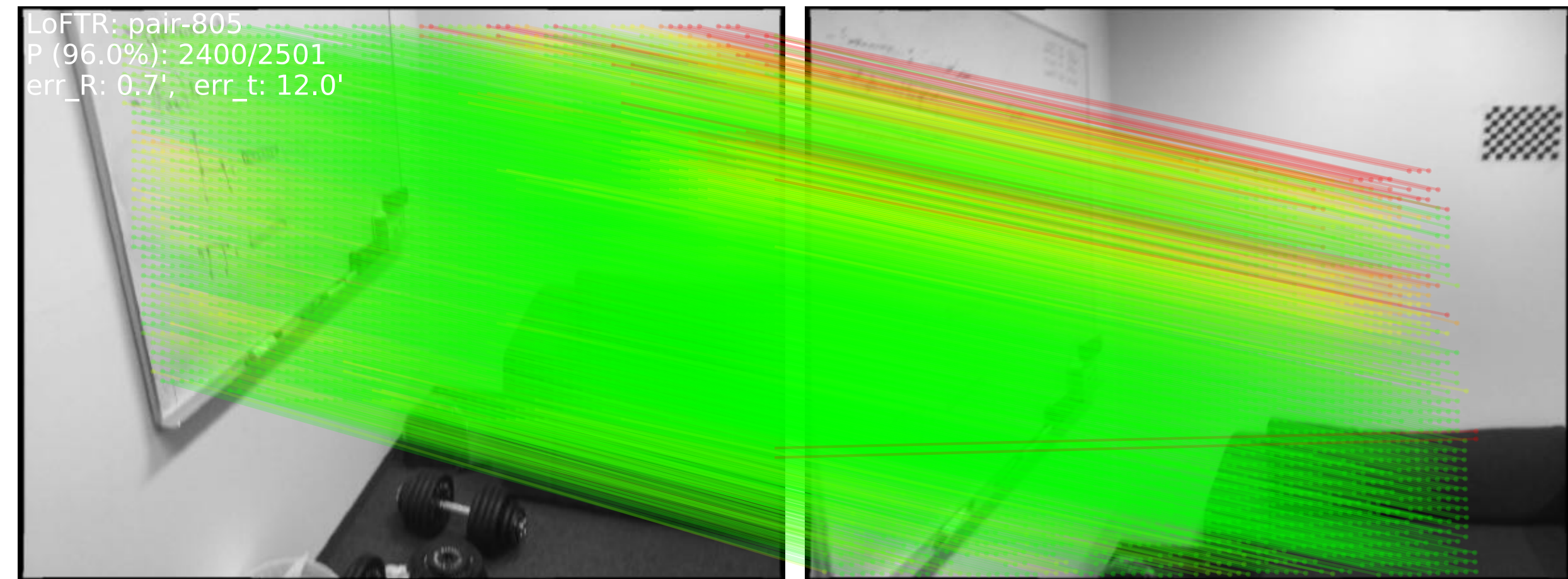
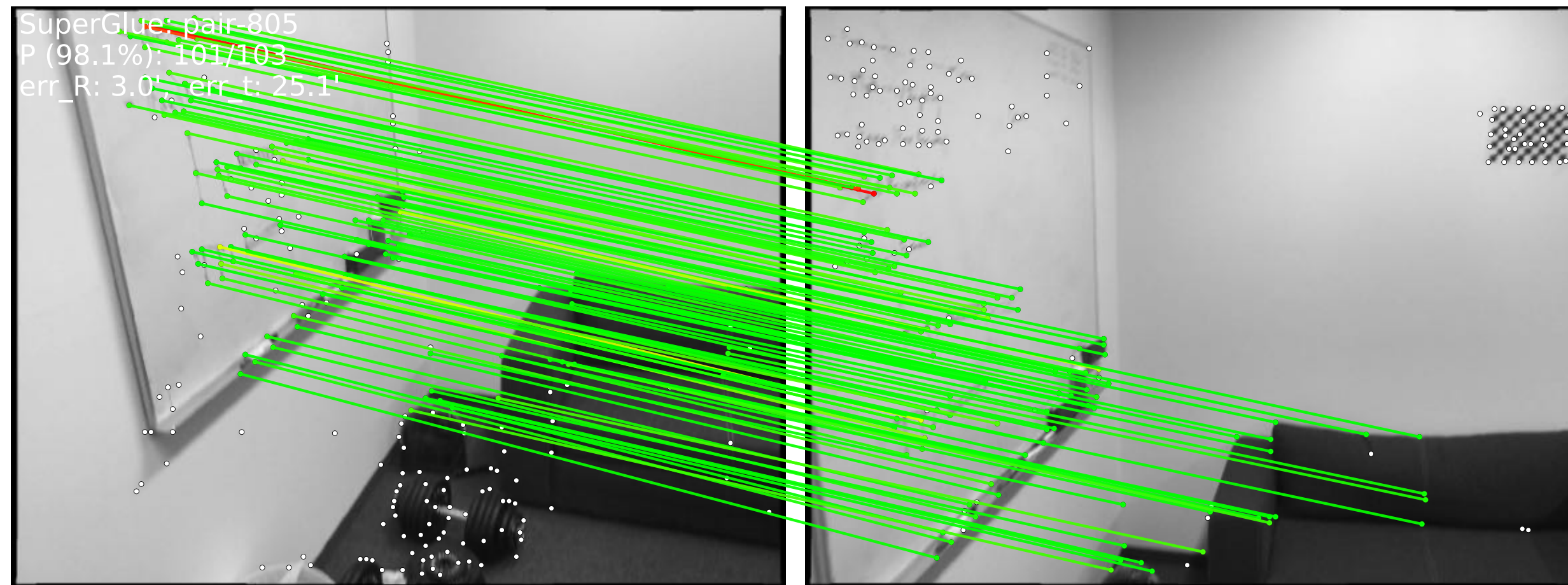


LoFTR

Comparison with SOTA detector-based method

SuperGlue

LoFTR



Pose error: $\Delta R = 3.0^\circ$, $\Delta t = 25.1^\circ$

Pose error: $\Delta R = 0.7^\circ$, $\Delta t = 12.0^\circ$

- 😊 Sparse matches
- 😞 Textured regions only

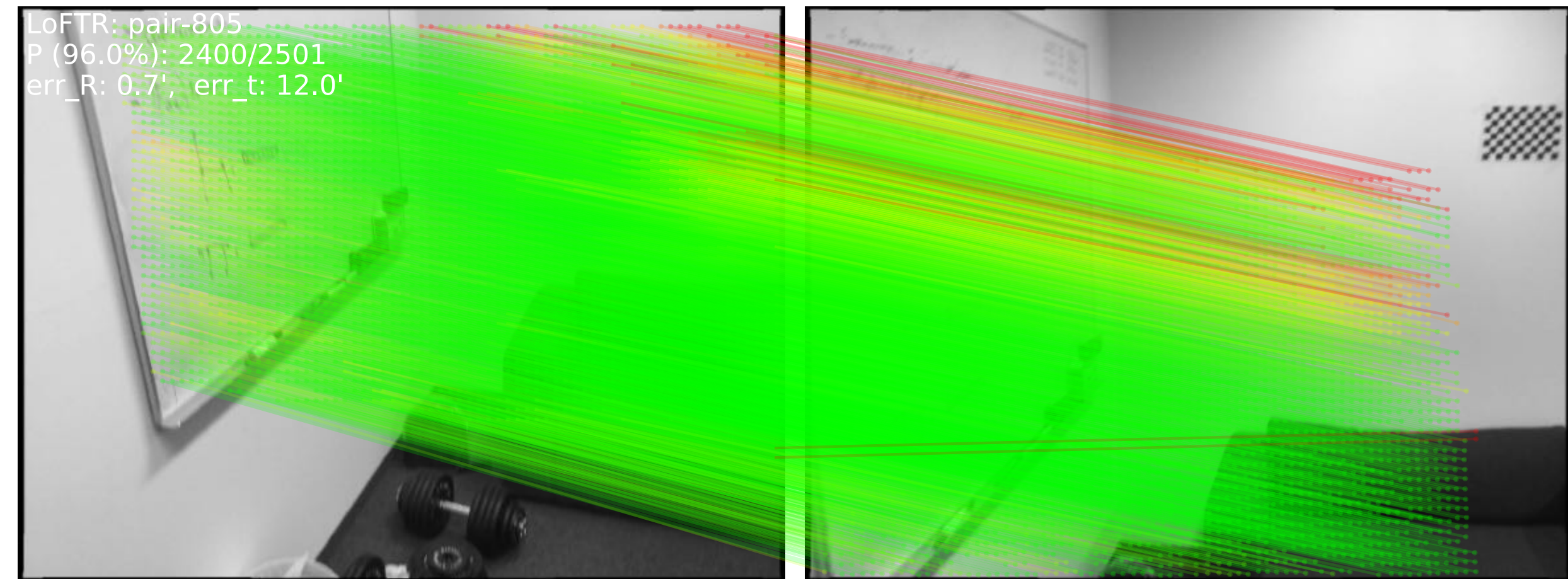
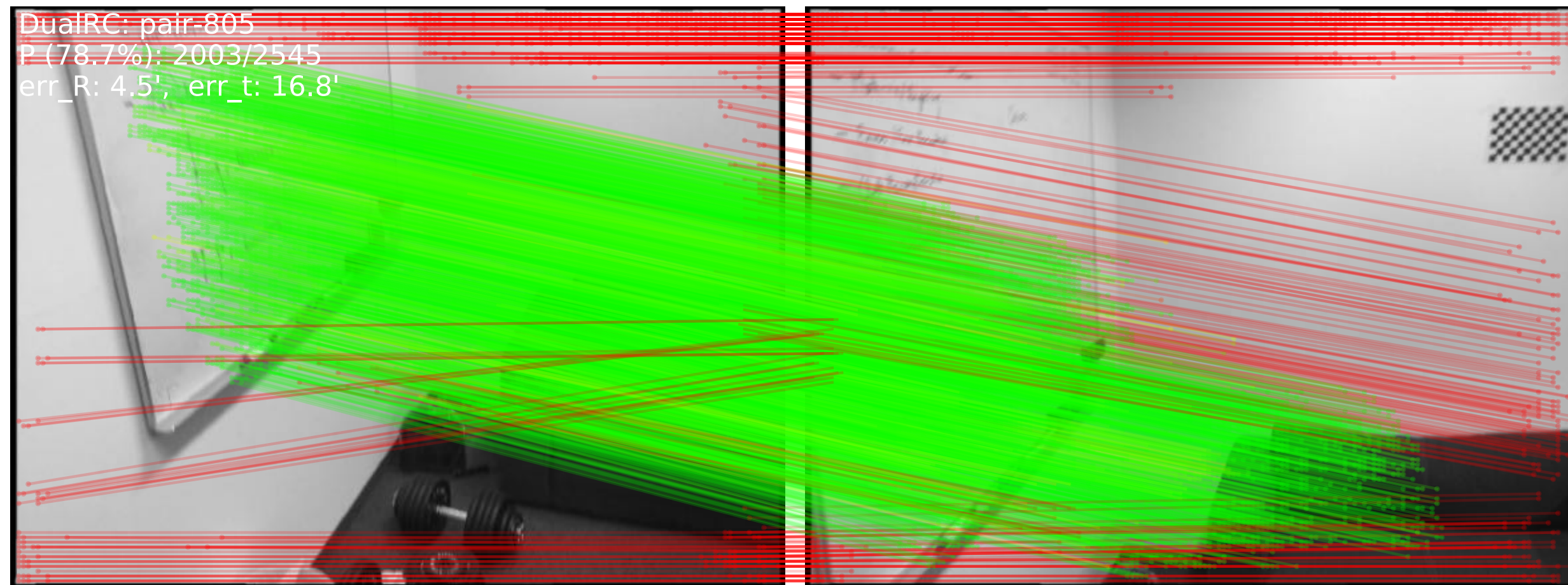
- 😊 Semi-dense matches
- 😊 Textured & Texture-less regions
- 😊 Well-distributed matches

LoFTR

Comparison with SOTA detector-free method

DRC-Net

LoFTR



Pose error: $\Delta R = 4.5^\circ$, $\Delta t = 16.8^\circ$

- 🙄 Local consensus only
 - 👉 many incorrect matches
- 🙄 No matches on purely texture-less regions

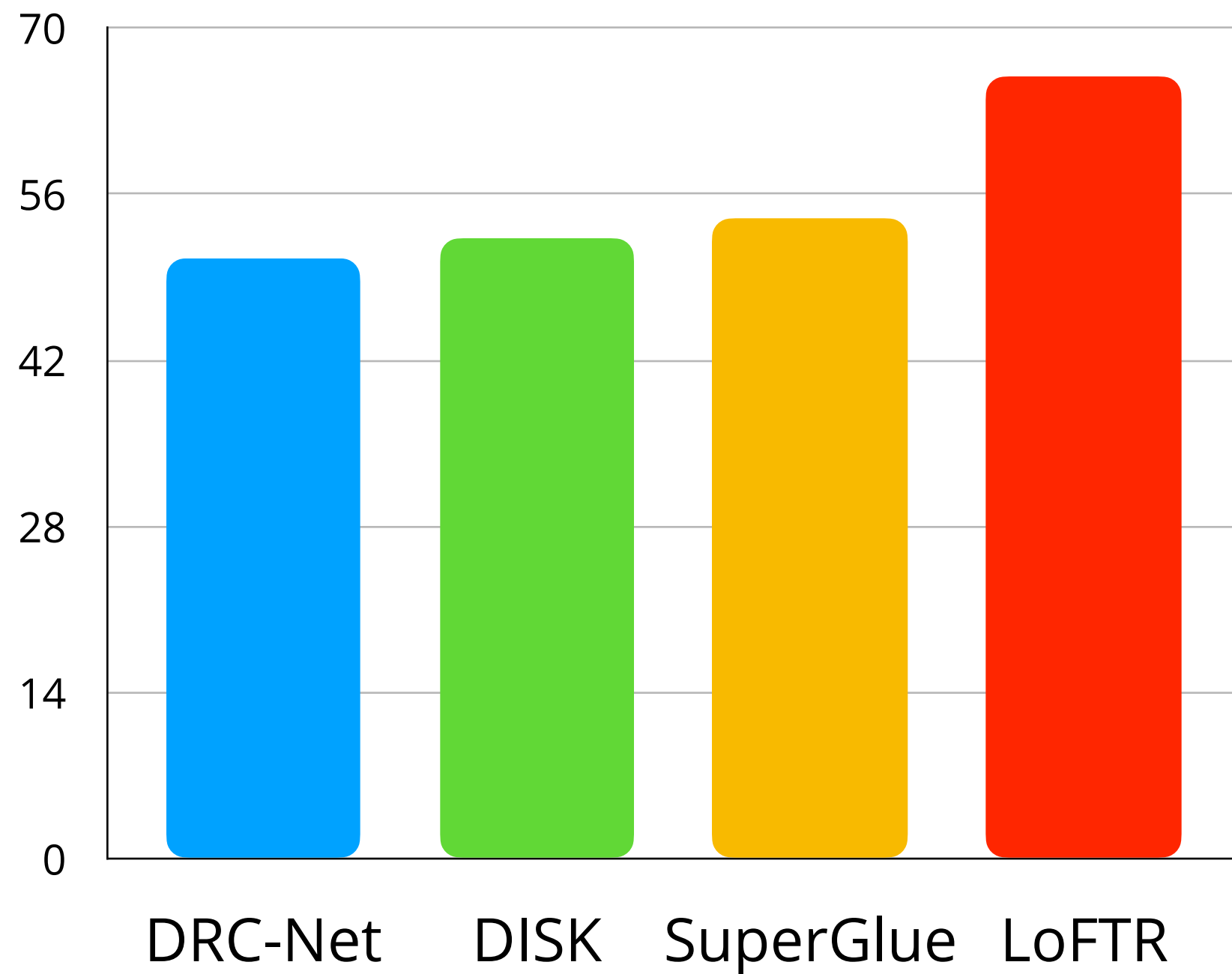
Pose error: $\Delta R = 0.7^\circ$, $\Delta t = 12.0^\circ$

- 😊 Global consensus
 - 👉 mostly correct matches
- 😊 Semi-dense matches on purely texture-less regions

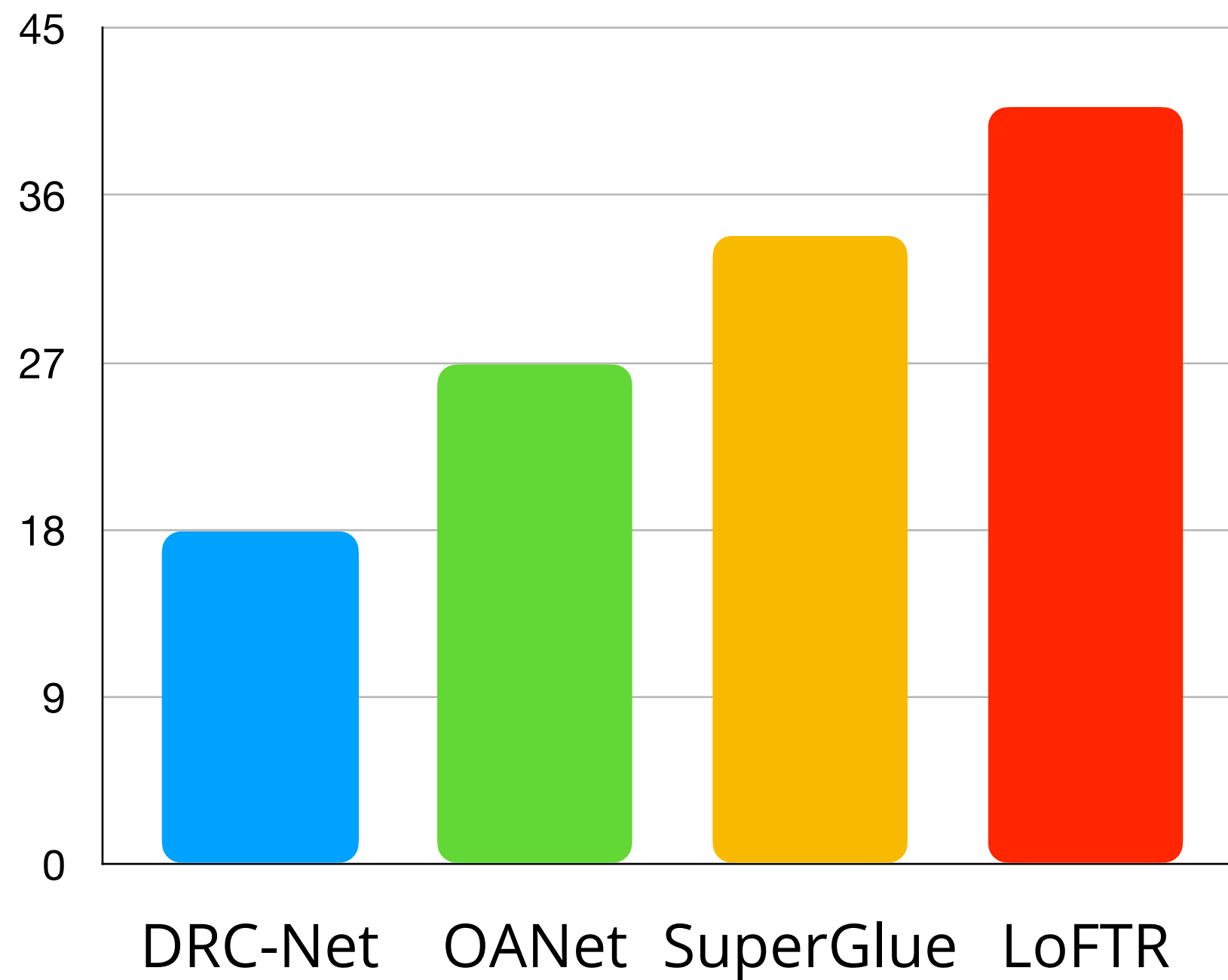
Experiments

Relative pose estimation

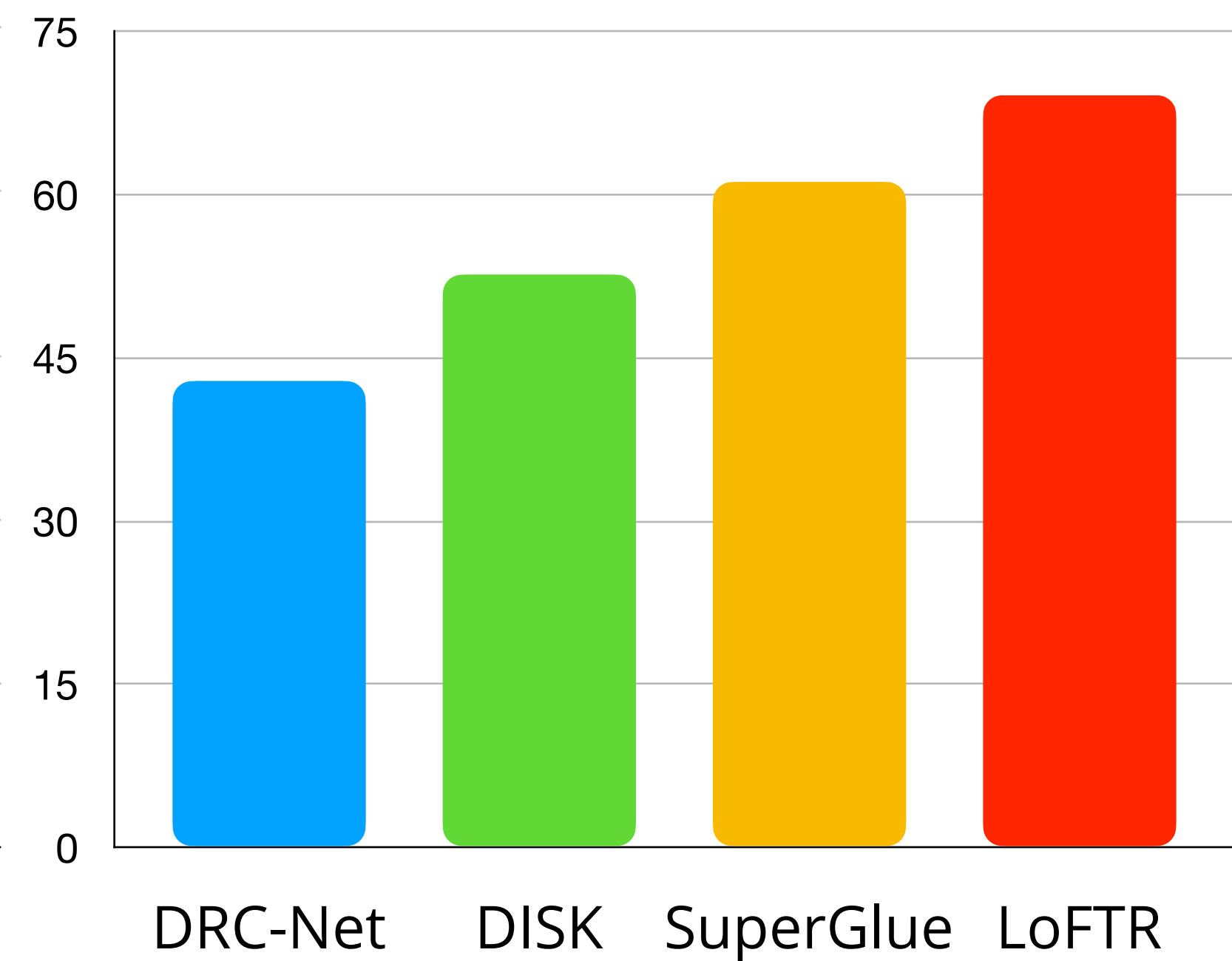
HPatches (AUC@3px)



ScanNet (AUC@10°)



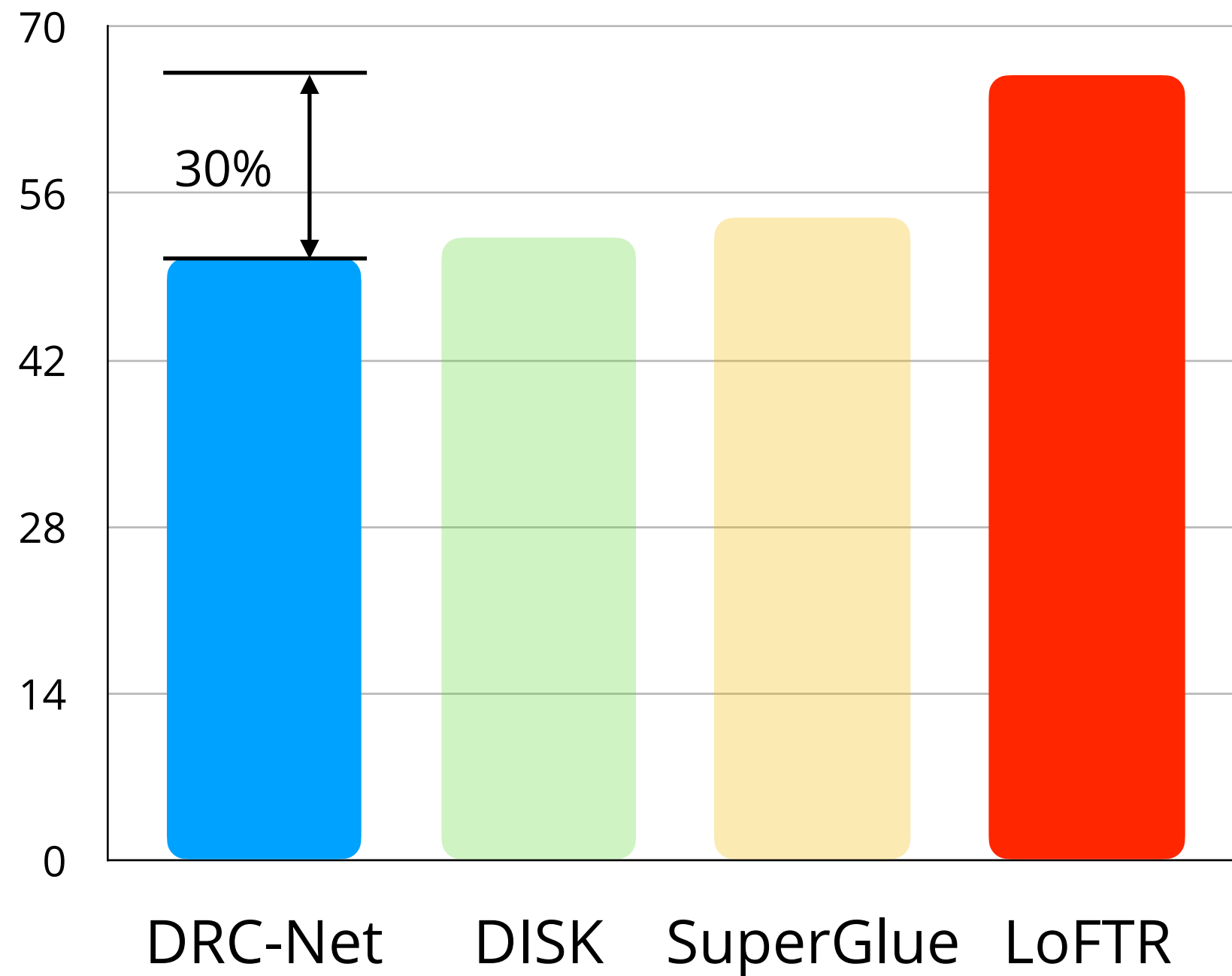
MegaDepth (AUC@10°)



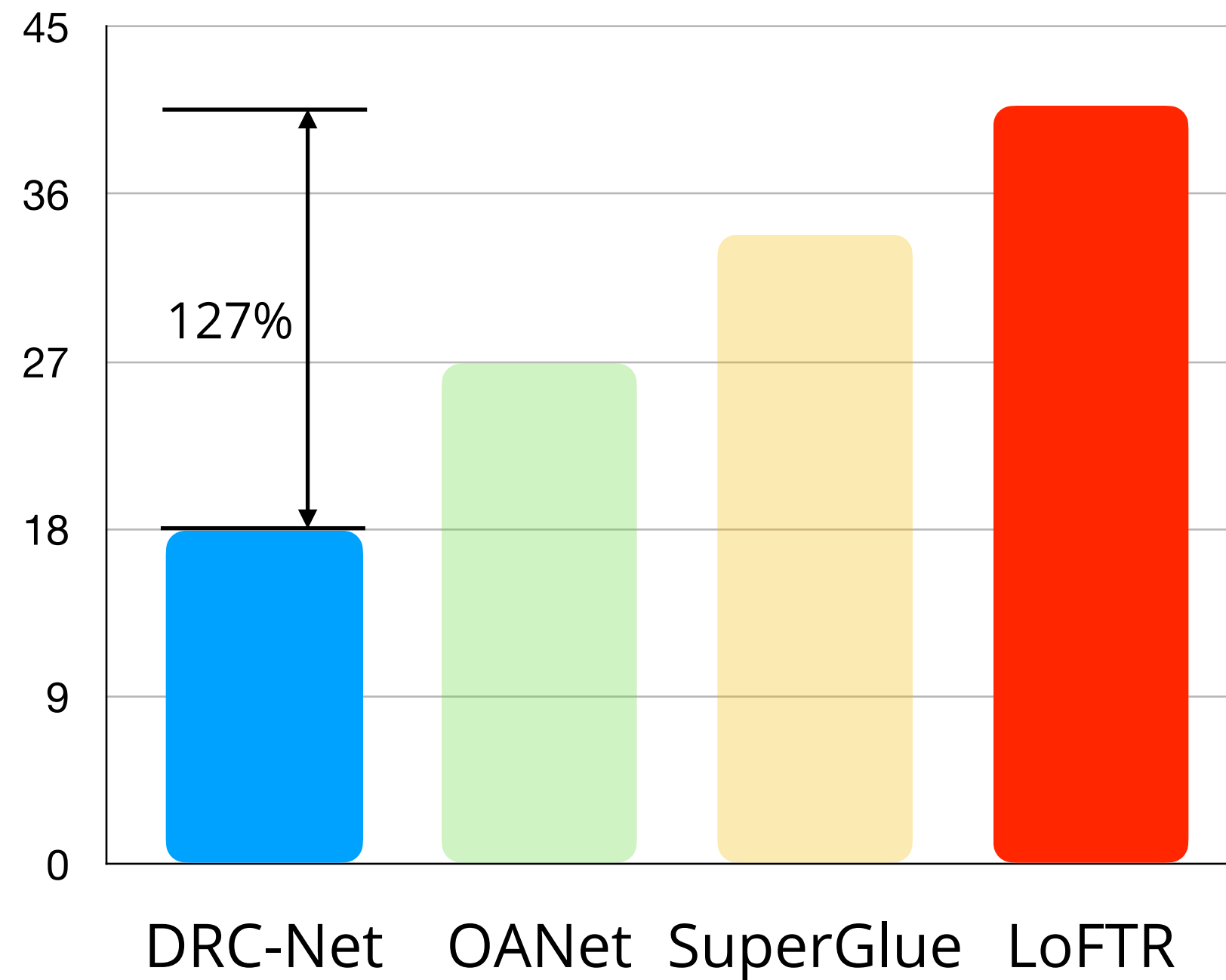
Experiments

Relative pose estimation

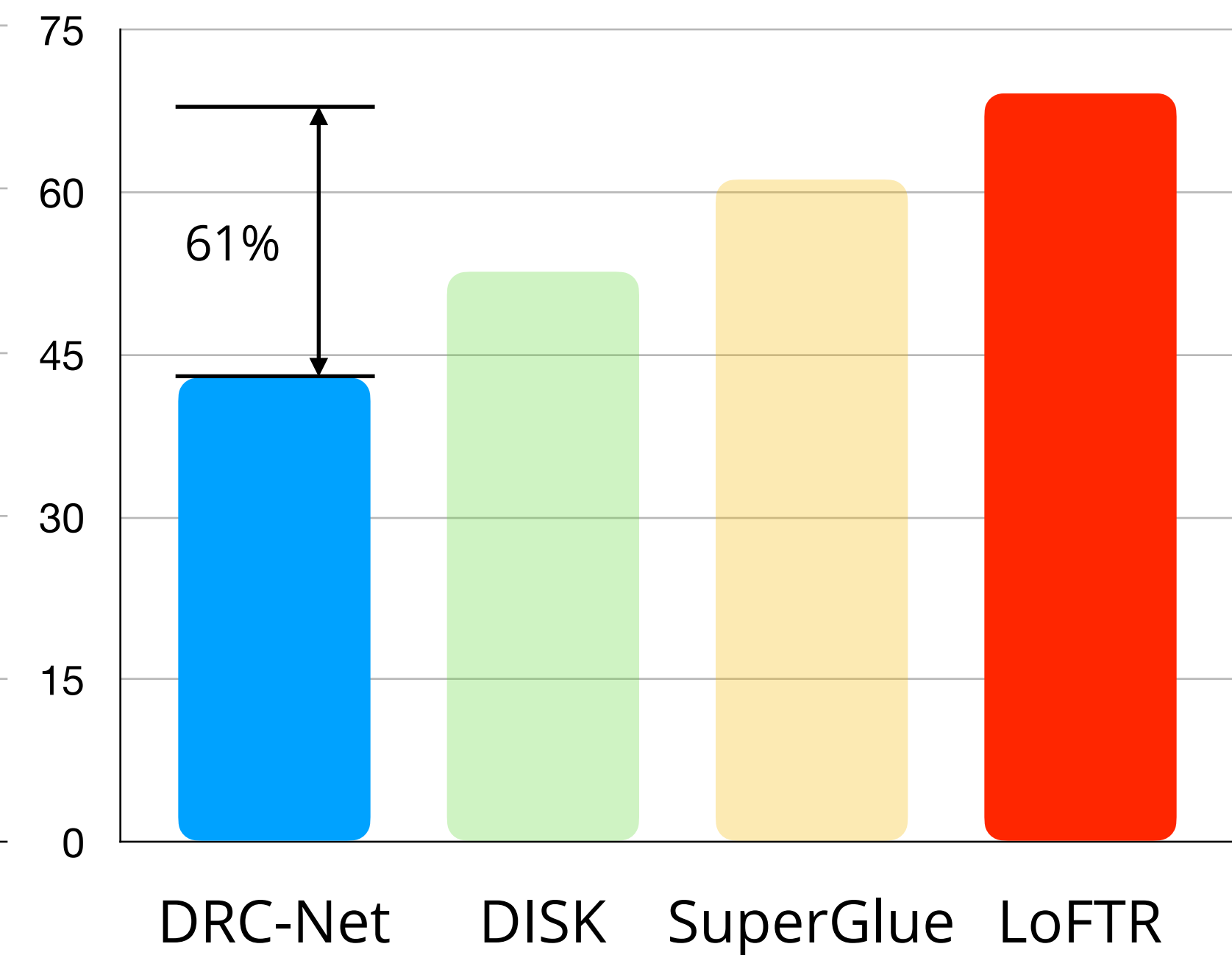
HPatches (AUC@3px)



ScanNet (AUC@10°)



MegaDepth (AUC@10°)

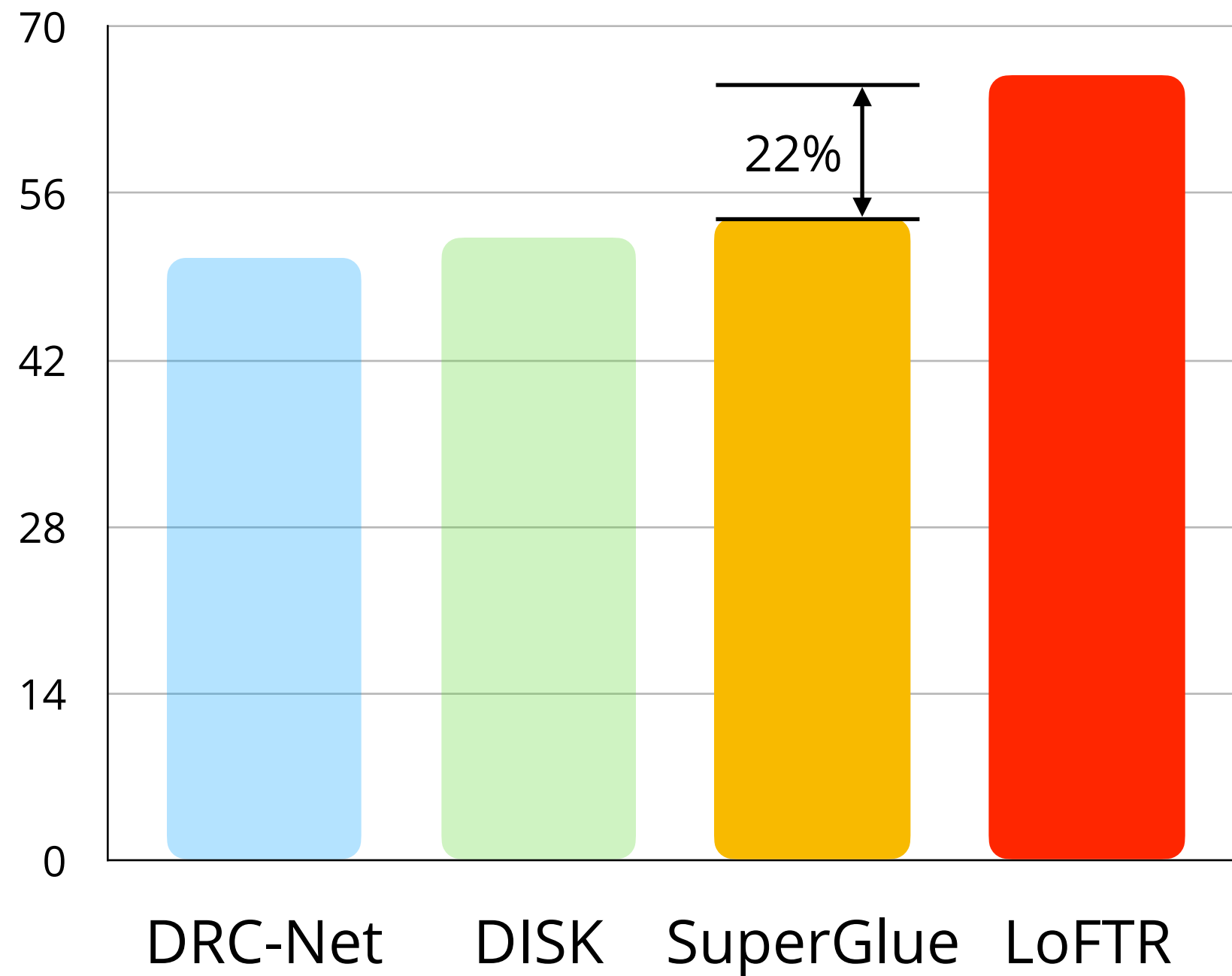


Improvements (in ratio) compared to the SOTA detector-free baseline DRC-Net

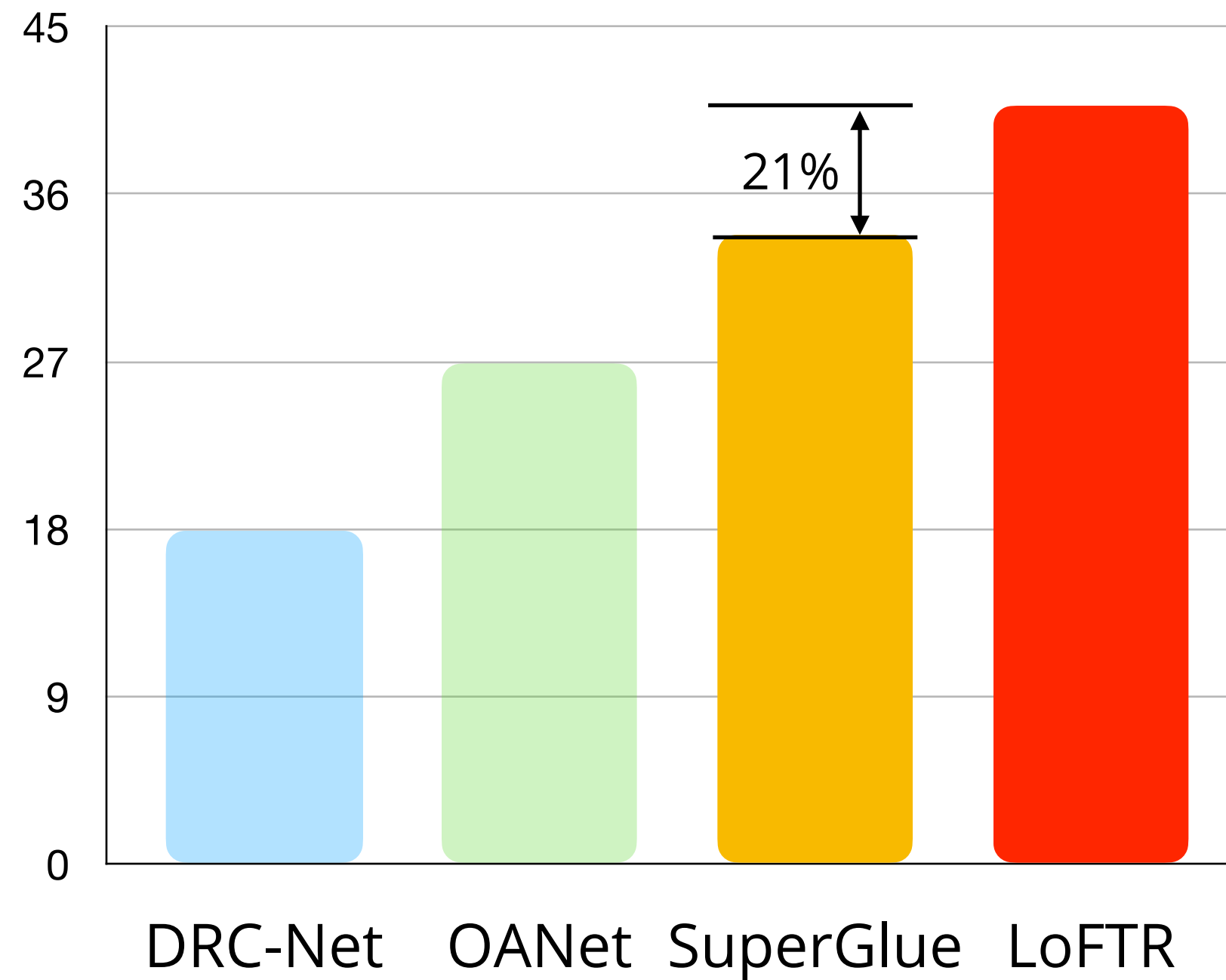
Experiments

Relative pose estimation

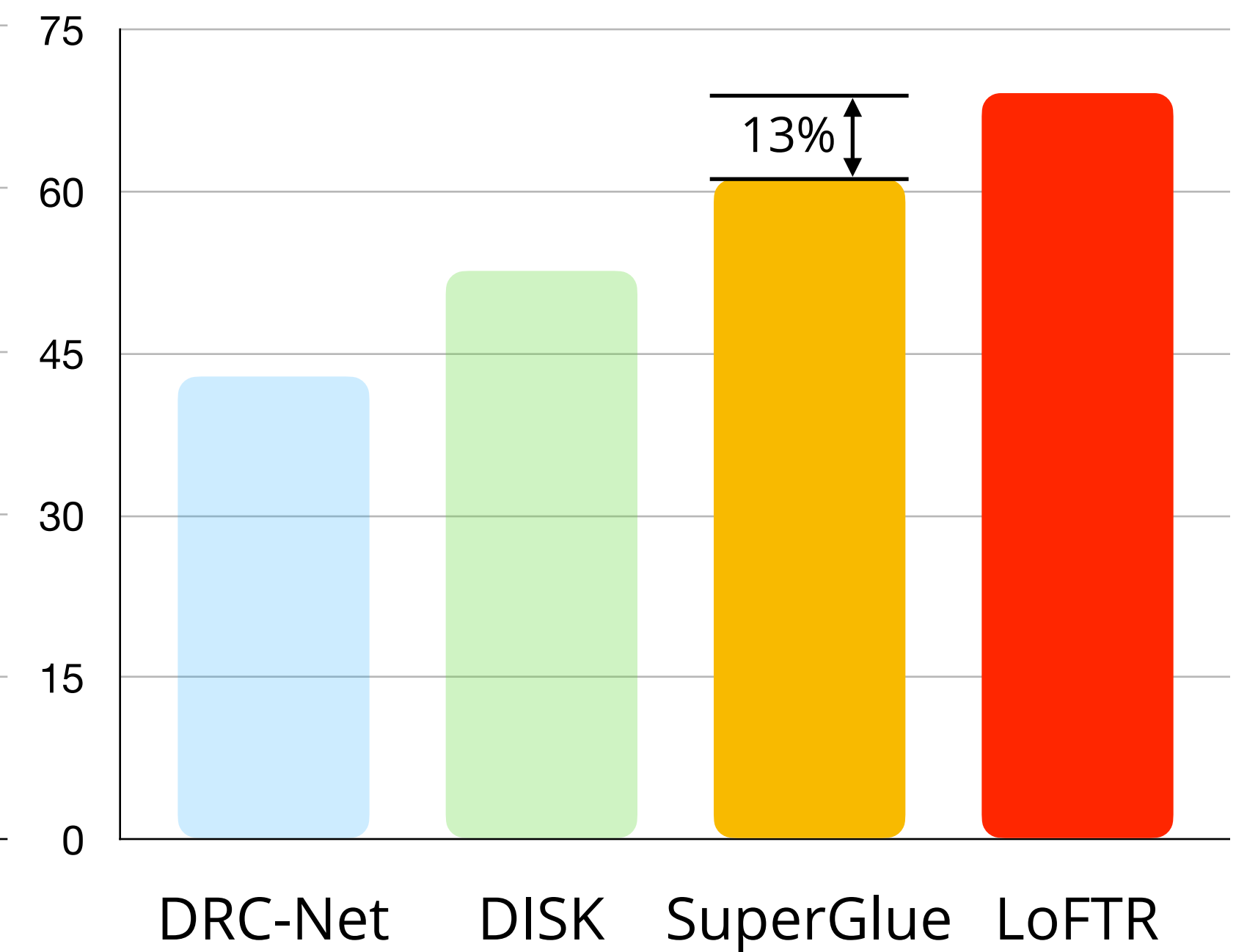
HPatches (AUC@3px)



ScanNet (AUC@10°)



MegaDepth (AUC@10°)

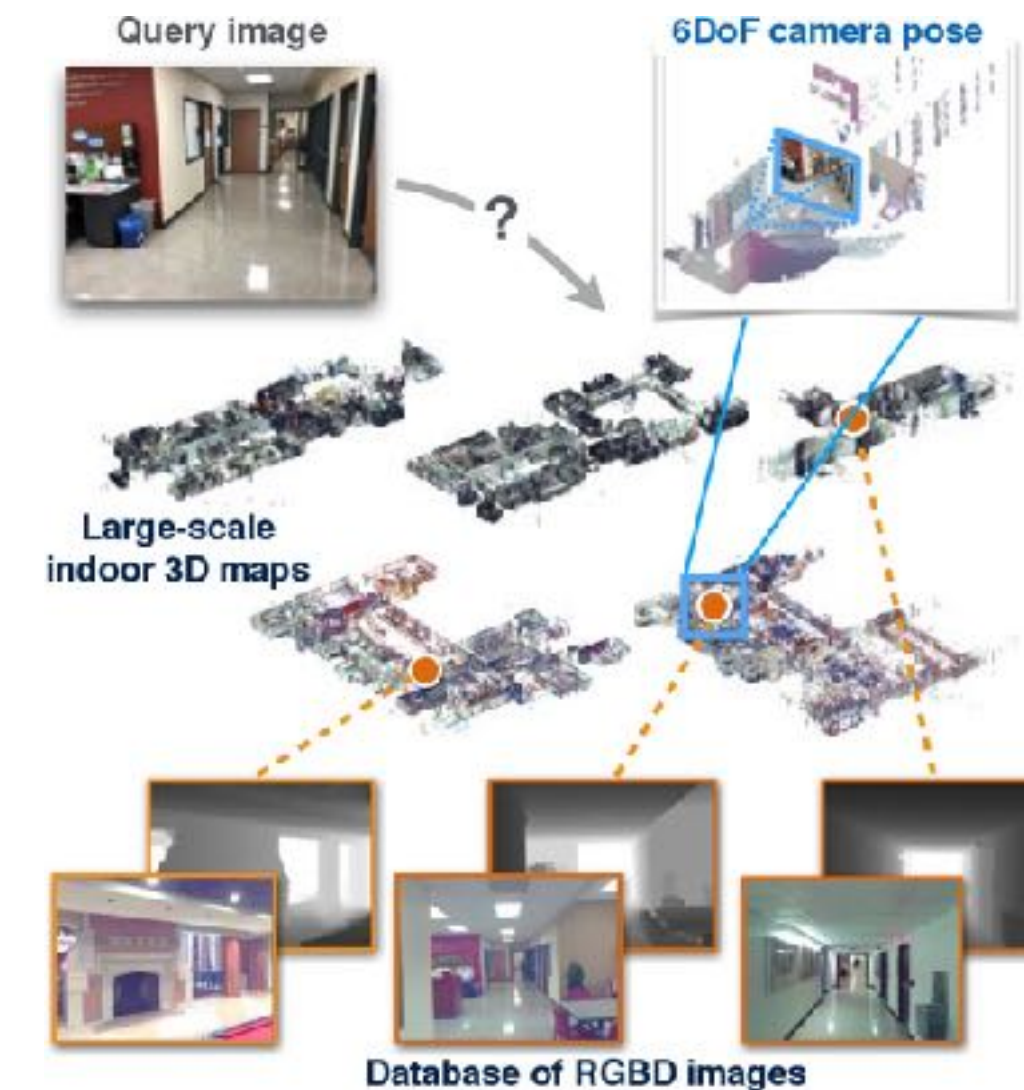


Improvements (in ratio) compared to the SOTA detector-based baseline SuperGlue






Experiments

Visual localization






Rank  on two VisLoc benchmarks* !

VisLoc Local Feature Challenge (Aachen)

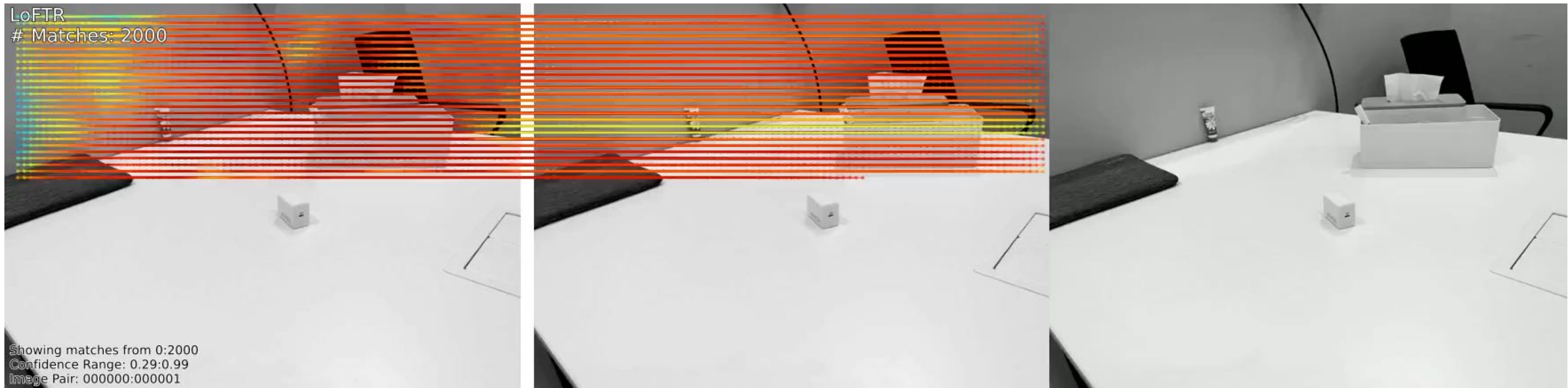
	Night			RANK
	0.25m, 2° / 0.5m, 5° / 1.0m, 10°			
LoFTR	72.8	88.5	99	
SuperGlue	73.3	88.0	98.4	
R2D2	71.2	86.9	98.9	

VisLoc Indoor Localization (InLoc)

	DUC1			DUC2			RANK
	0.25m, 2° / 5° / 10°			0.25m, 2° / 5° / 10°			
LoFTR + HLOC	47.5	72.2	84.8	54.2	74.8	85.5	
SuperGlue + HLOC	49.0	68.7	80.8	53.4	77.1	82.4	
R2D2 + KAPTURE	41.4	60.1	73.7	47.3	67.2	73.3	

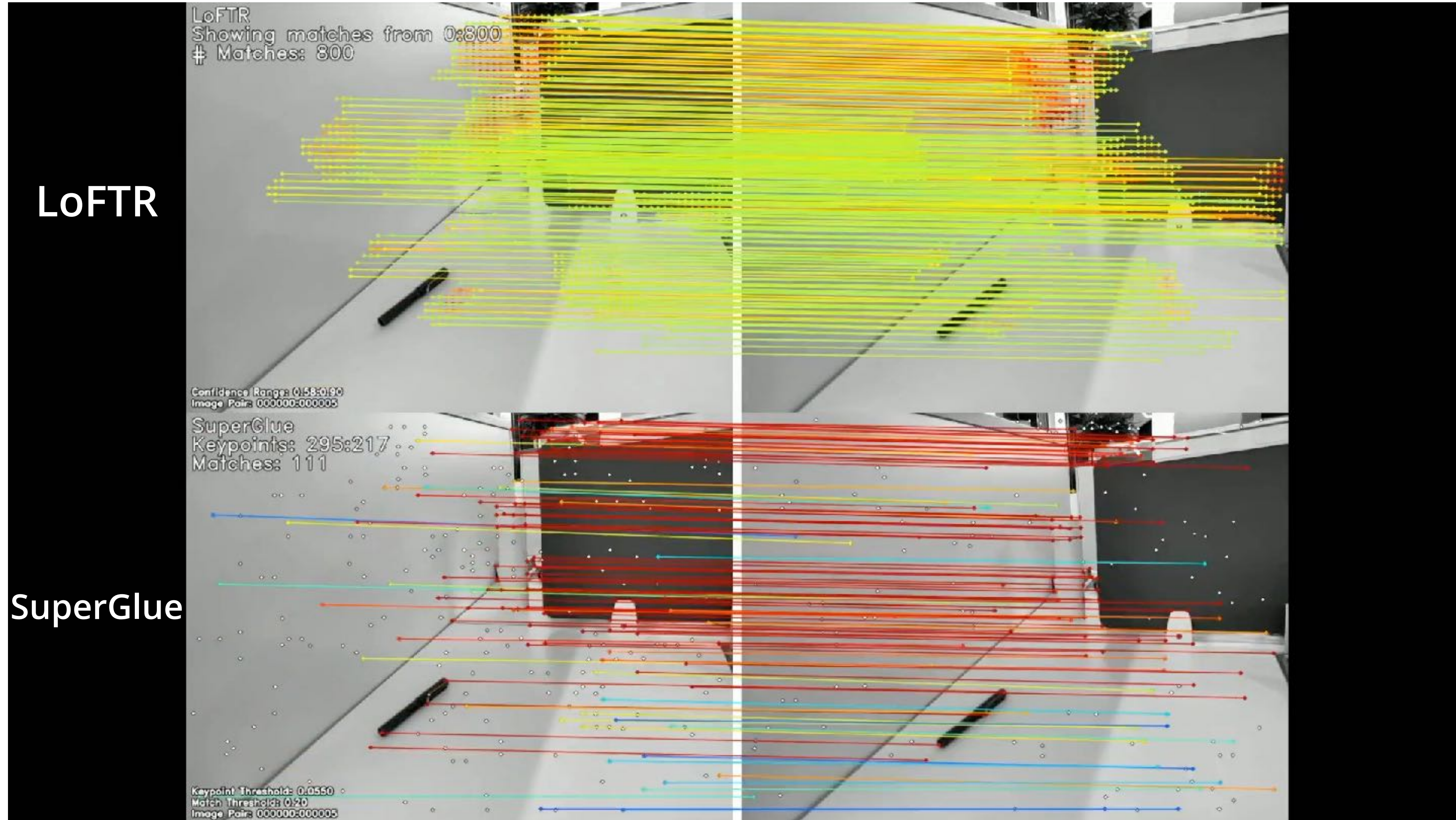
Demo

Matching on a low-texture area



Demo

Comparison with SuperGlue



LoFTR: **Detector-Free** Local Feature Matching with **Transformers**

Jiaming Sun* Zehong Shen* Yu'ang Wang* Hujun Bao Xiaowei Zhou
CVPR 2021

Project page: <https://zju3dv.github.io/loftr>

Code: <https://github.com/zju3dv/LoFTR>

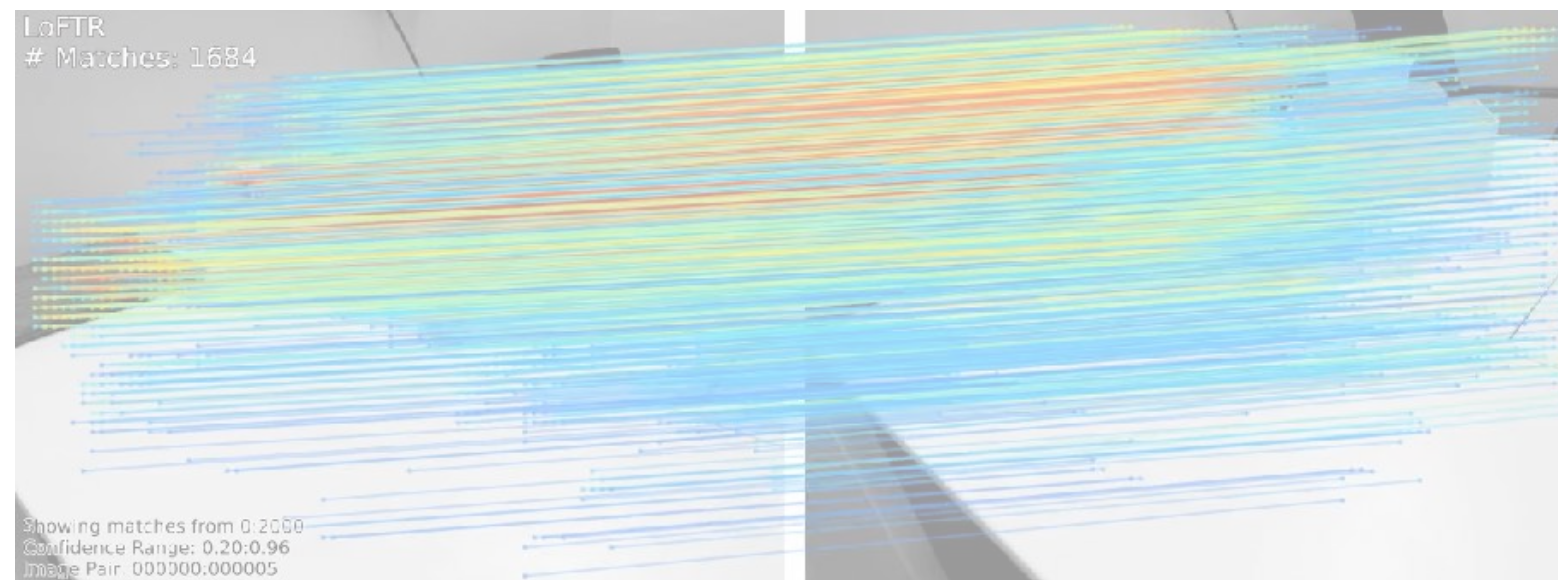
Paper link: <https://arxiv.org/pdf/2104.00680.pdf>

Summary on Localization

- Learning to find correspondences instead of directly regress poses
 - Leverage the **geometric structure** of the problem (BA, PnP, RANSAC)
 - Let CNNs/Transformers do what they do best (learning feature representations)
- Learning feature matching
 - Looking at both image at same time, instead of detecting interests points independently, learns **position-dependent** feature representations
 - Leverage the **inductive bias** of different network architectures (CNNs: local, Transformers: global)
 - Learn the matches **end-to-end** directly

Today's Topic

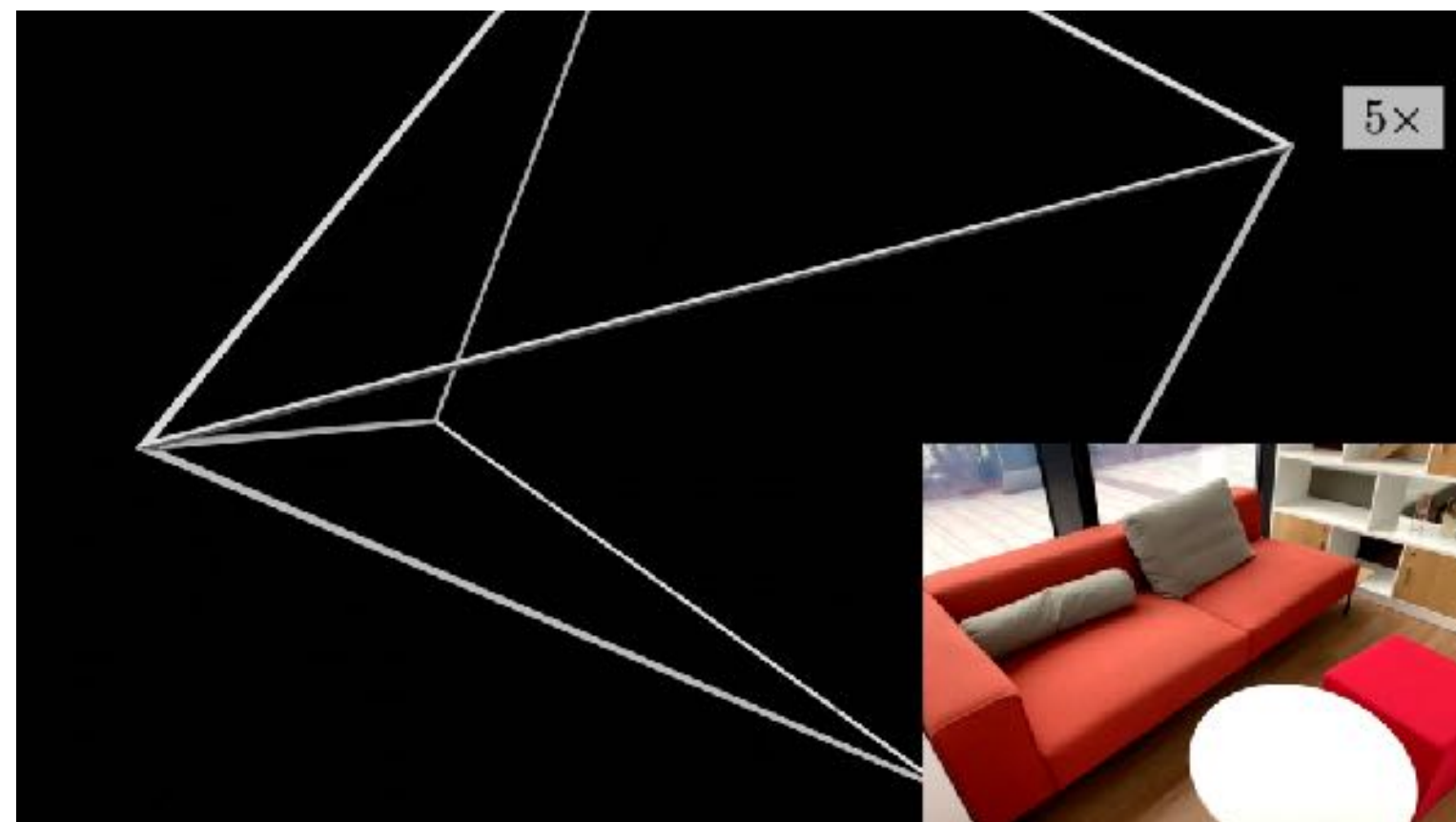
Visual localization



LoFTR: **Detector-Free** Local Feature Matching with **Transformers**

Jiaming Sun* Zehong Shen* Yu'ang Wang* Hujun Bao Xiaowei Zhou
CVPR 2021

3D Reconstruction

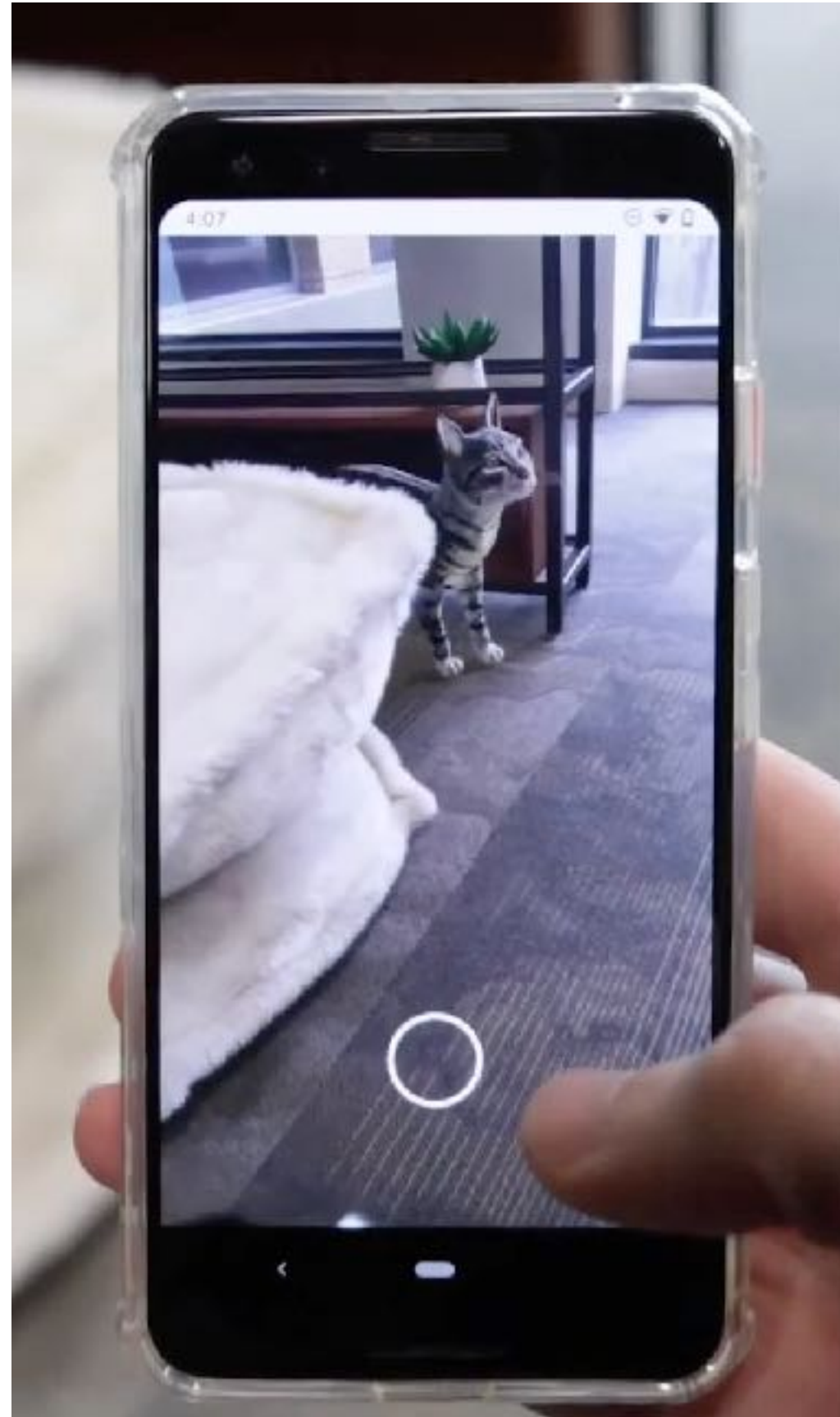


NeuralRecon: **Real-Time Coherent** 3D Reconstruction from Monocular Video

Jiaming Sun* Yiming Xie* Linghao Chen Xiaowei Zhou Hujun Bao
CVPR 2021 (Oral)

Motivation

Real-time 3D reconstruction is crucial for immersive AR effects



Motivation

Depth sensor v.s. Monocular camera

With a depth sensor



- 😊 Accurate depth measurement
- 😞 Takes a lot of energy
- 😞 Only available to a few high-end products

Motivation

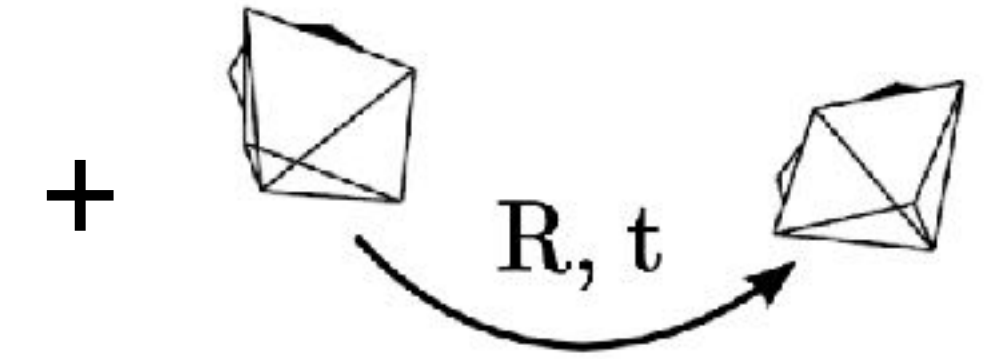
Depth sensor v.s. Monocular camera

With a depth sensor



- 😊 Accurate depth measurement
- 😞 Takes a lot of energy
- 😞 Only available to a few high-end products

With a monocular camera



Motivation

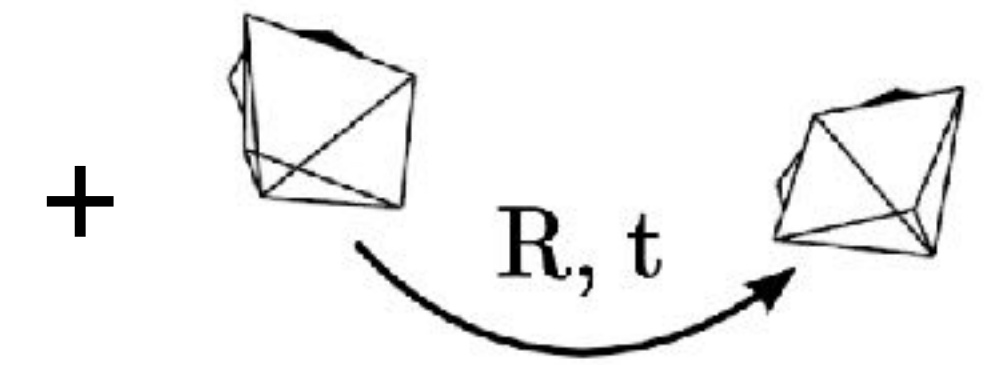
Depth sensor v.s. Monocular camera

With a depth sensor



- 😊 Accurate depth measurement
- 😞 Takes a lot of energy
- 😞 Only available to a few high-end products

With a monocular camera



- 😊 Immediately available to many phones
- 😞 Not as accurate as depth sensors
- 😞 Poor 3D reconstruction quality

Motivation

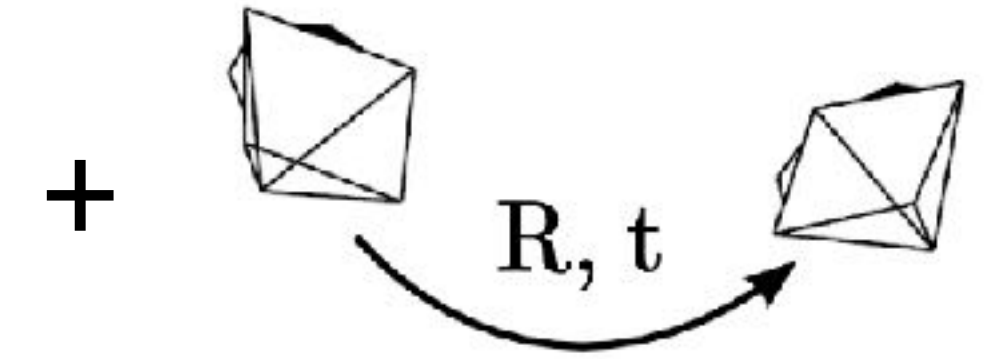
Depth sensor v.s. Monocular camera

With a depth sensor



- 😊 Accurate depth measurement
- 😞 Takes a lot of energy
- 😞 Only available to a few high-end products

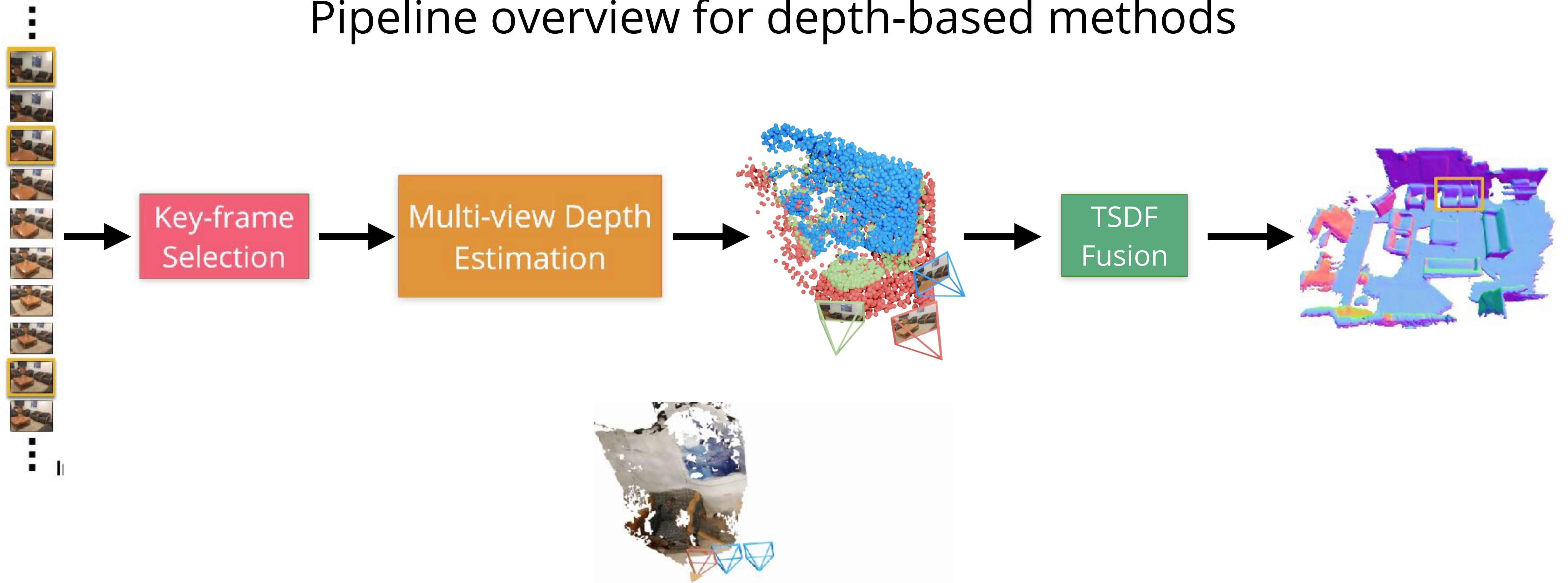
With a monocular camera



- 😊 Immediately available to many phones
- 😞 Not as accurate as depth sensors
- 😞 Poor 3D reconstruction quality

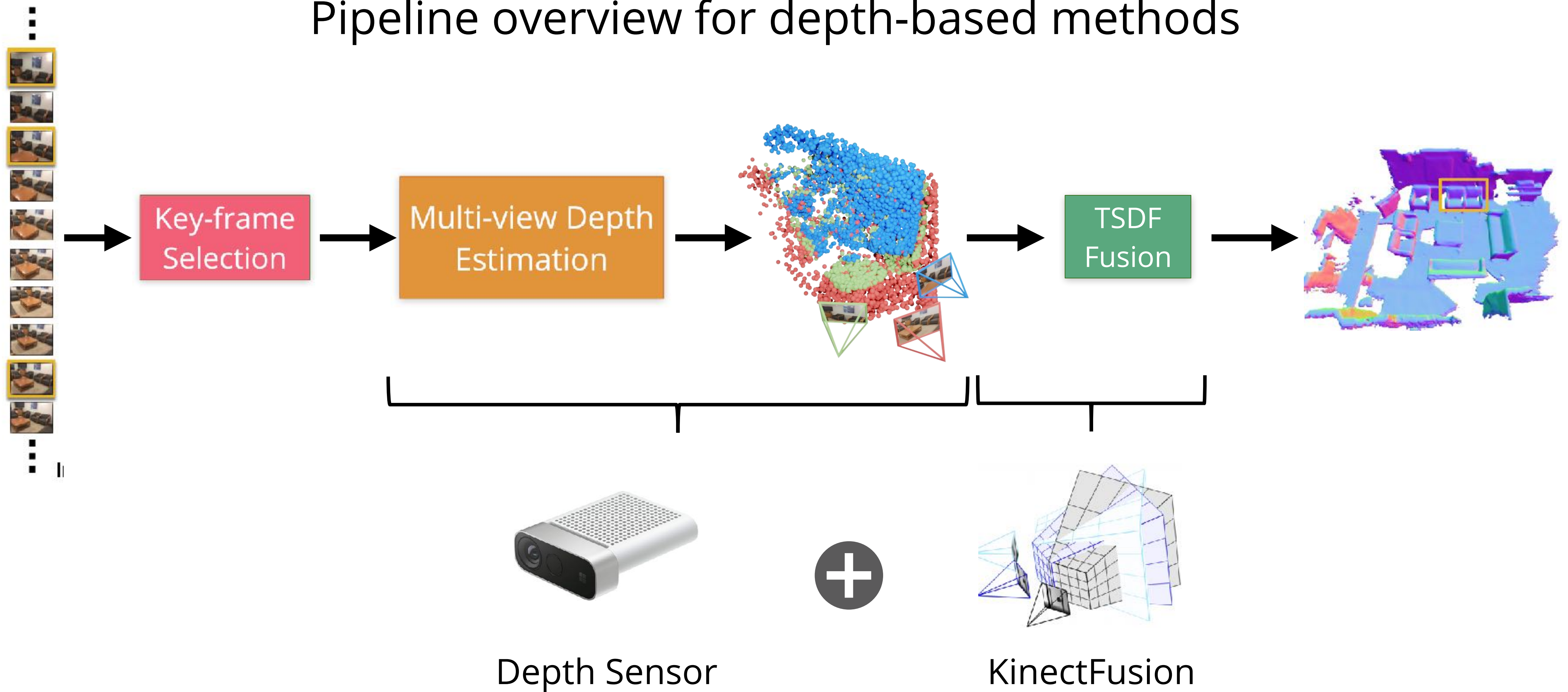
Motivation

Pipeline overview for depth-based methods



Motivation

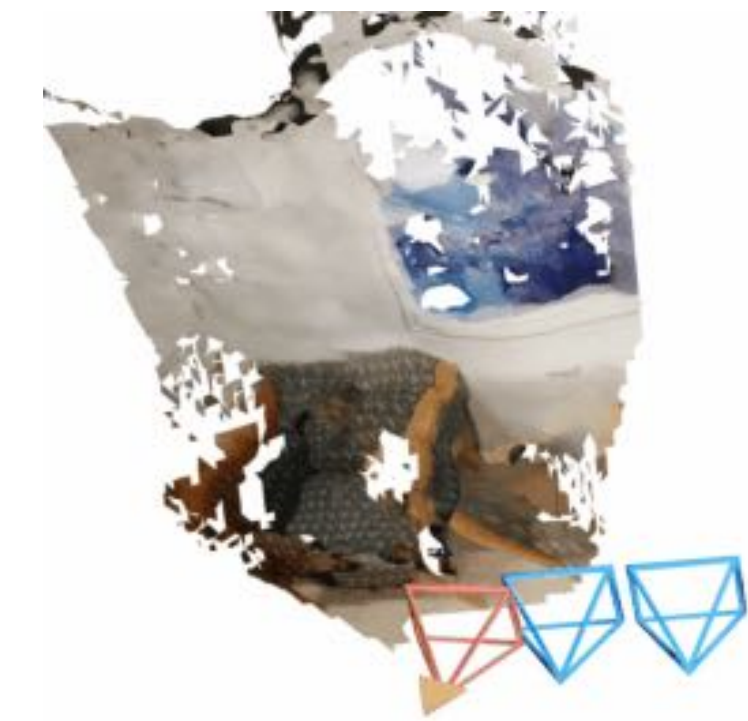
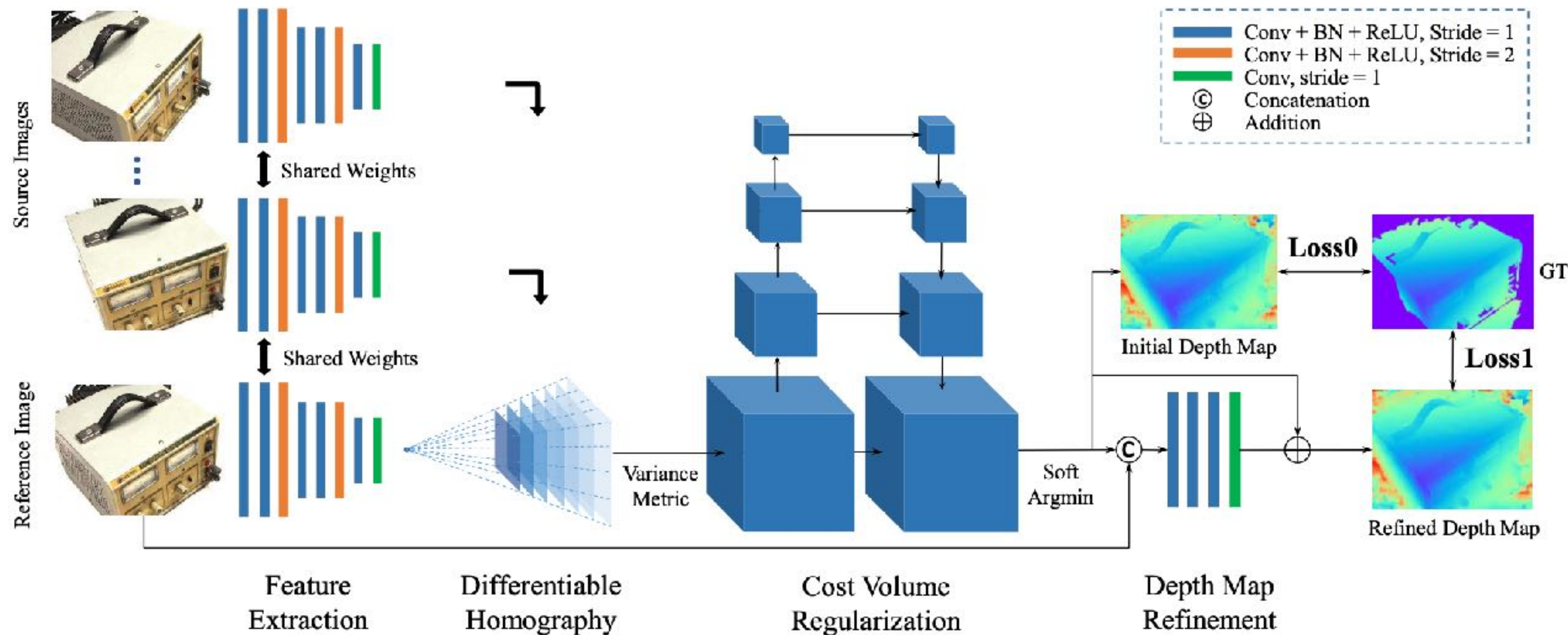
Pipeline overview for depth-based methods



Motivation

MVSNet: Depth Inference for Unstructured Multi-view Stereo

Yao Yao¹, Zixin Luo¹, Shiwei Li¹, Tian Fang², and Long Quan¹



Motivation

DeepTAM: Deep Tracking and Mapping

Huizhong Zhou* Benjamin Ummenhofer* Thomas Brox

DEEPV2D: VIDEO TO DEPTH WITH DIFFERENTIABLE STRUCTURE FROM MOTION

Zachary Teed
Princeton University
zteed@cs.princeton.edu

Jia Deng
Princeton University
jiadeng@cs.princeton.edu

BA-NET: DENSE BUNDLE ADJUSTMENT NETWORKS

Chengzhou Tang
School of Computer Science
Simon Fraser University
chengzhou_tang@sfu.ca

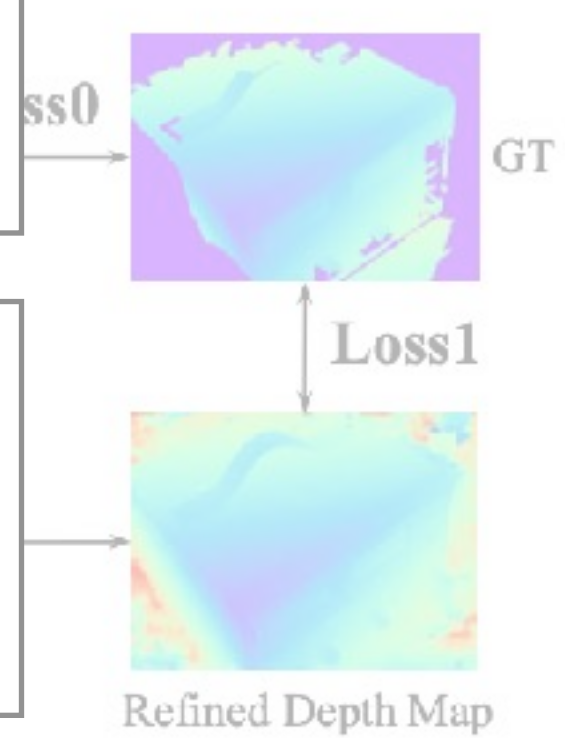
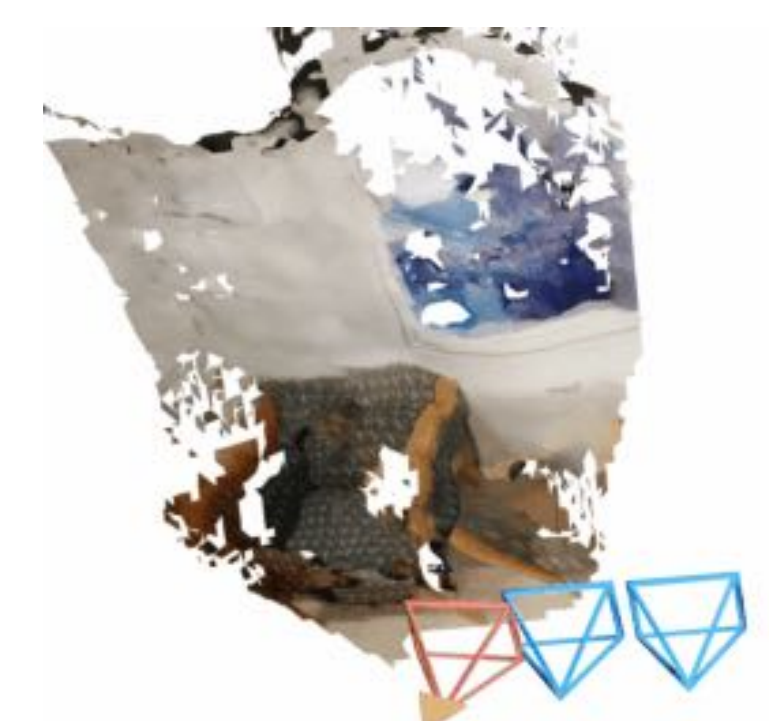
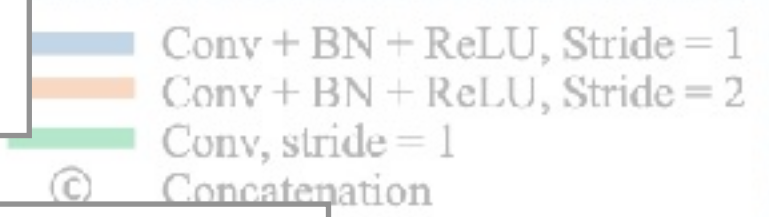
Ping Tan
School of Computer Science
Simon Fraser University
pingtan@sfu.ca

MVDepthNet: Real-time Multiview Depth Estimation Neural Network

Kaixuan Wang Shaojie Shen
Hong Kong University of Science and Technology

Neural RGB→D Sensing: Depth and Uncertainty from a Video Camera

Chao Liu^{1,2*} Jinwei Gu^{1,3*} Kihwan Kim¹ Srinivasa Narasimhan² Jan Kautz¹
¹NVIDIA ²Carnegie Mellon University ³SenseTime



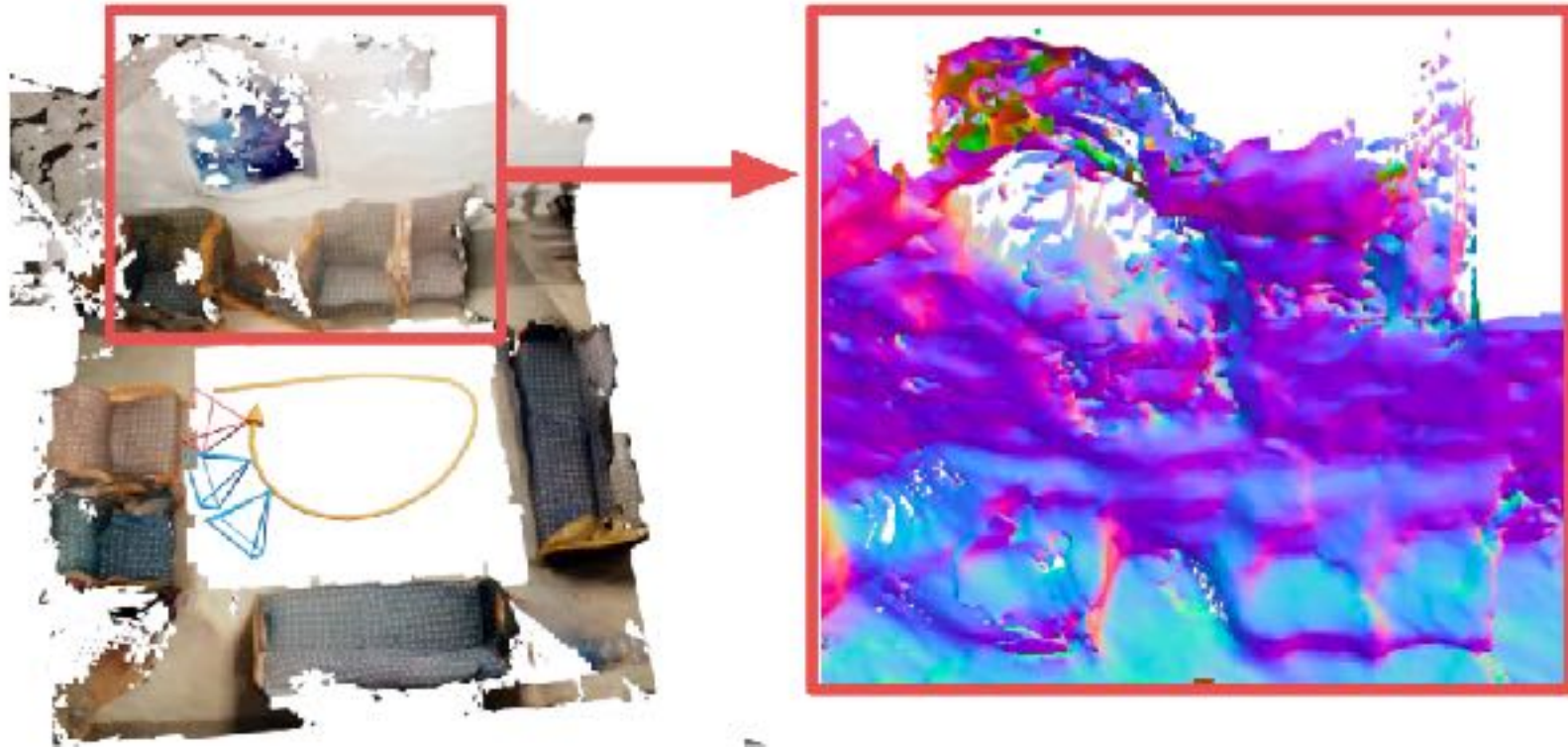
Feature Extraction Differentiable Homography Cost Volume Regularization Depth Map Refinement

Recently: Cascade-Stereo, DeepSFM, CNMNet, Consistent Depth...

Motivation

Depth-based methods v.s. NeuralRecon

Depth-based methods

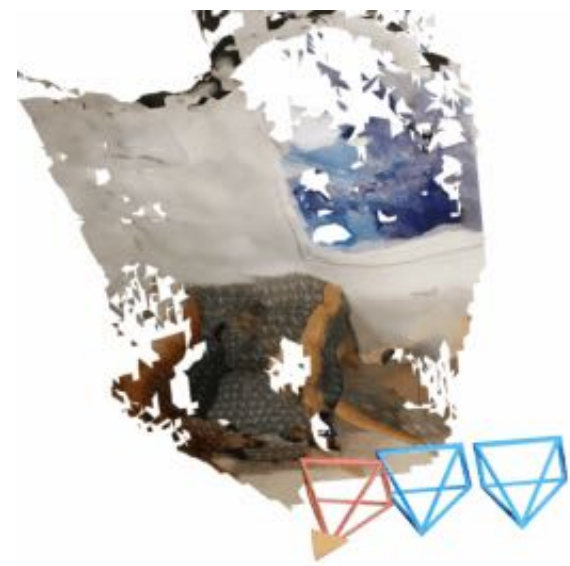


😞 Either layered or scattered results

Motivation

Depth-based methods v.s. NeuralRecon

Depth-based methods



😞 Either layered or scattered results

😞 Redundant computation

Motivation

Depth-based methods v.s. NeuralRecon

Depth-based methods



- 😞 Either layered or scattered results
- 😞 Redundant computation

Our solution: NeuralRecon



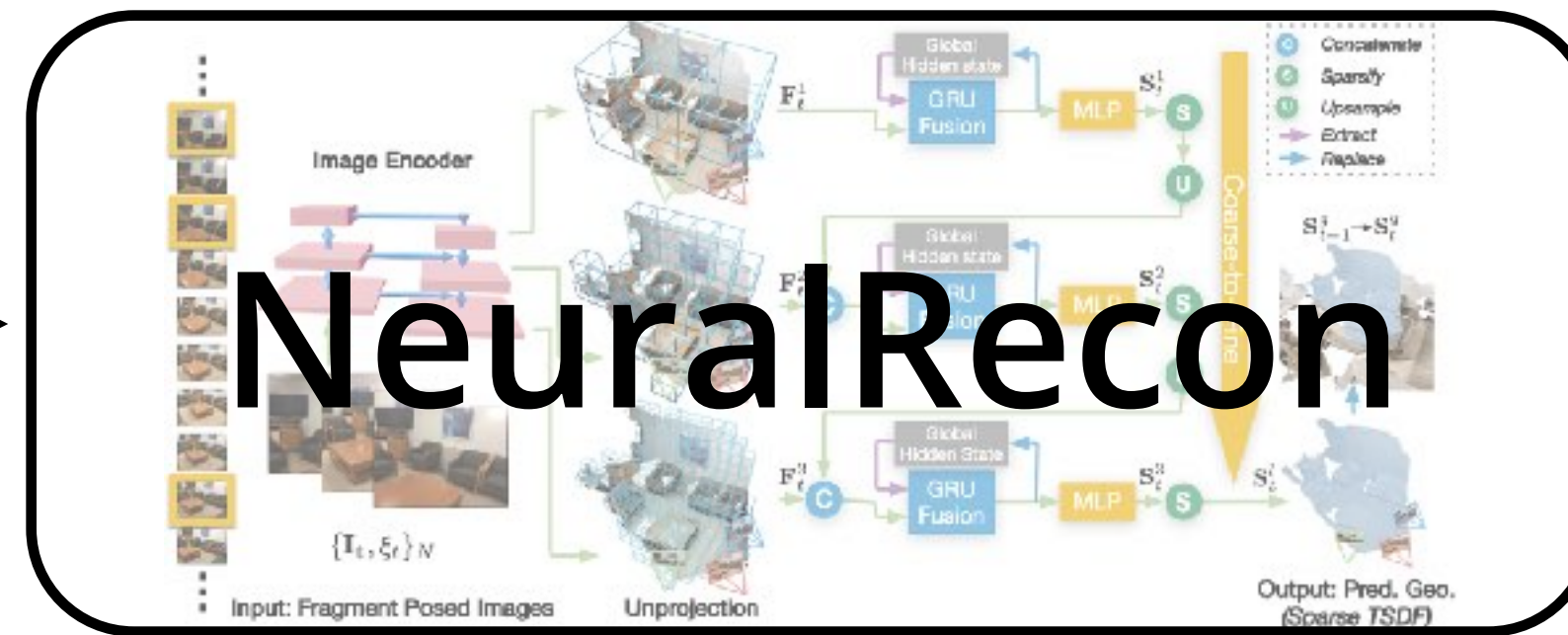
- 😊 Reconstruct local surfaces directly in TSDF
- 😊 Joint fragment reconstruction and fusion
- 😊 Better quality and faster speed

NeuralRecon

Overview



Input:
Posed Images



End-to-End System



Output:
Scene Geometry
(*sparse TSDF volume*)

NeuralRecon

Overview



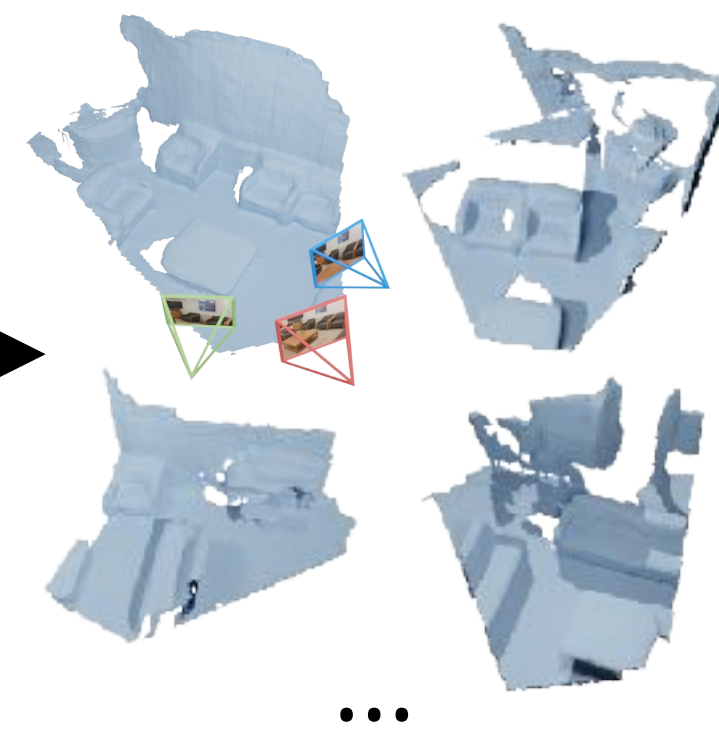
Input:
Posed Images



Key-frame
Selection



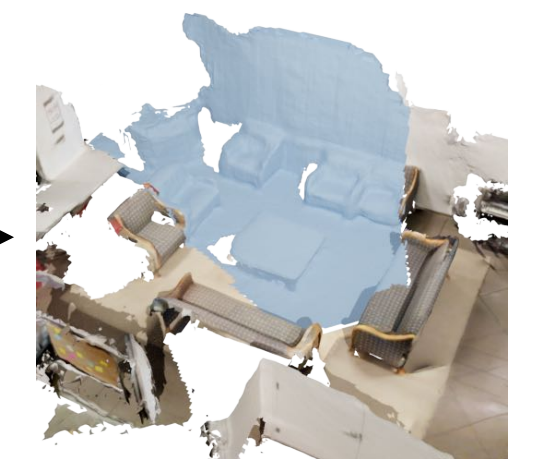
Fragment
Reconstruction



Fragment Geometry
(sparse TSDF volume)



Fuse to Global
Volume

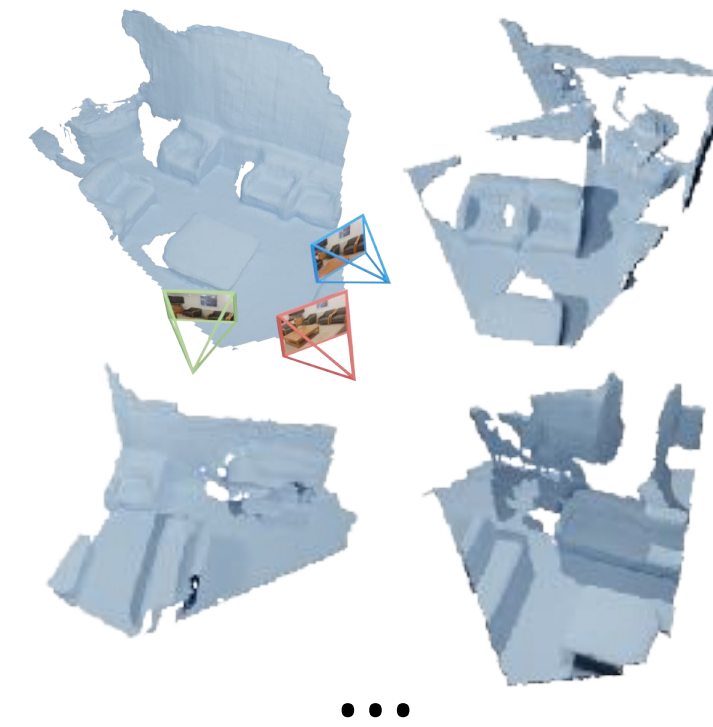


Output:
Scene Geometry
(sparse TSDF volume)

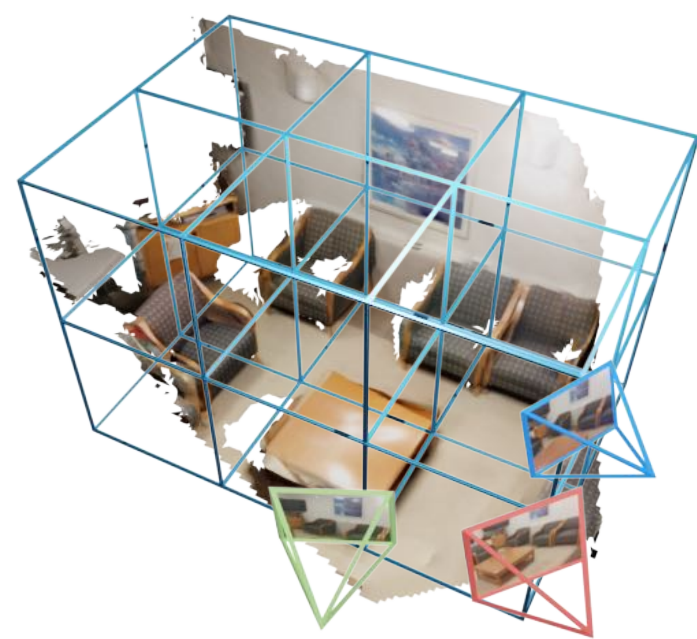
NeuralRecon

Fragment reconstruction overview

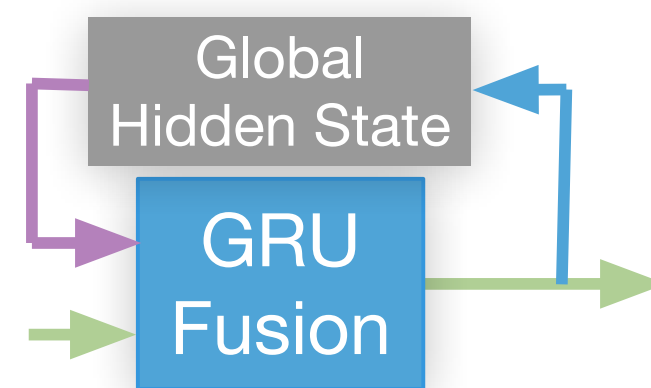
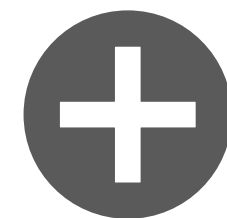
Fragment
Reconstruction



Fragment Geometry
(sparse TSDF volume)



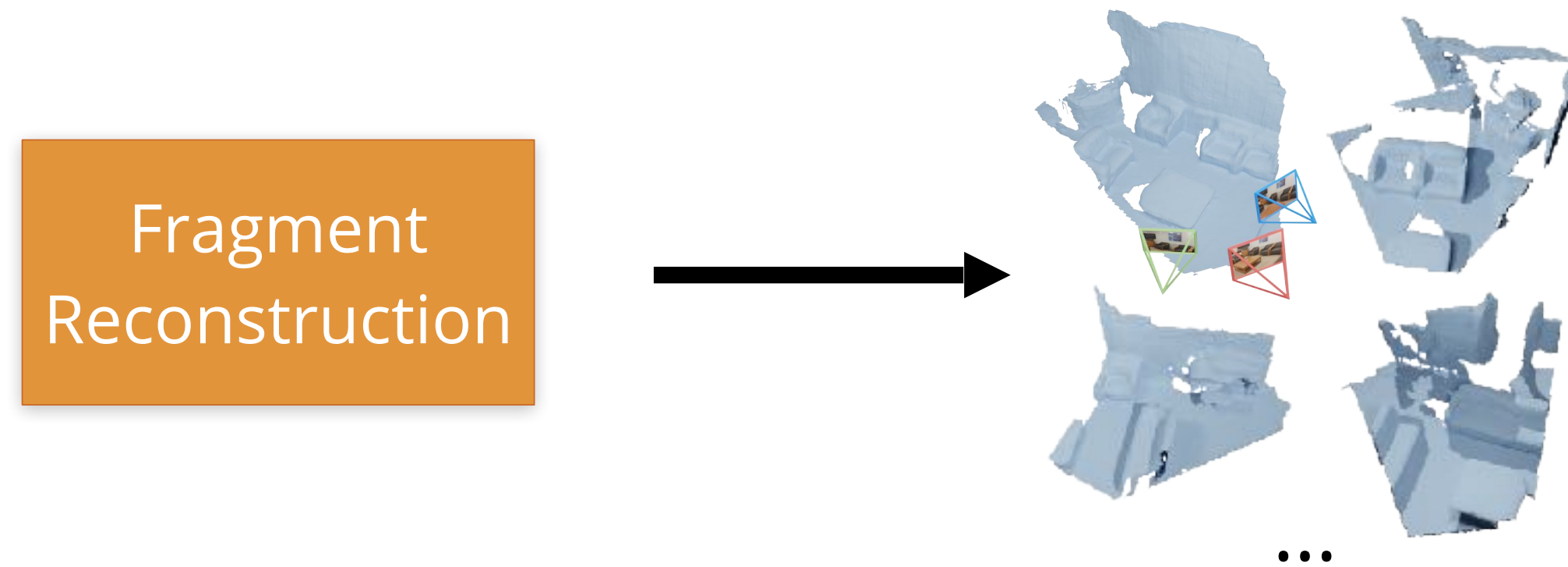
View-Independent
Volume



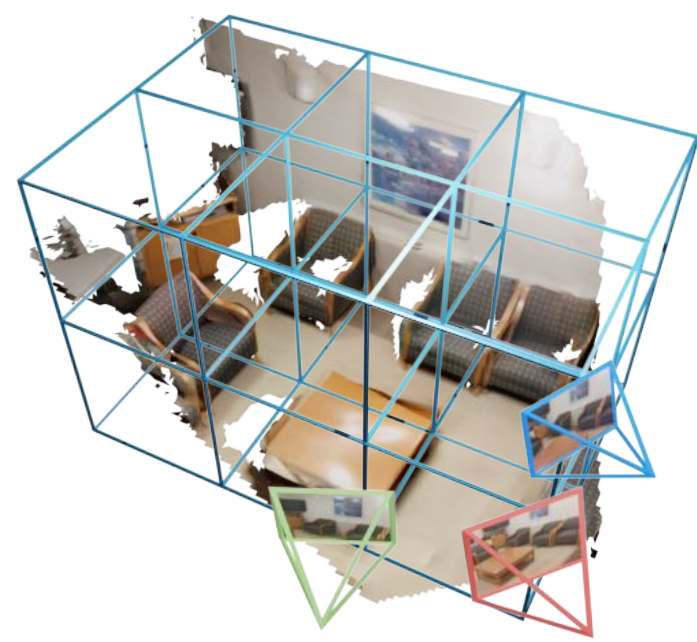
Joint TSDF Reconstruction
and Fusion

NeuralRecon

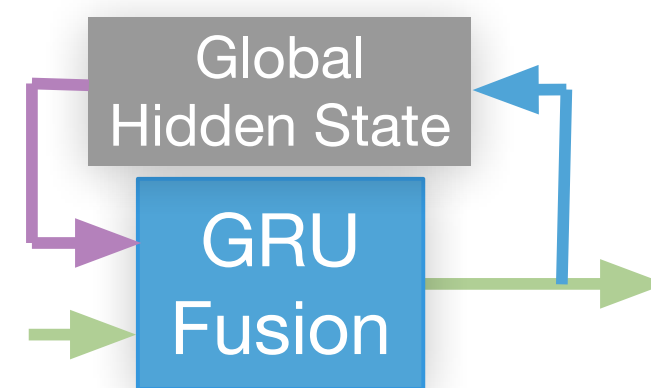
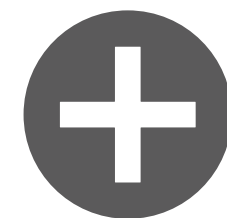
Fragment reconstruction overview



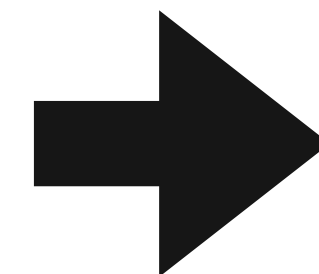
Fragment Geometry
(*sparse TSDF volume*)



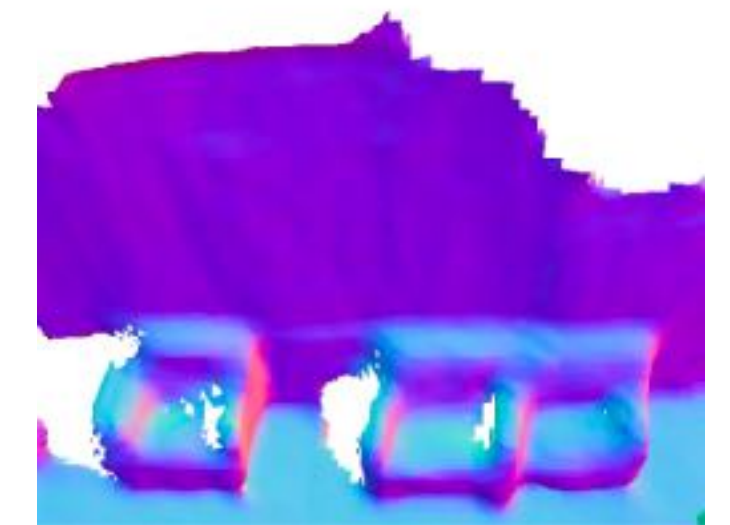
View-Independent
Volume



Joint TSDF Reconstruction
and Fusion



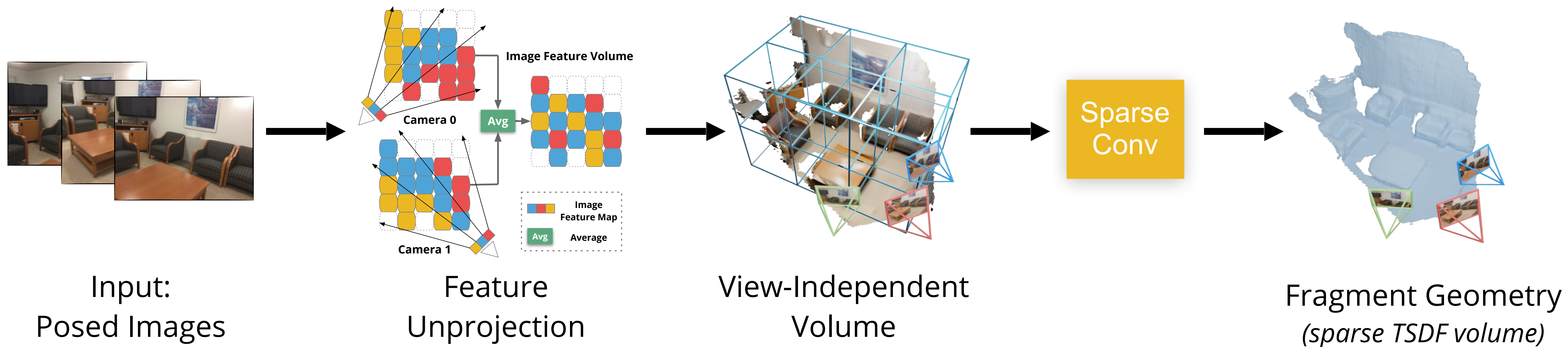
Real-time



Coherent

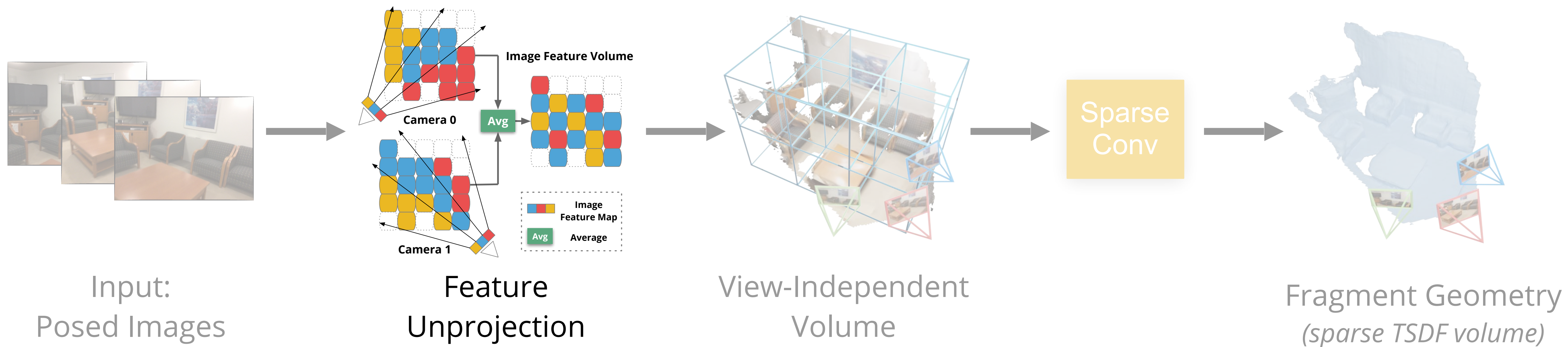
NeuralRecon

Fragment reconstruction



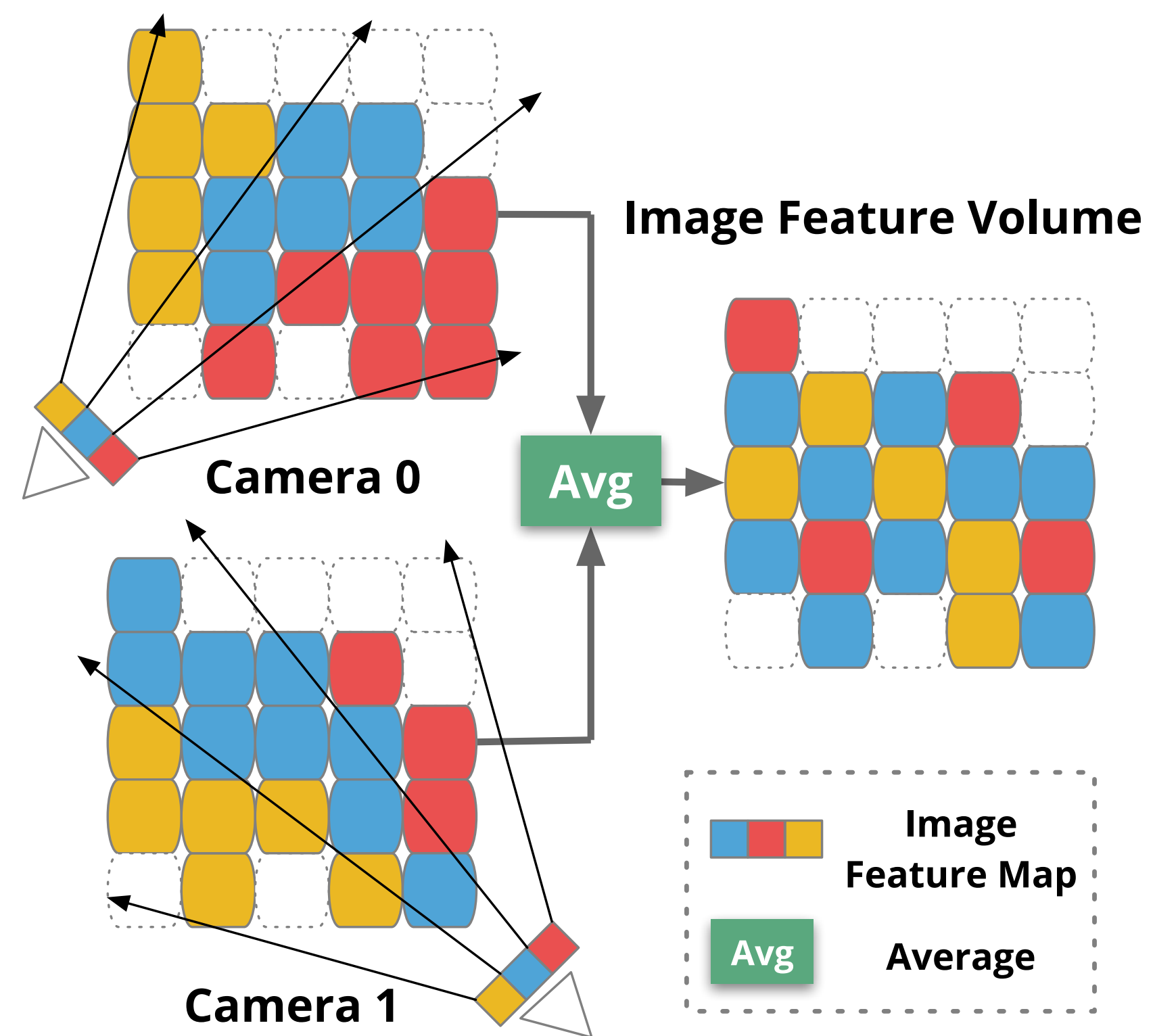
NeuralRecon

Fragment reconstruction



NeuralRecon

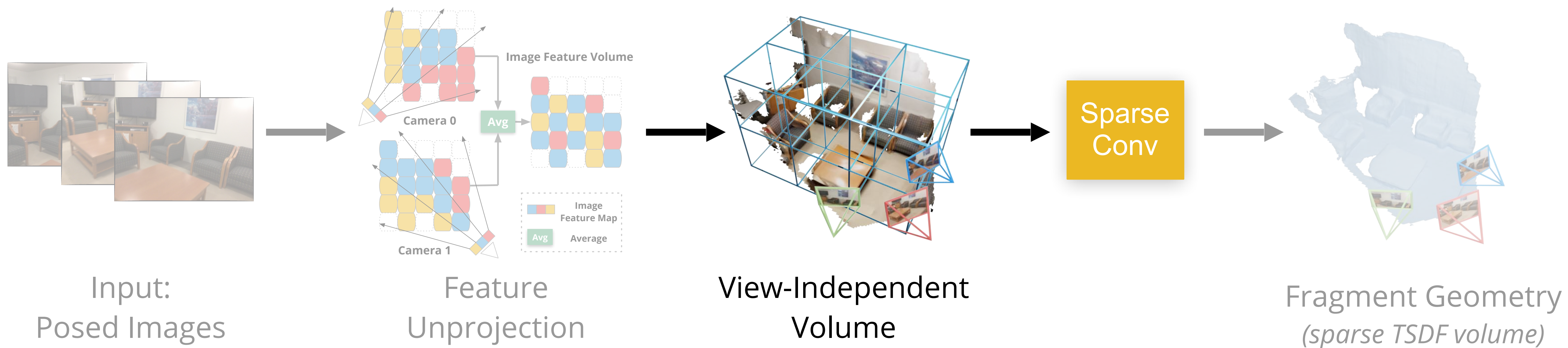
Fragment reconstruction



Feature
Unprojection

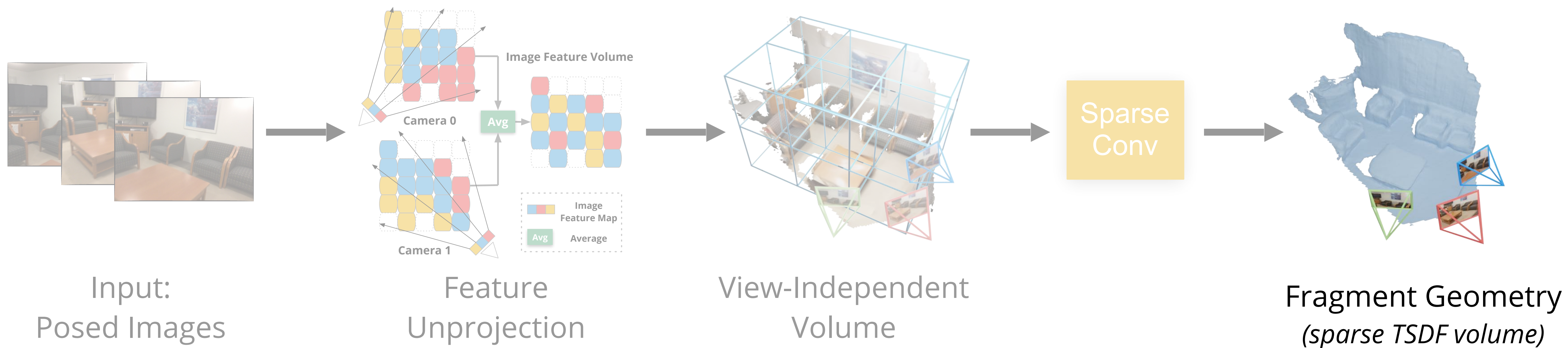
NeuralRecon

Fragment reconstruction



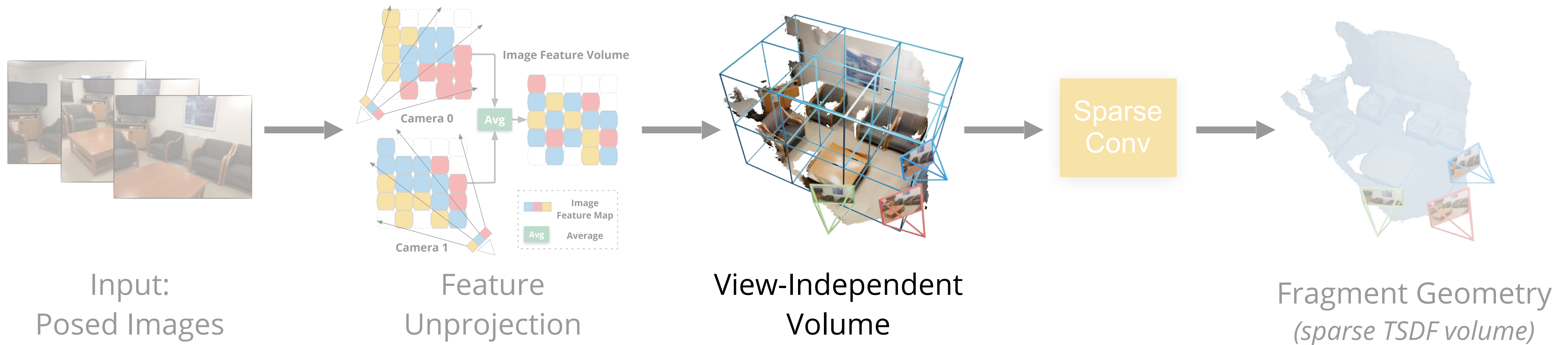
NeuralRecon

Fragment reconstruction



NeuralRecon

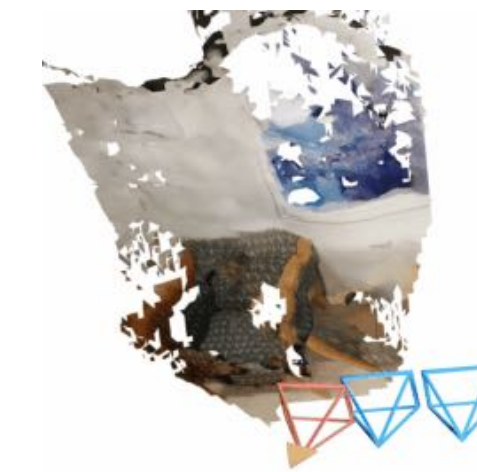
Fragment reconstruction



Why is it better?



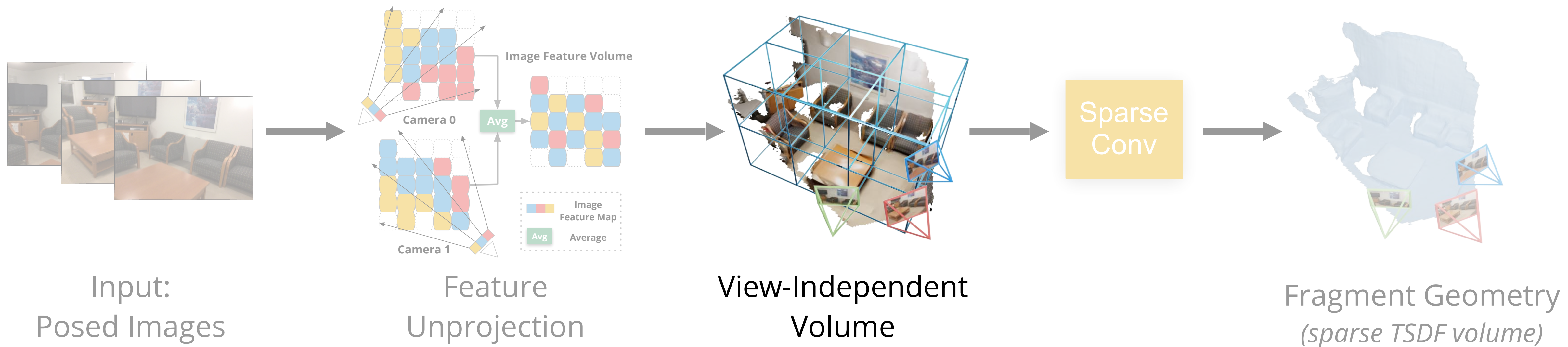
Volume-based



Depth-based

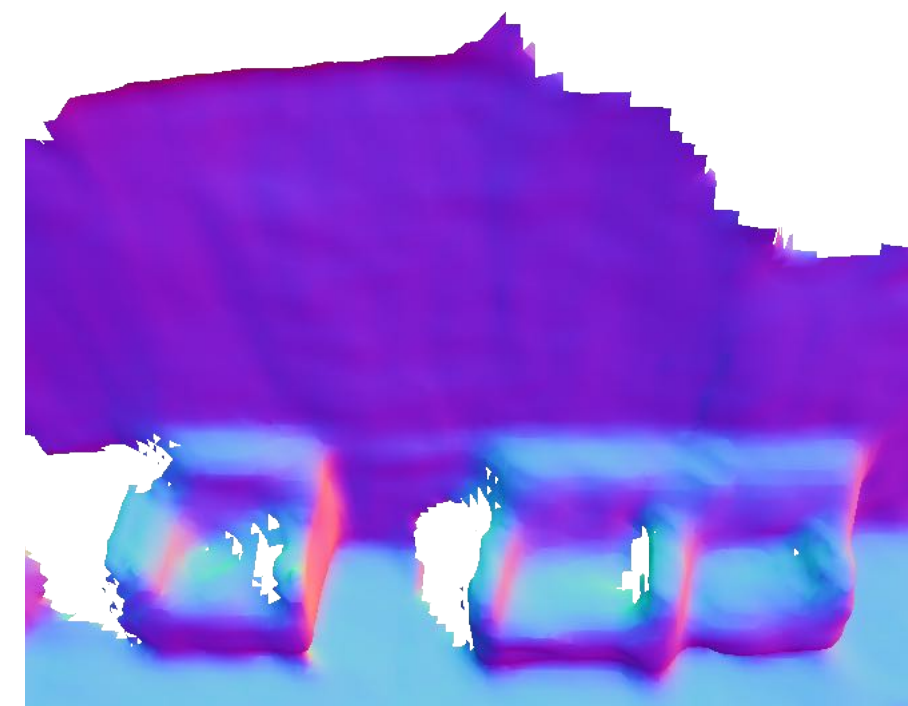
NeuralRecon

Fragment reconstruction

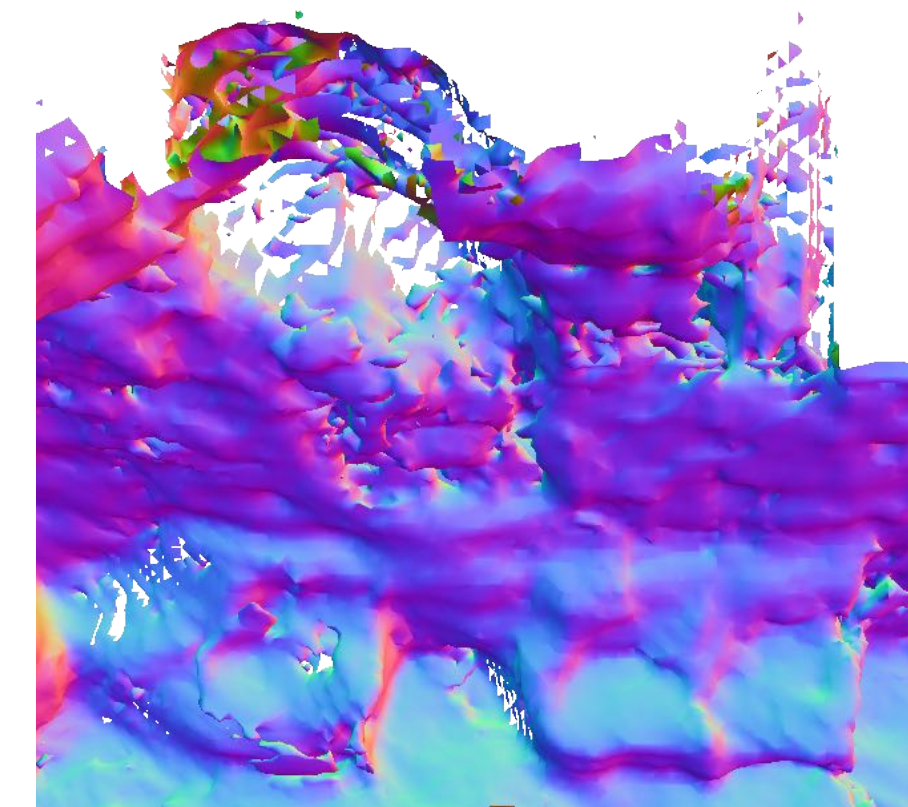


Why is it better?

1. Directly predicts the TSDF rather than fusing single-view depth maps
==> *learns the shape prior of 3D surfaces, produces locally coherent geometry*



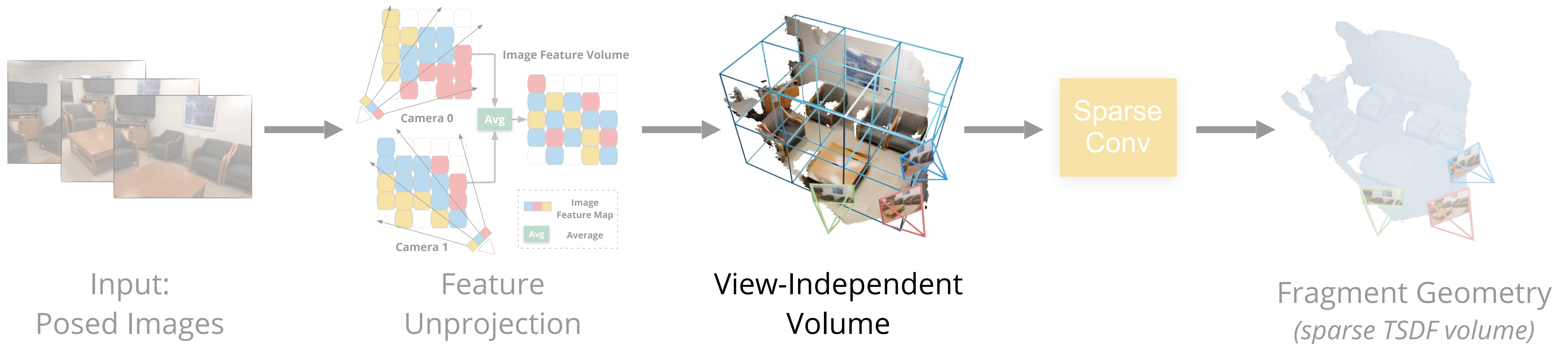
Volume-based



Depth-based

NeuralRecon

Fragment reconstruction

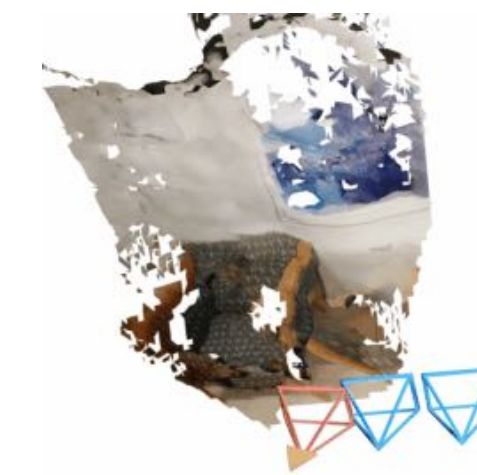


Why is it better?

1. Directly predicts the TSDF rather than fusing single-view depth maps
==> *learns the shape prior of 3D surfaces, produces locally coherent geometry*
2. View-independent volume
==> *reduces redundant computation, faster*



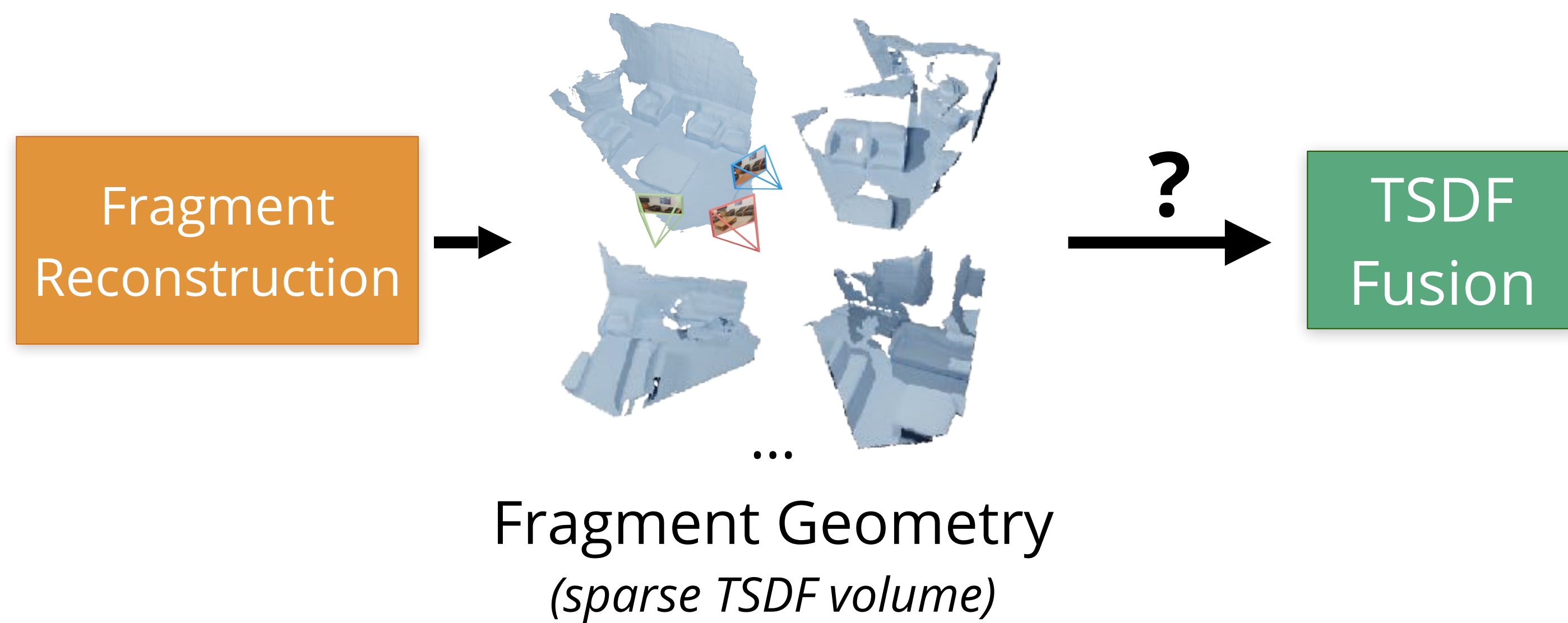
Volume-based



Depth-based

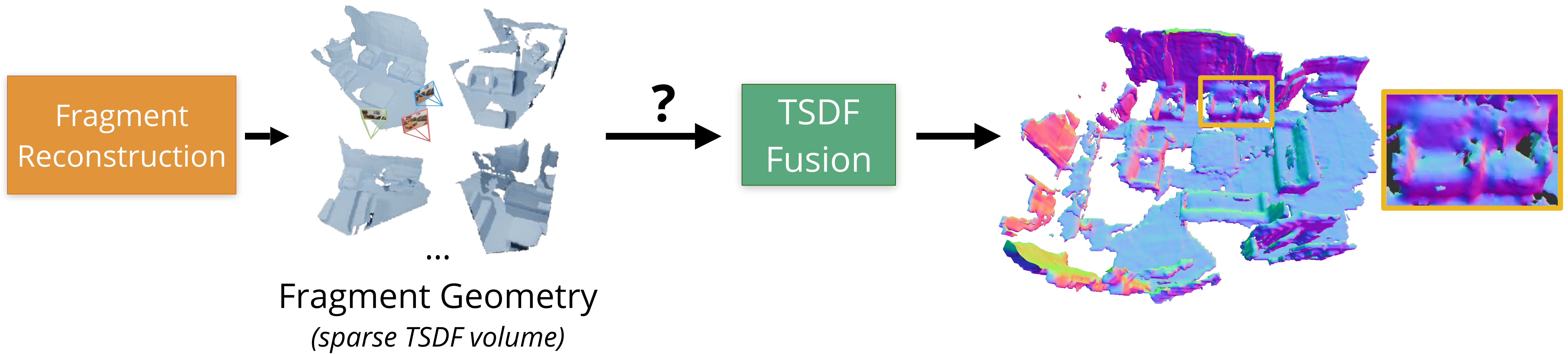
NeuralRecon

TSDF Fusion?



NeuralRecon

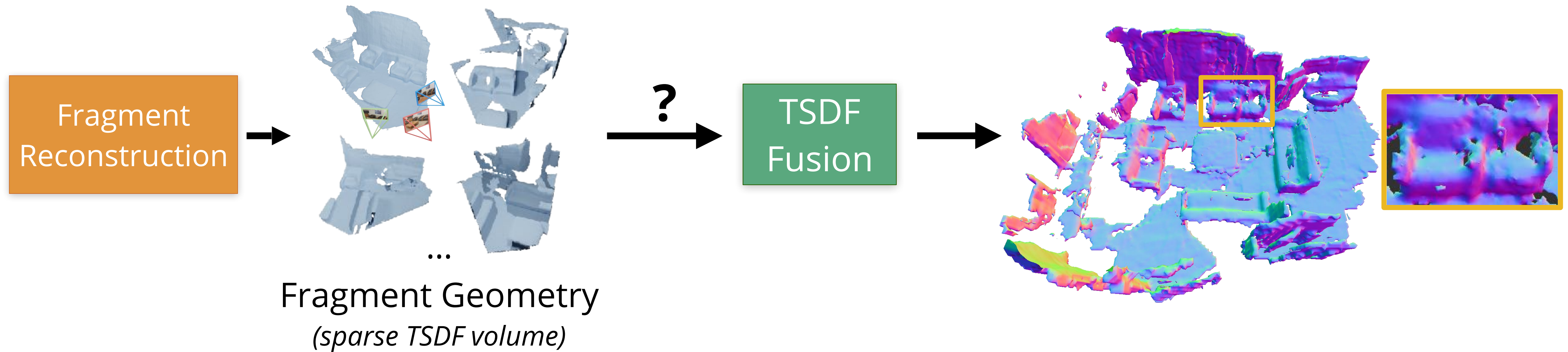
TSDF Fusion?



☹️ Too many artifacts!

NeuralRecon

TSDF Fusion?



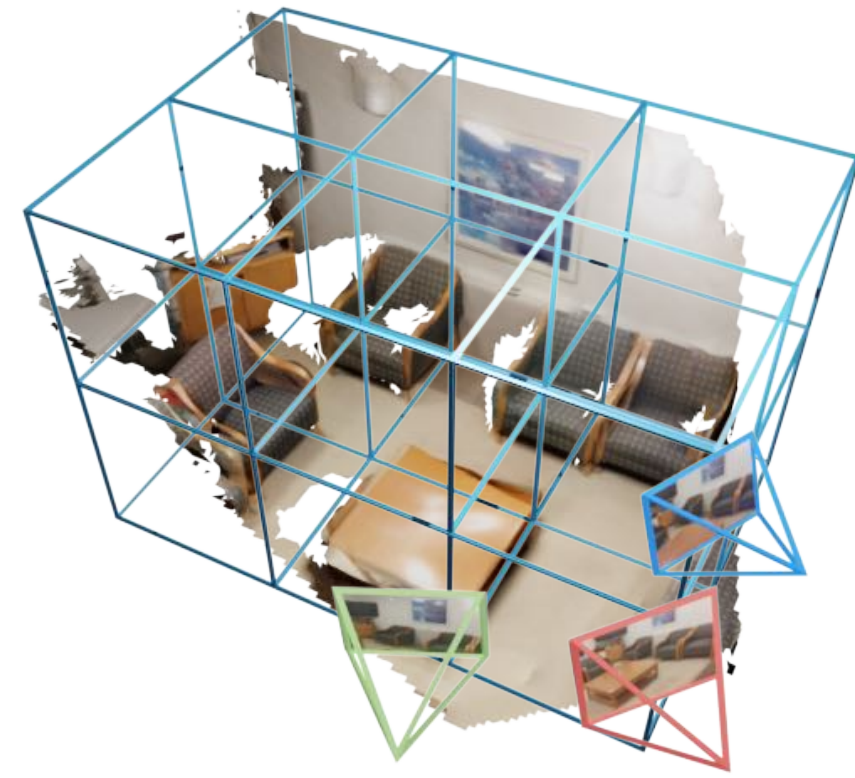
Reason: predicted TSDFs are not consistent between fragments

NeuralRecon

Joint reconstruction and fusion



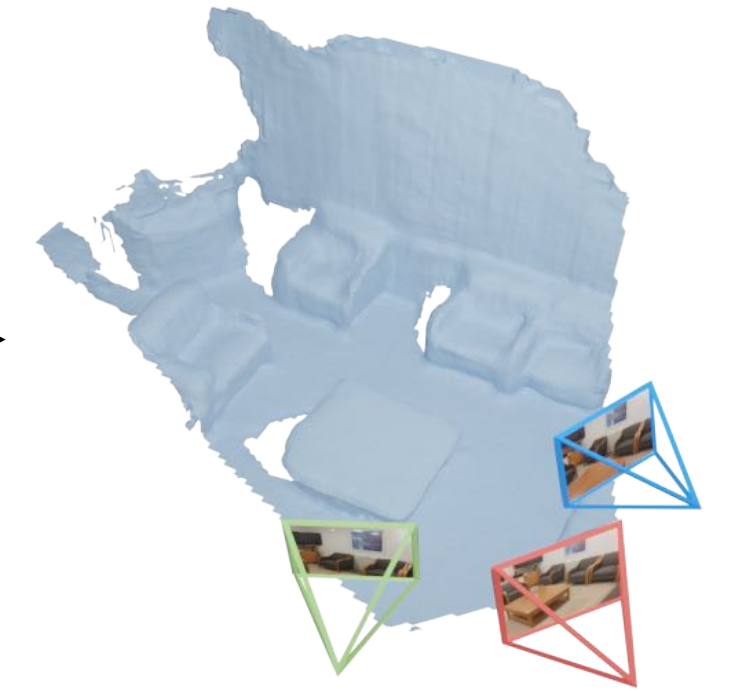
Input: Posed Images



View-Independent
Volume



Sparse
Conv



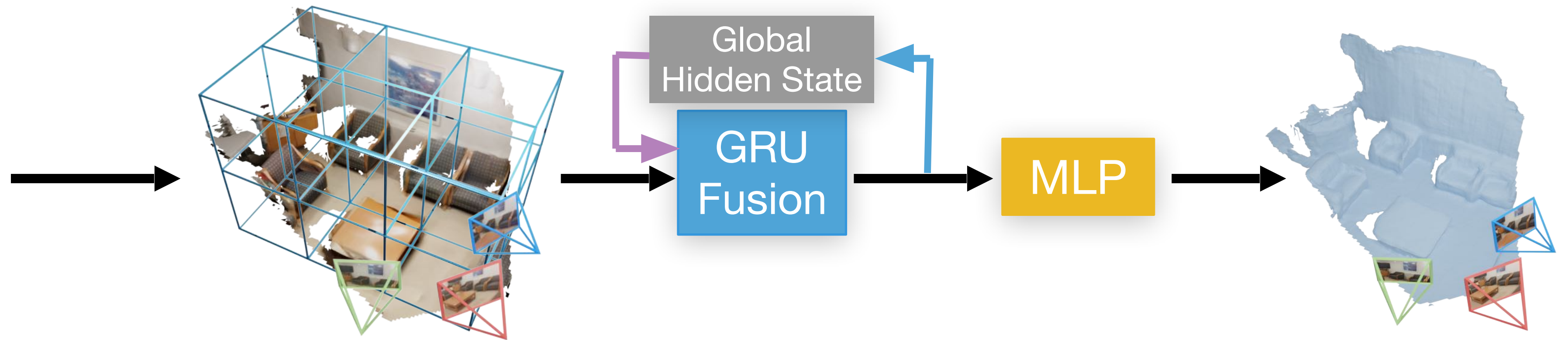
Fragment Geometry
(sparse TSDF volume)

NeuralRecon

Joint reconstruction and fusion



Input: Posed Images

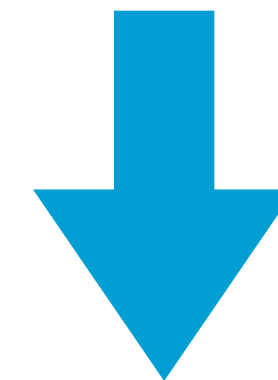
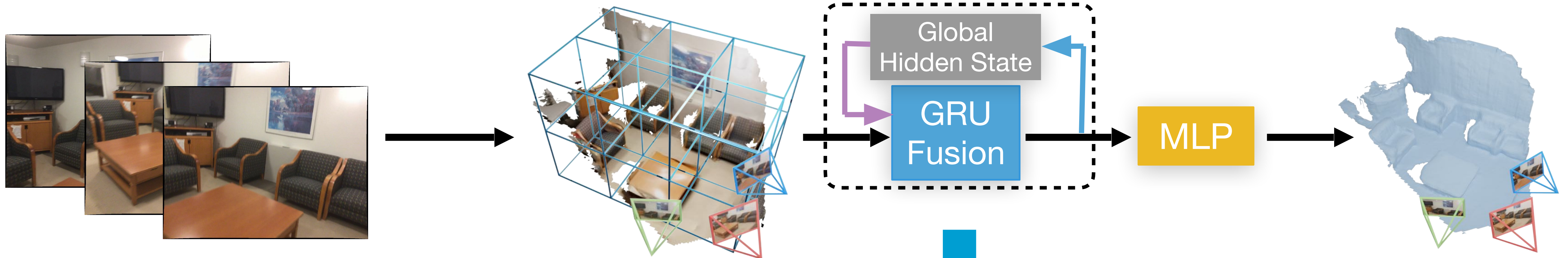


View-Independent
Volume

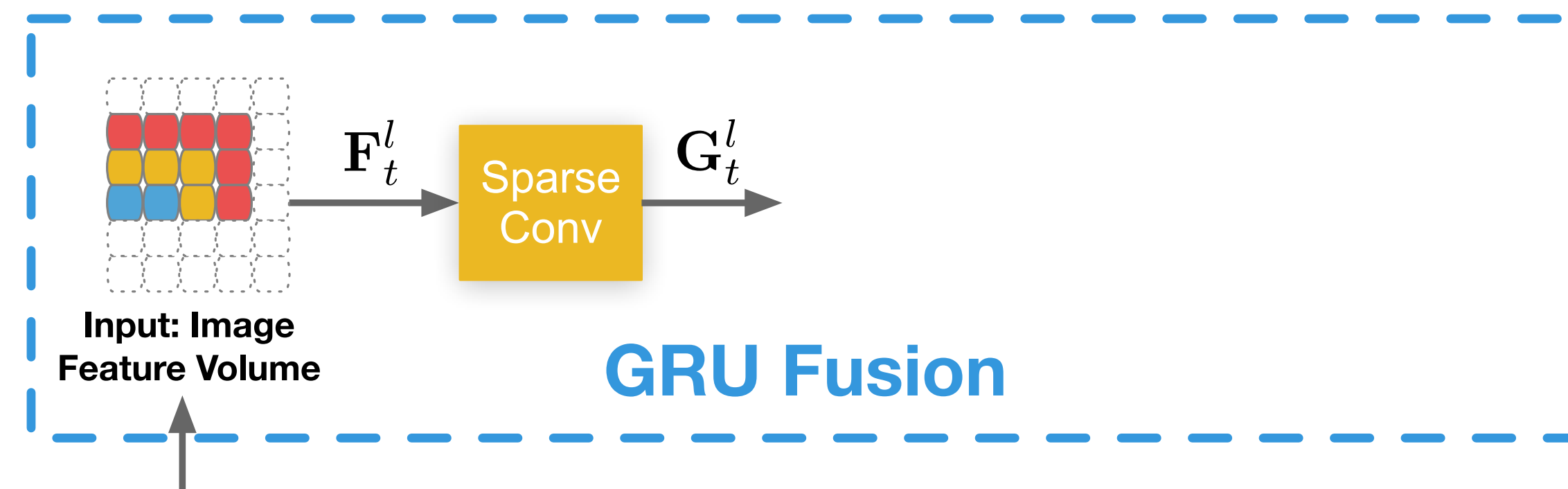
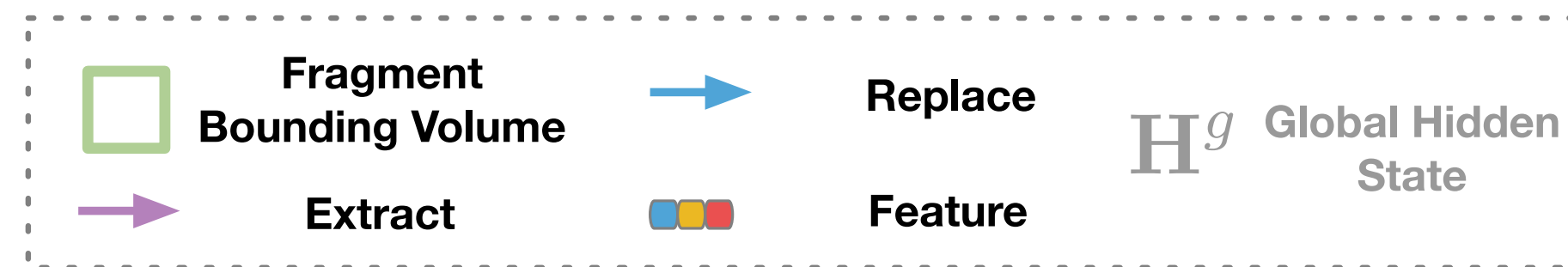
Fragment Geometry
(sparse TSDF volume)

NeuralRecon

Joint reconstruction and fusion

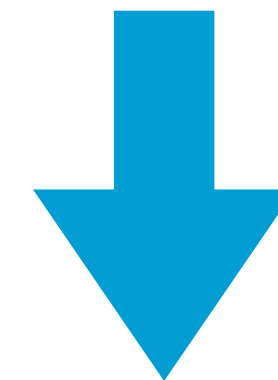
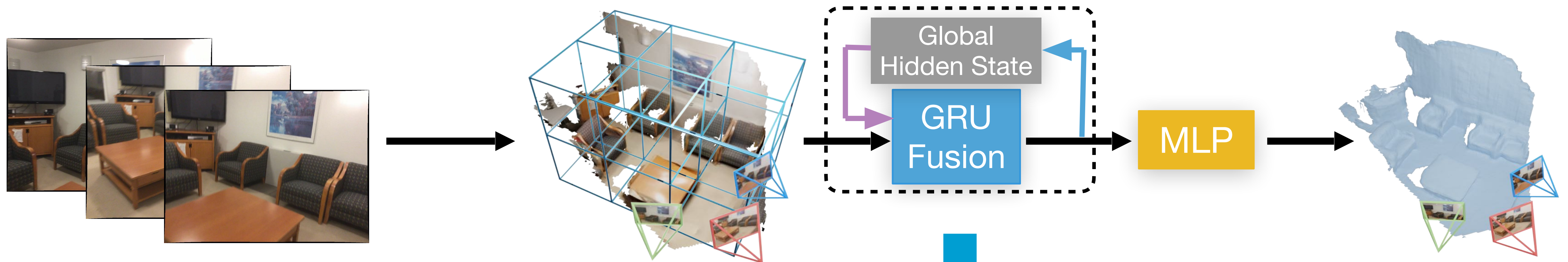


Directly fusing the features with GRU

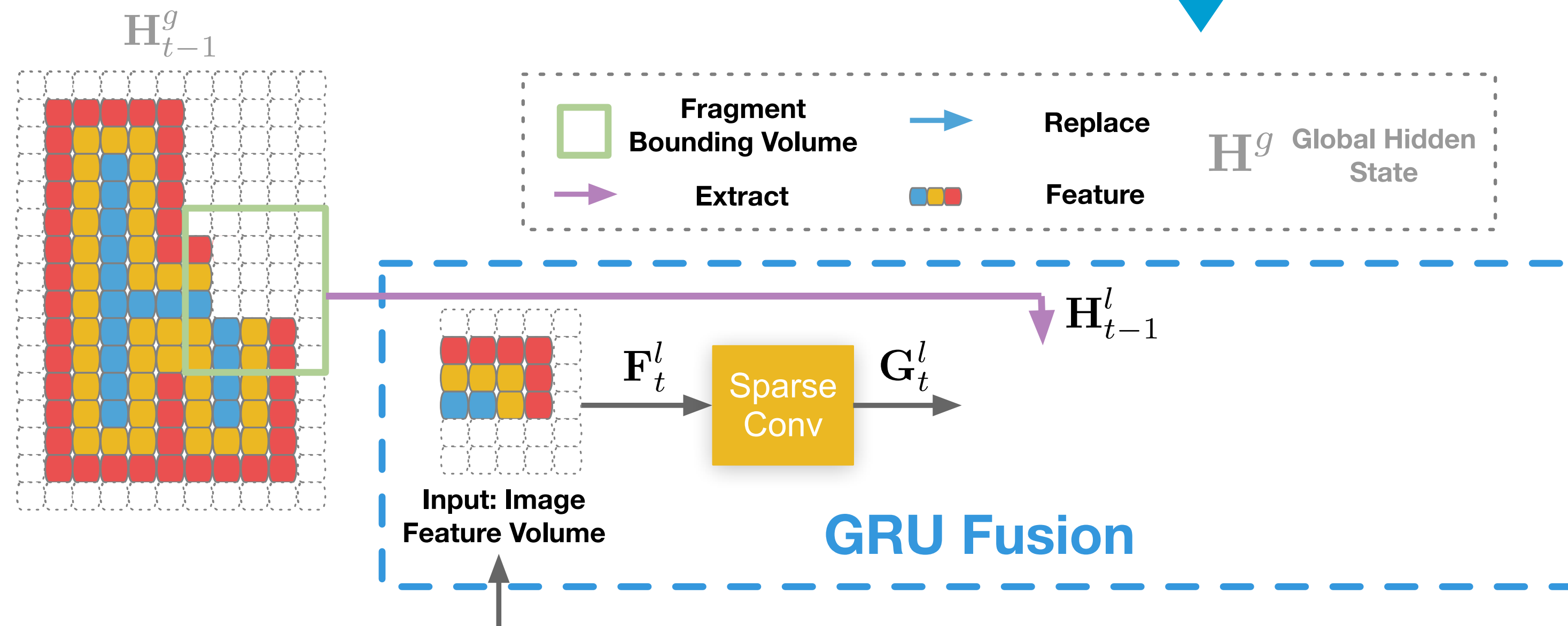


NeuralRecon

Joint reconstruction and fusion

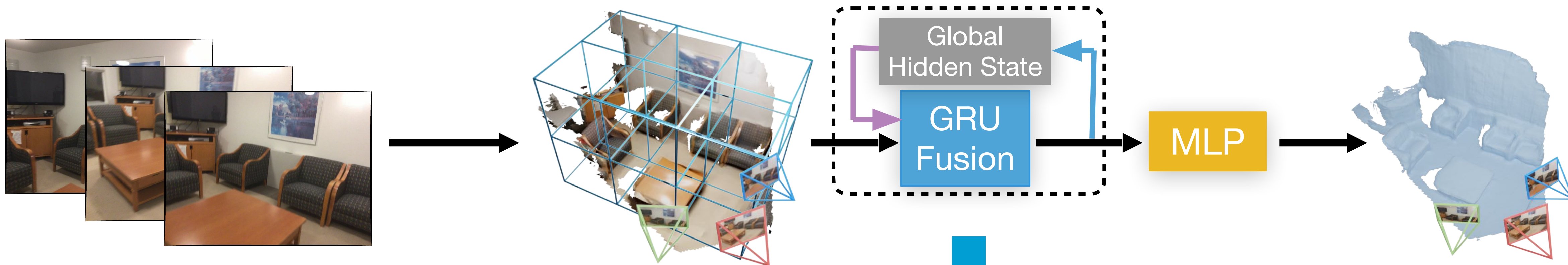


Directly fusing the features with GRU

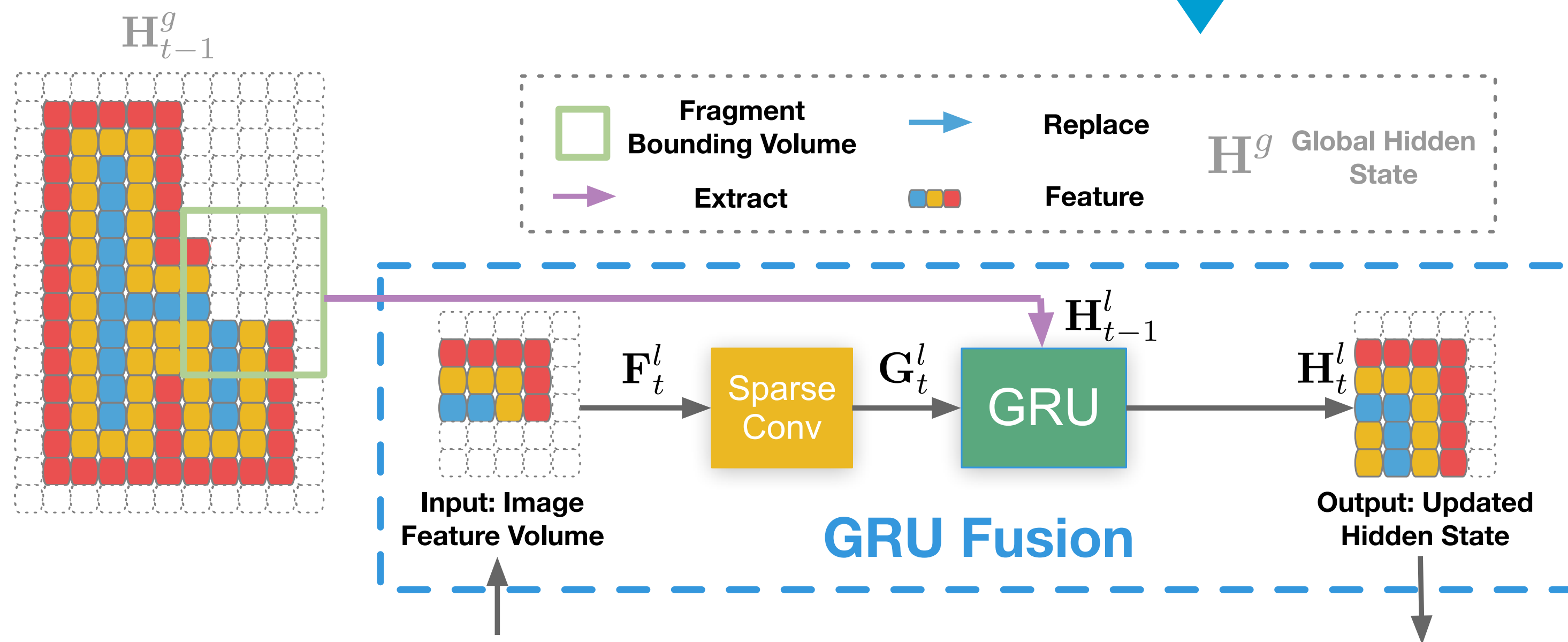


NeuralRecon

Joint reconstruction and fusion

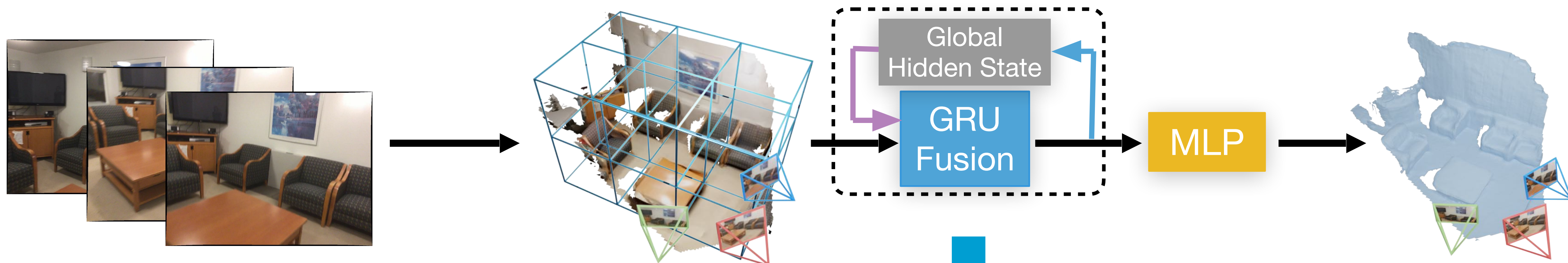


Directly fusing the features with GRU

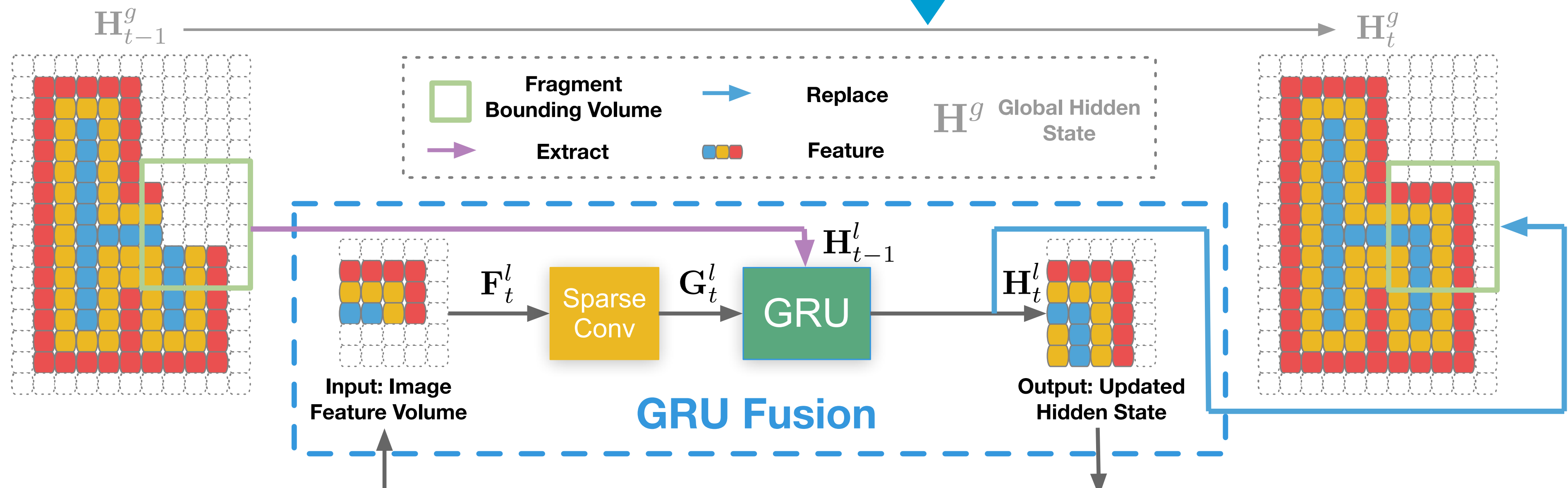


NeuralRecon

Joint reconstruction and fusion

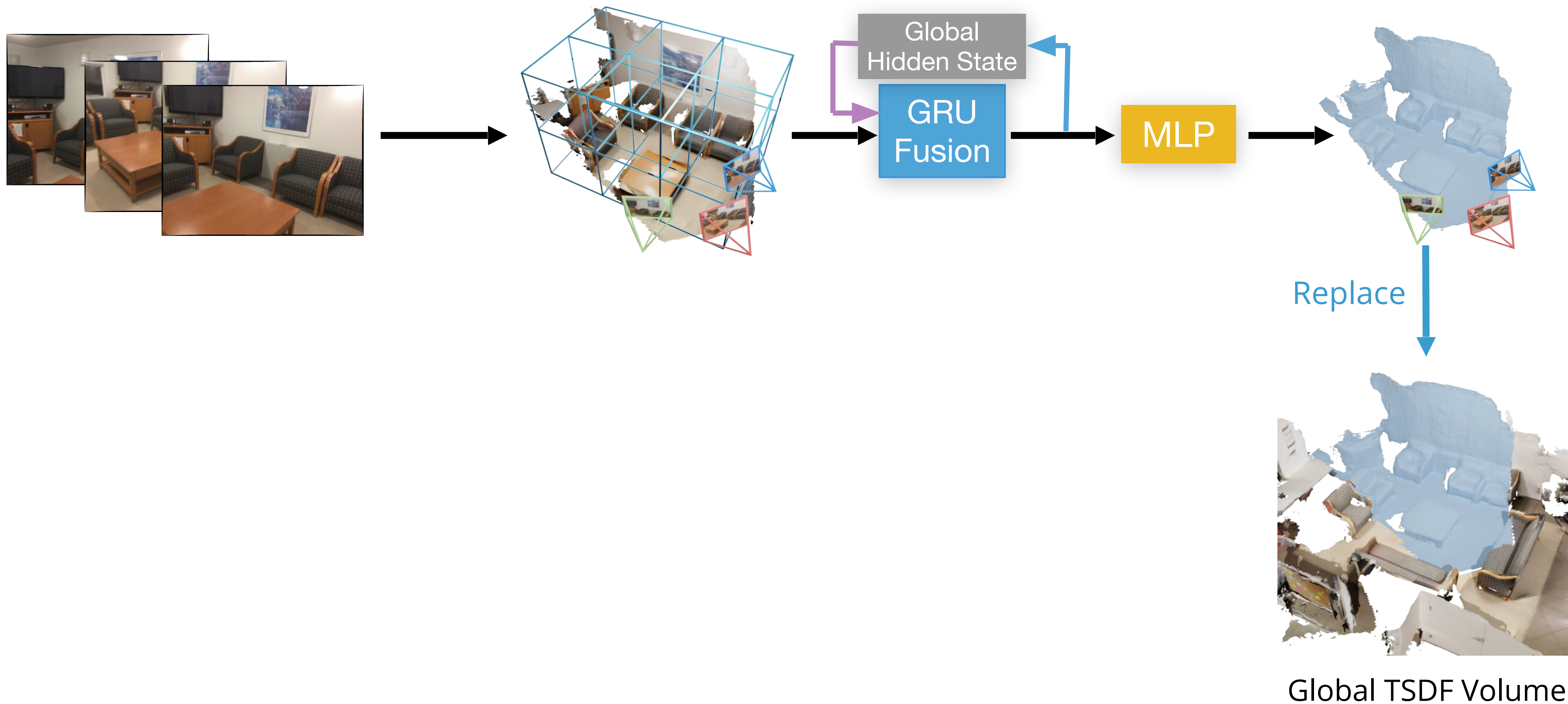


Directly fusing the features with GRU



NeuralRecon

Joint reconstruction and fusion



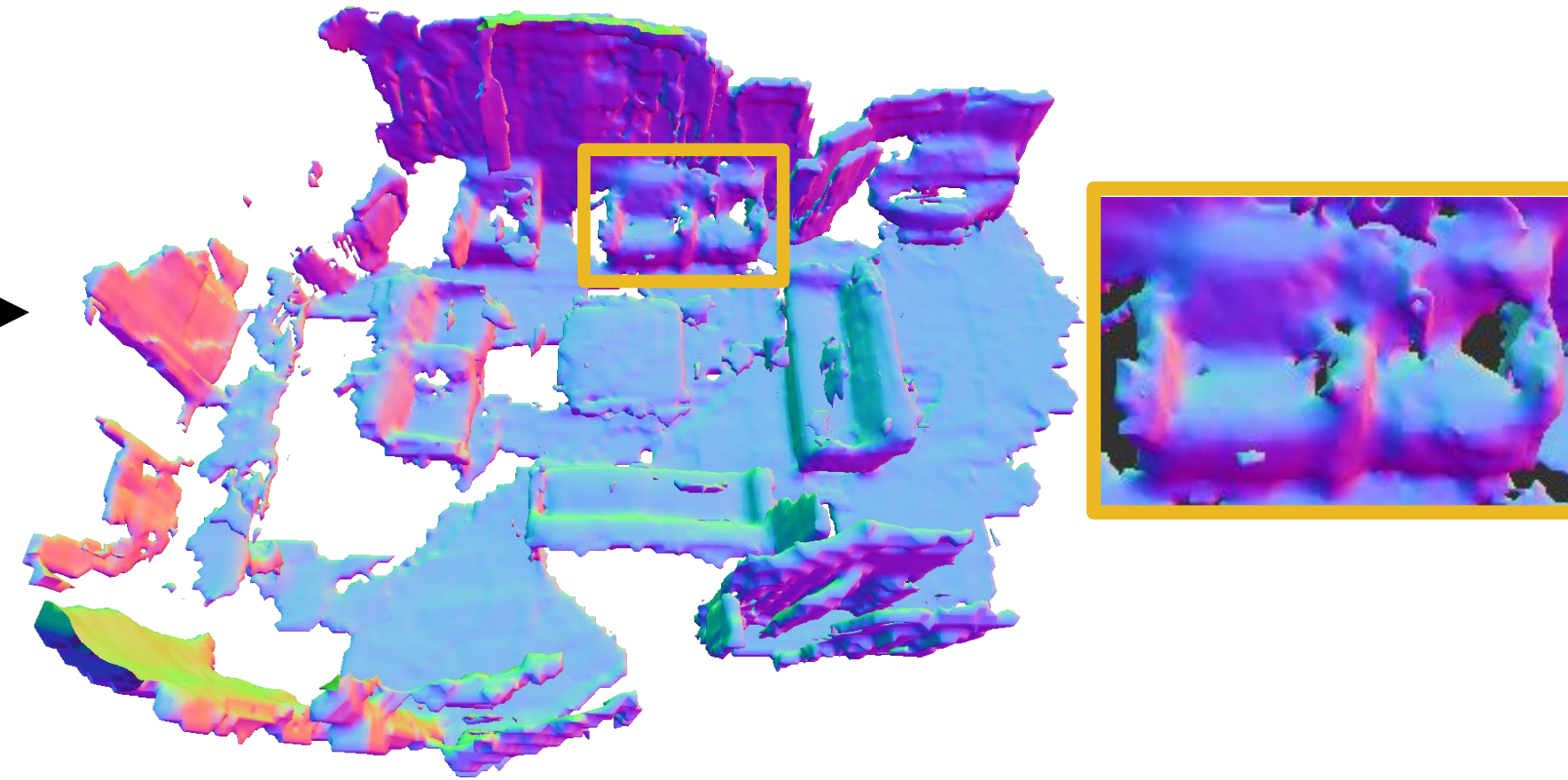
NeuralRecon

Joint reconstruction and fusion

Fragment
Reconstruction



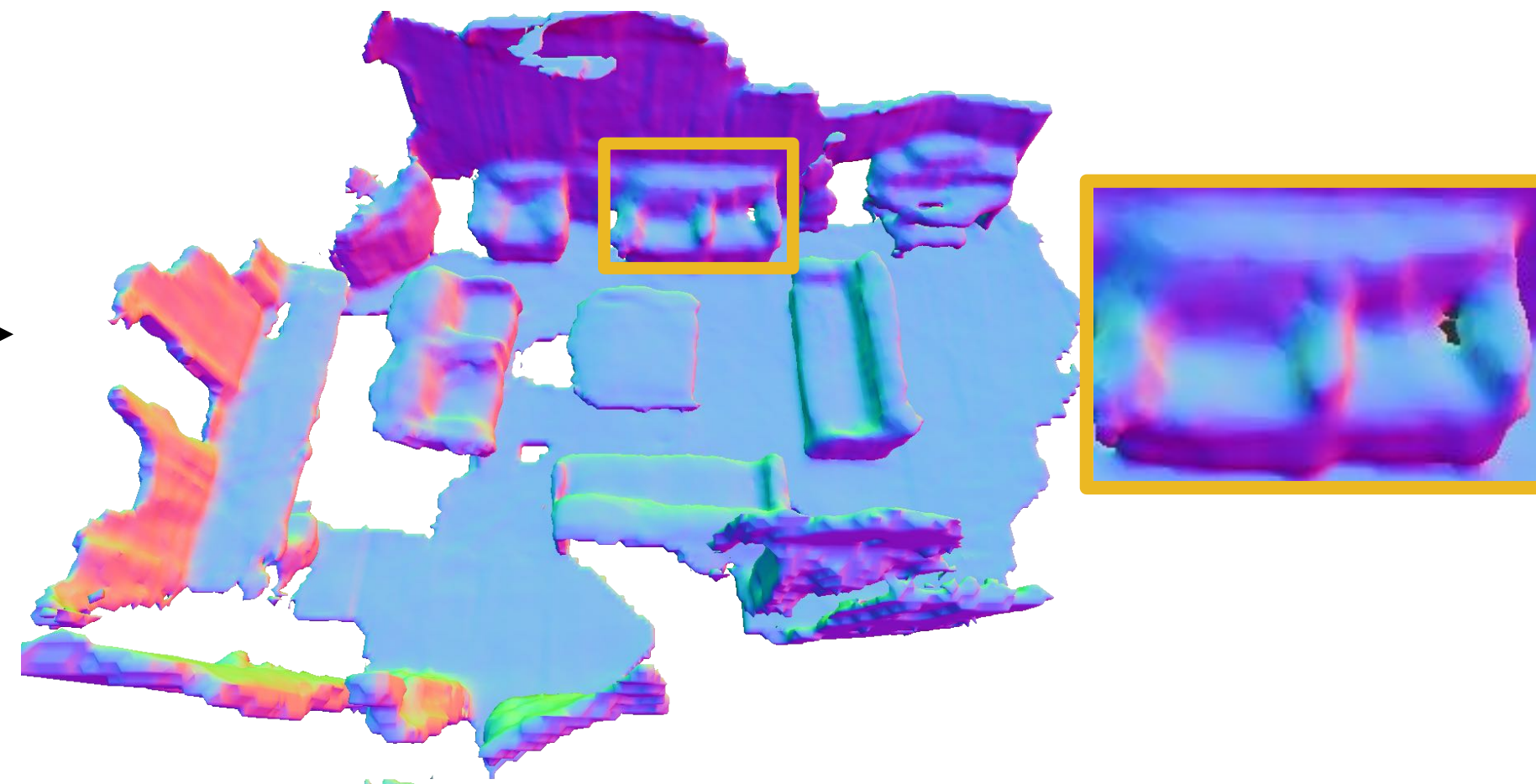
TSDF
Fusion



VS

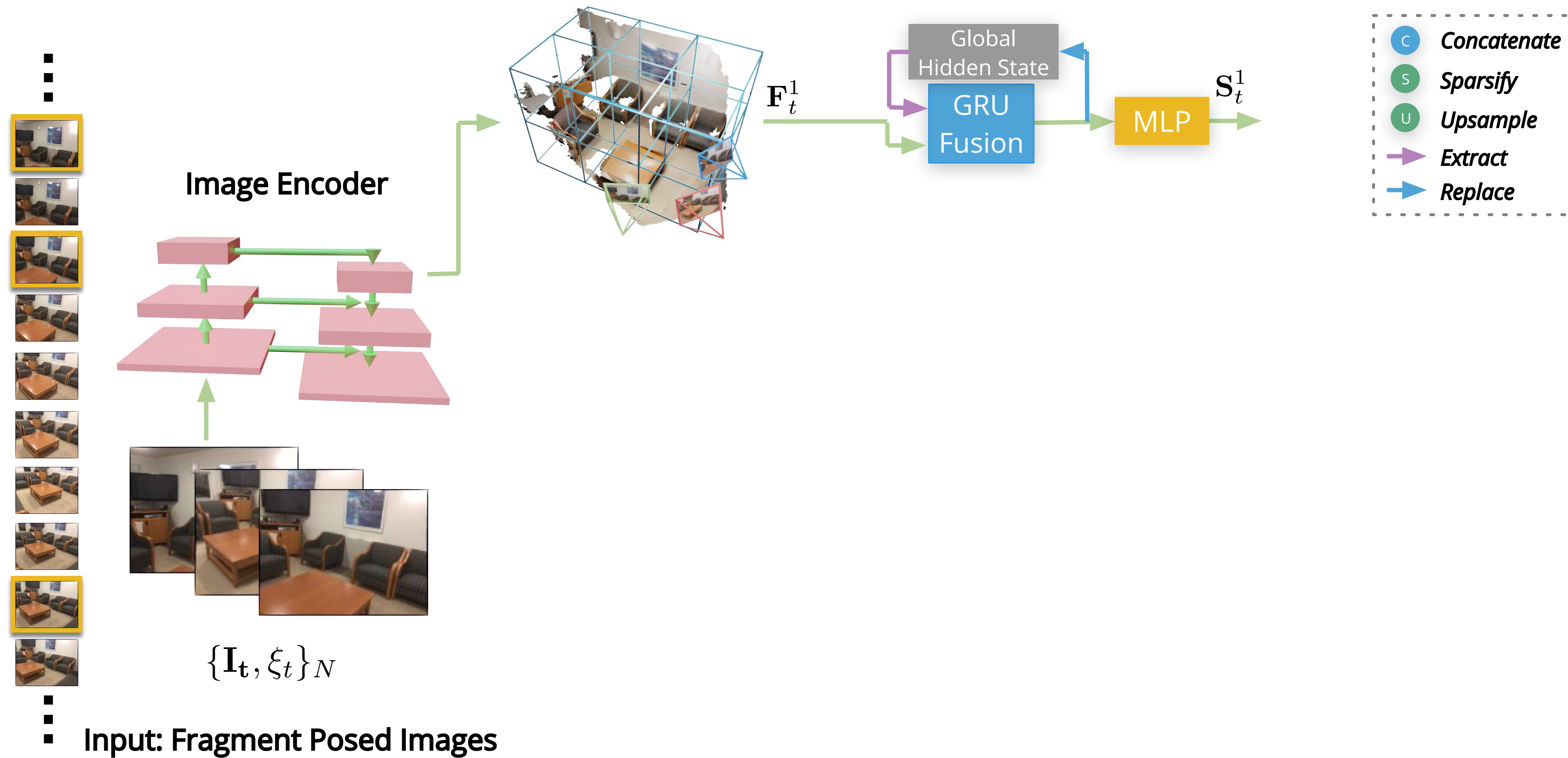
😊 Globally coherent

Joint
Reconstruction
and Fusion



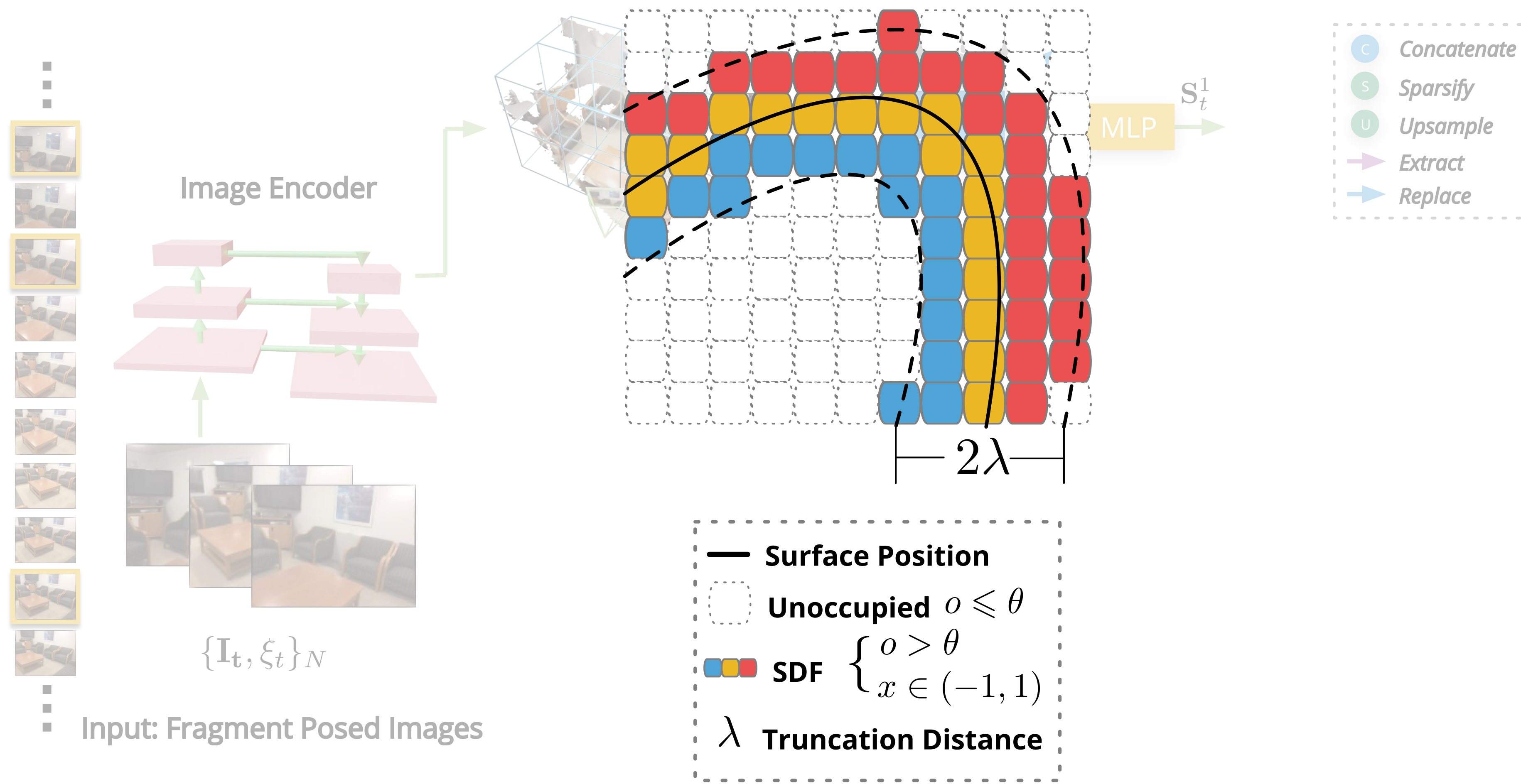
NeuralRecon

Coarse-to-fine architecture



NeuralRecon

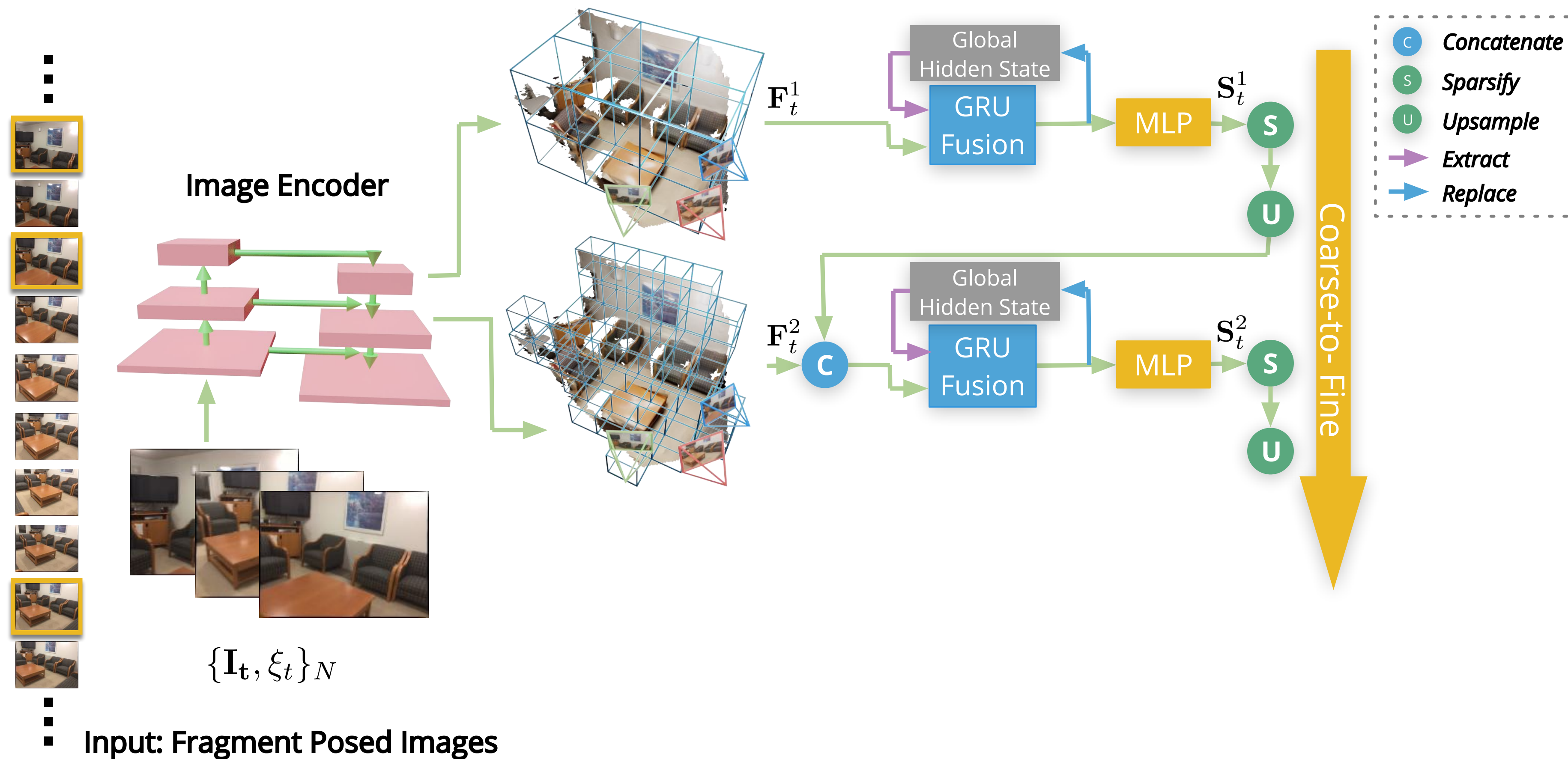
Coarse-to-fine architecture



Output of MLP : **Occupancy Score** and **SDF**

NeuralRecon

Coarse-to-fine architecture

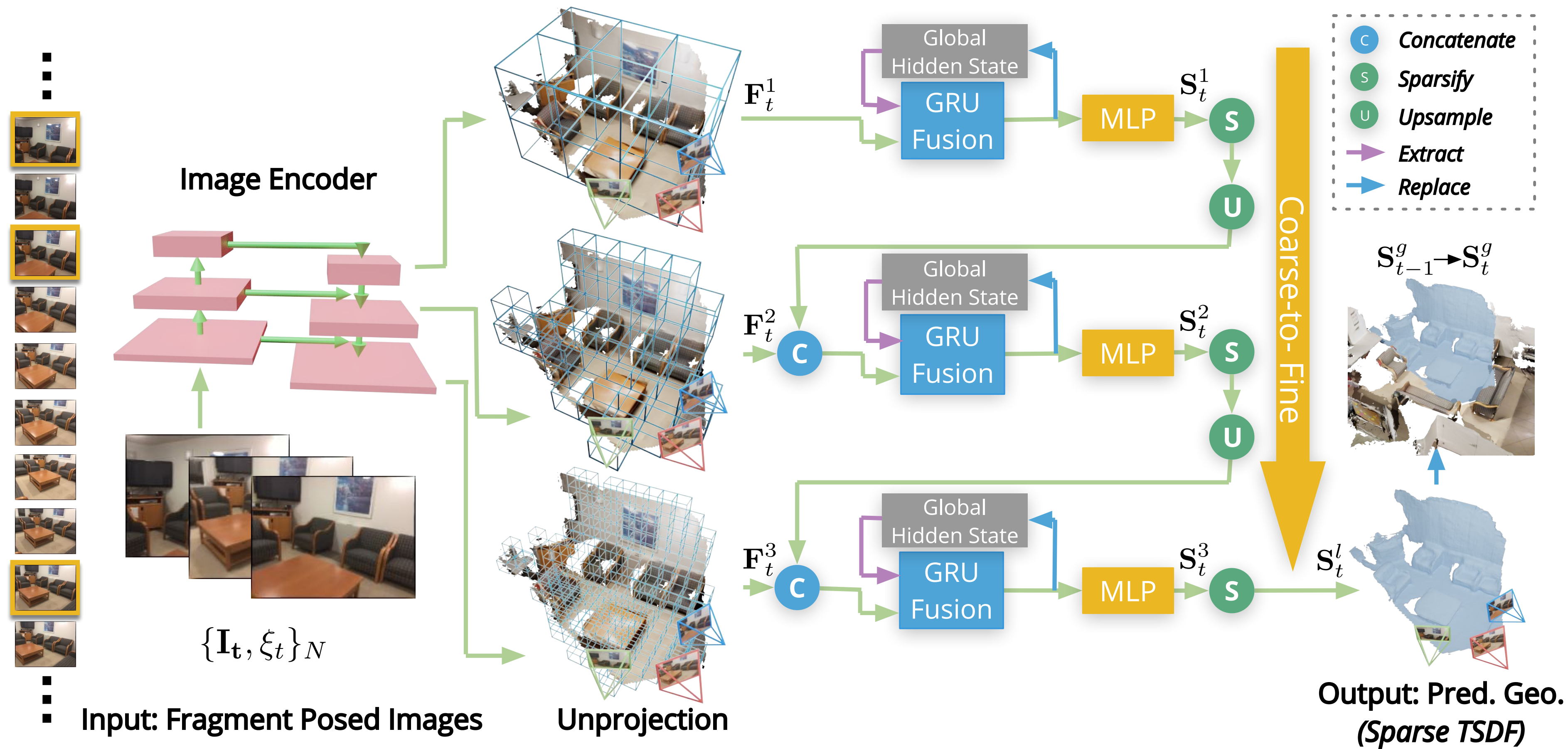


Output of *MLP* : **Occupancy Score** and **SDF**

S Filter by Occupancy Score > 0

NeuralRecon

Coarse-to-fine architecture

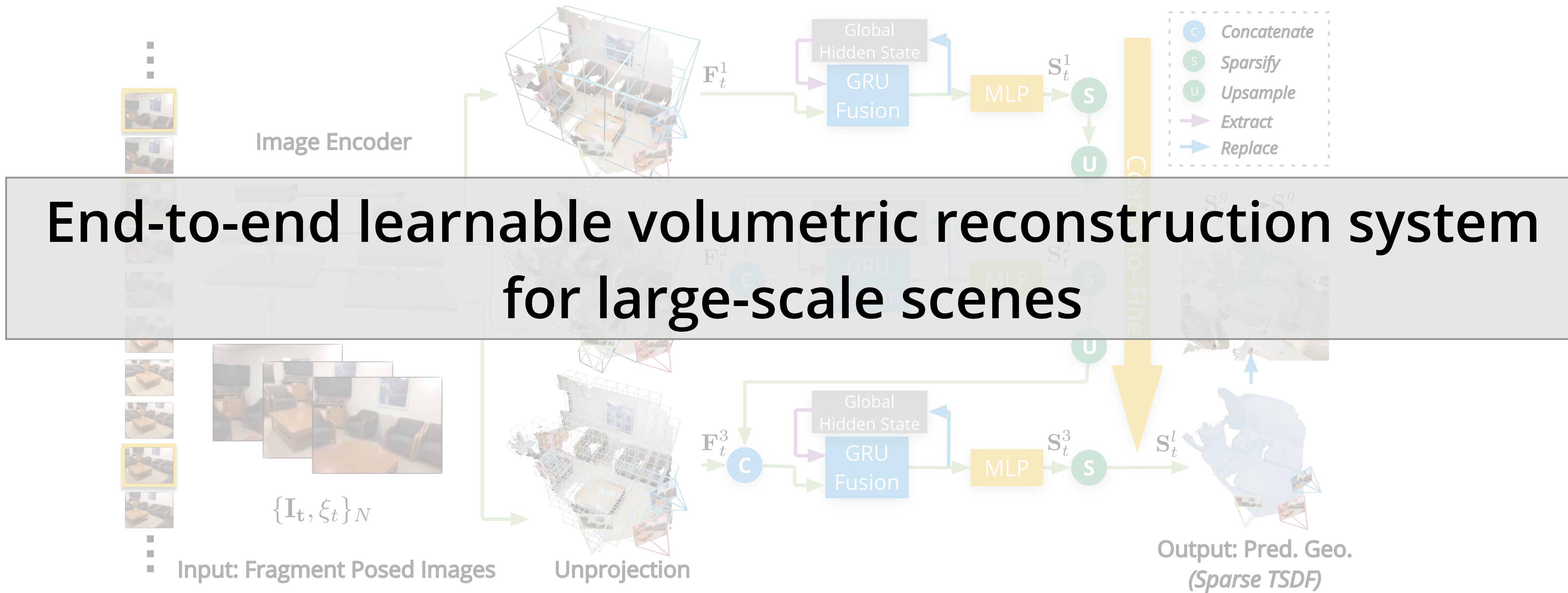


Output of *MLP* : **Occupancy Score** and **SDF**

S Filter by Occupancy Score > 0

NeuralRecon

Coarse-to-fine architecture



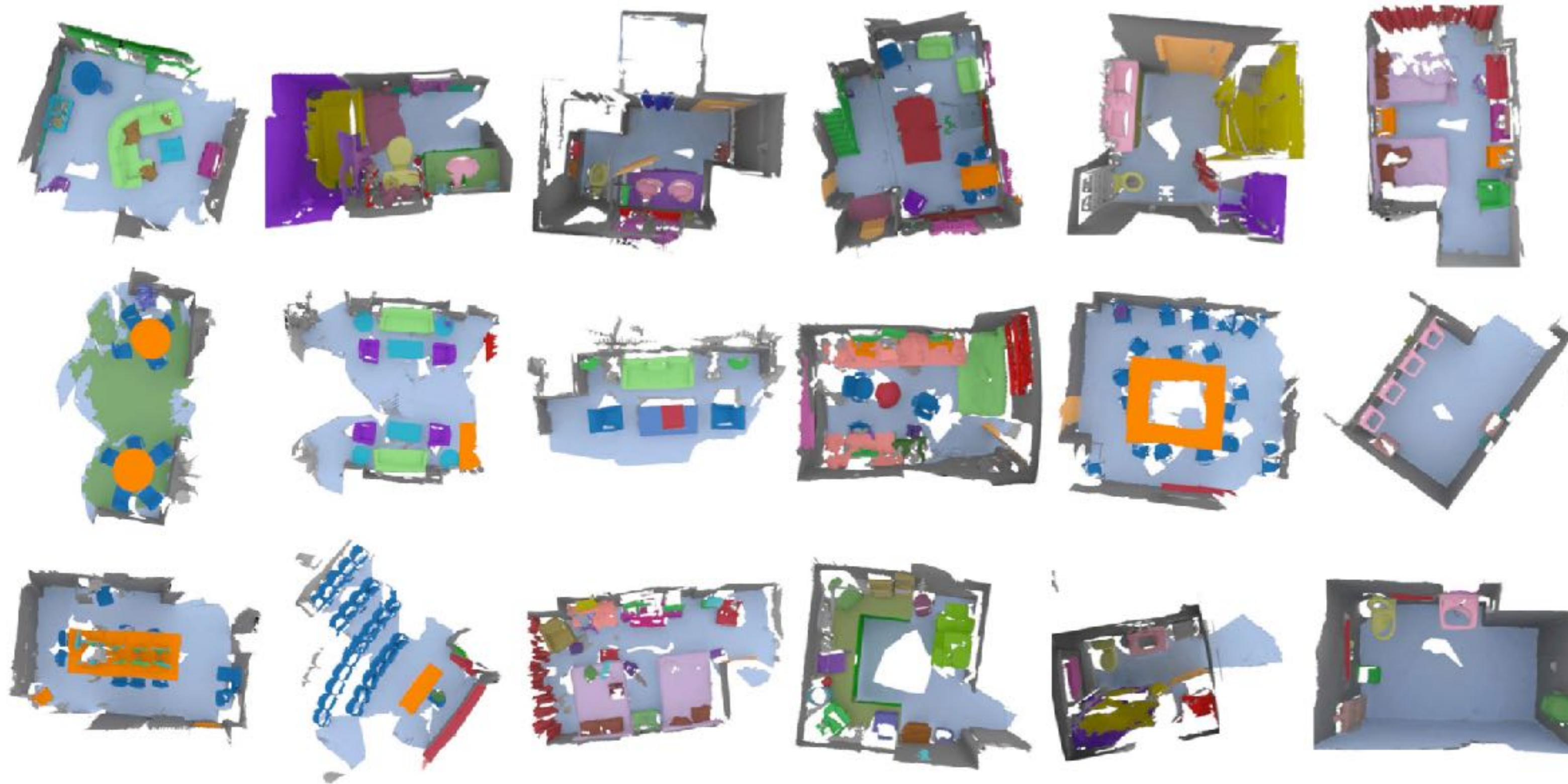
End-to-end learnable volumetric reconstruction system for large-scale scenes

Output of *MLP* : **Occupancy Score** and **SDF**

S Filter by Occupancy Score > 0

NeuralRecon

Training



ScanNet dataset

Contains 2.5M RGB images captured in 707 indoor scenes with ground-truth TSDF

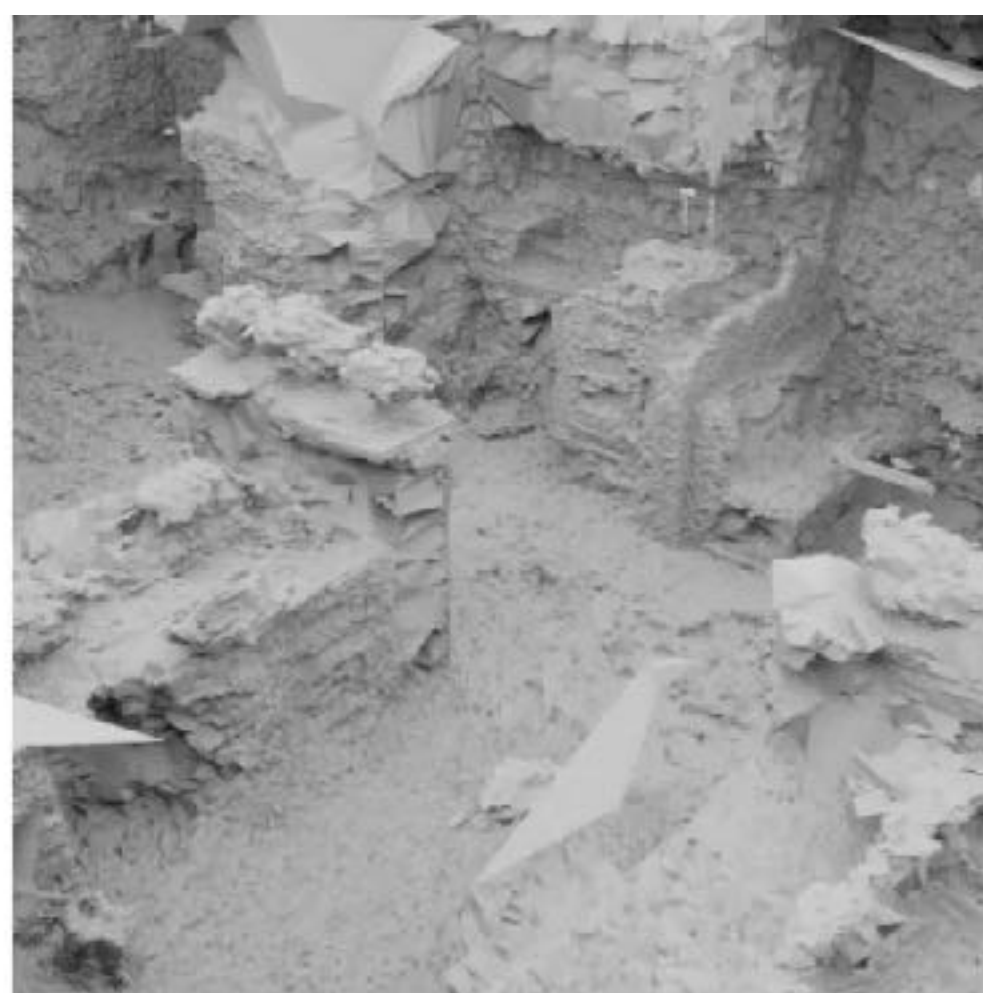
Binary cross-entropy (BCE) and L1 loss are used for training

Experiments

Qualitative results: office 1



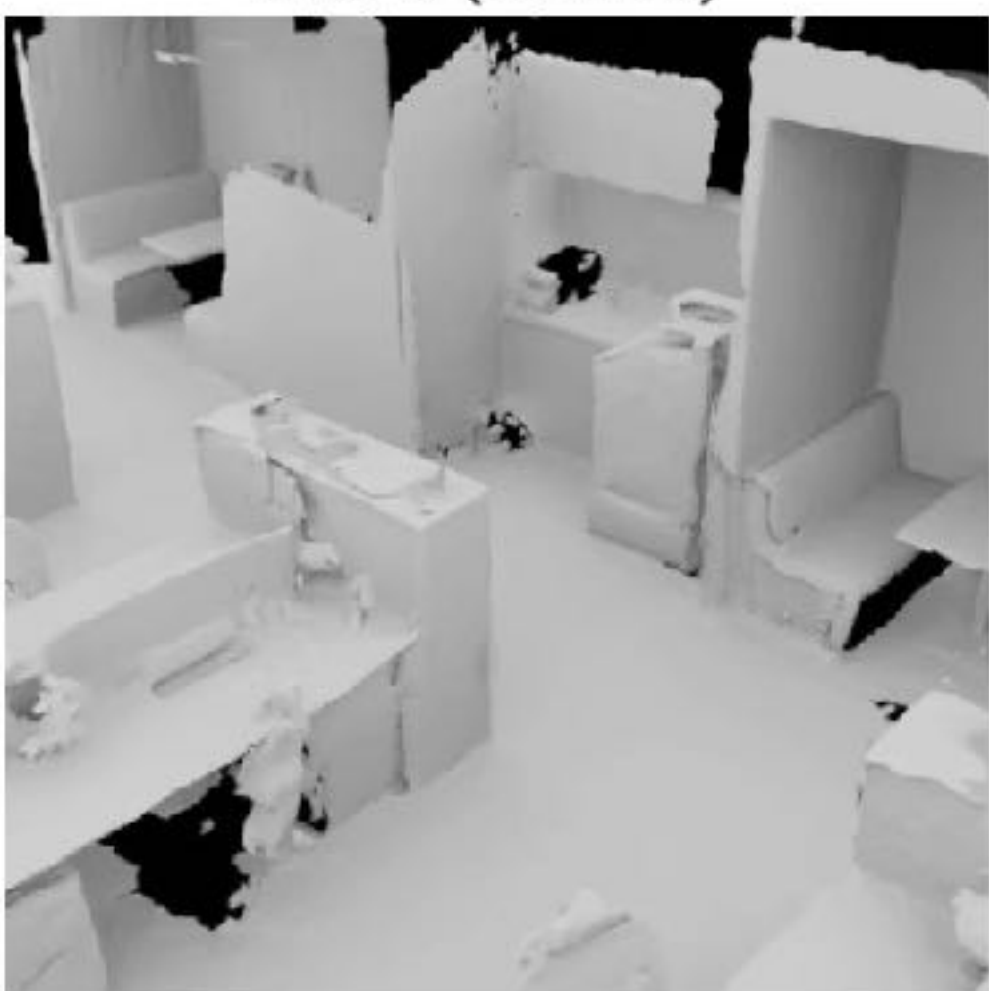
Ours (30ms)



COLMAP (2076ms)



DeepV2D (347ms)



Ground Truth



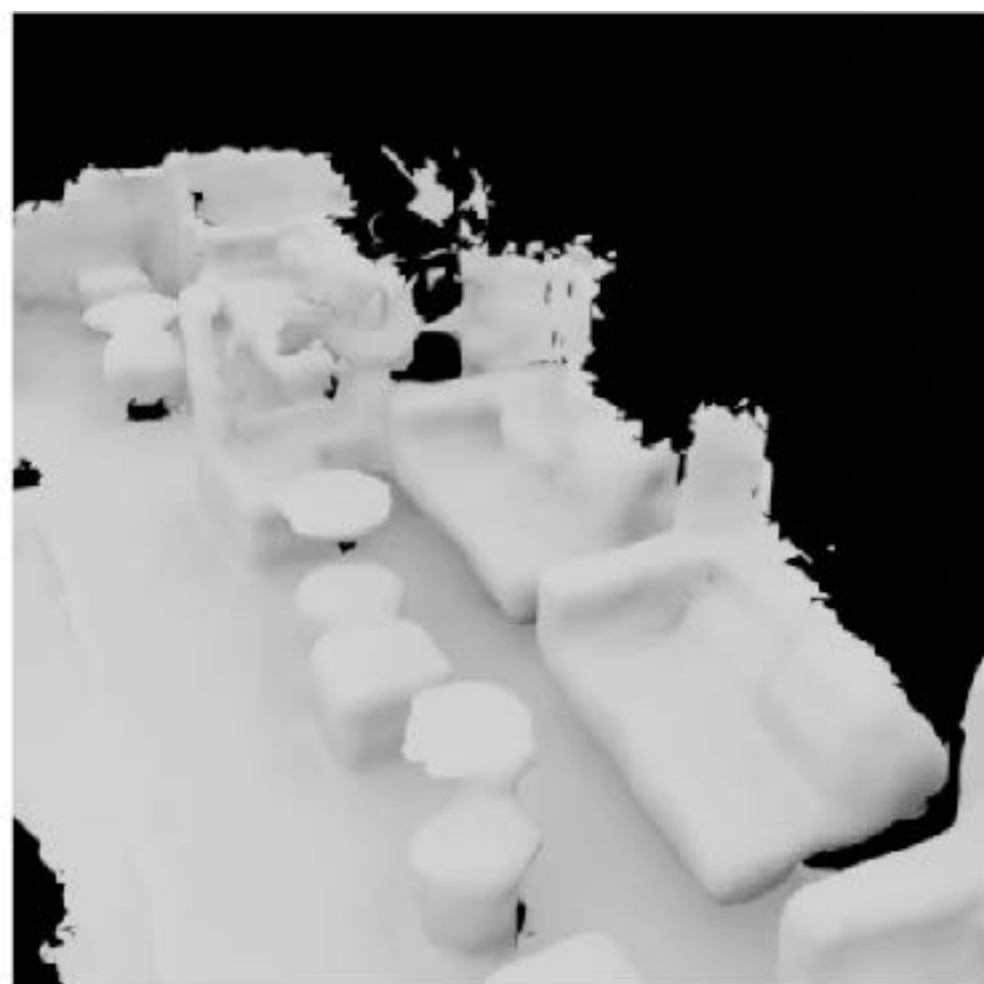
CNMNet (80ms)



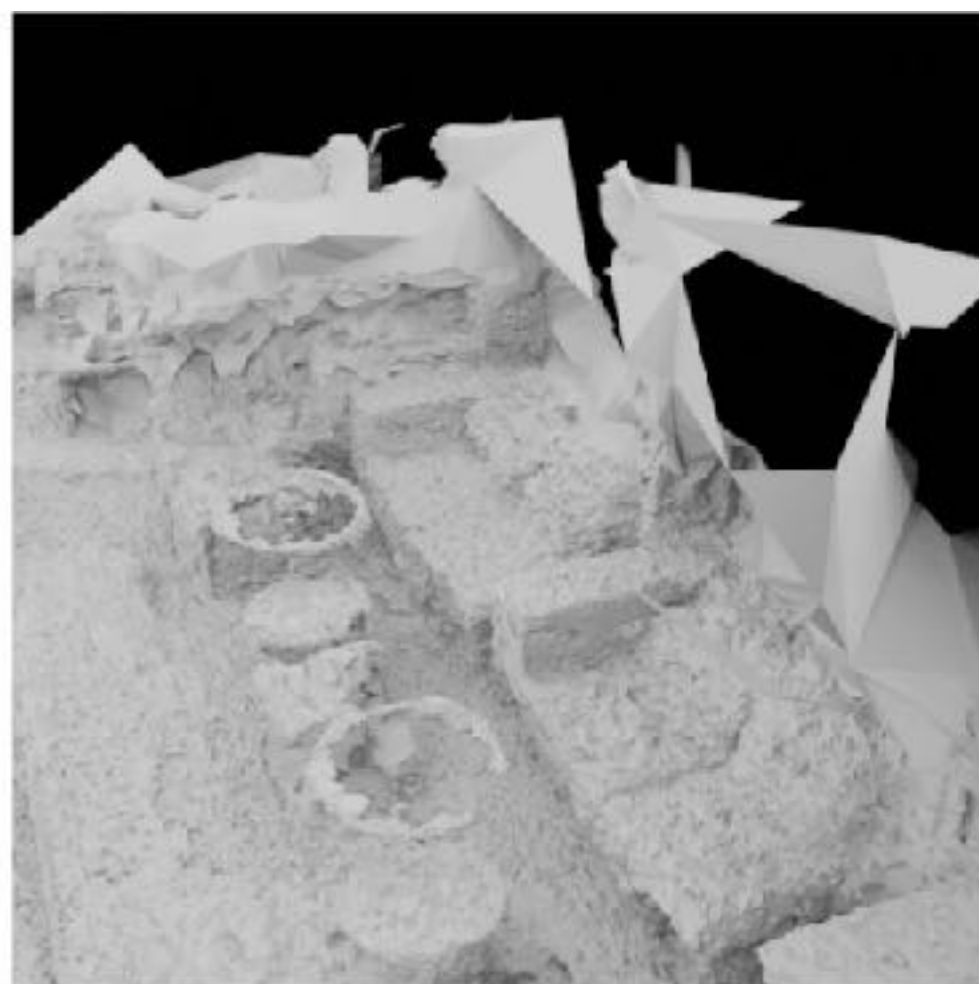
Atlas (292ms)

Experiments

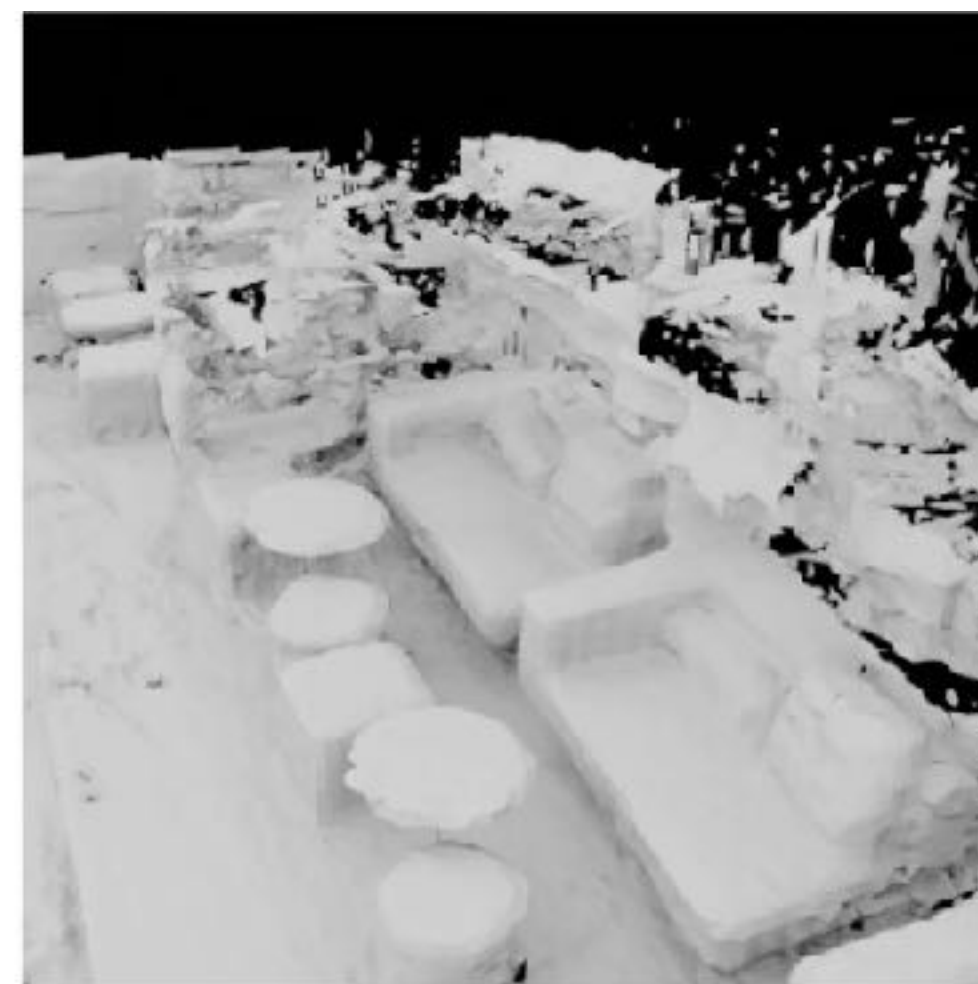
Qualitative results: office 2



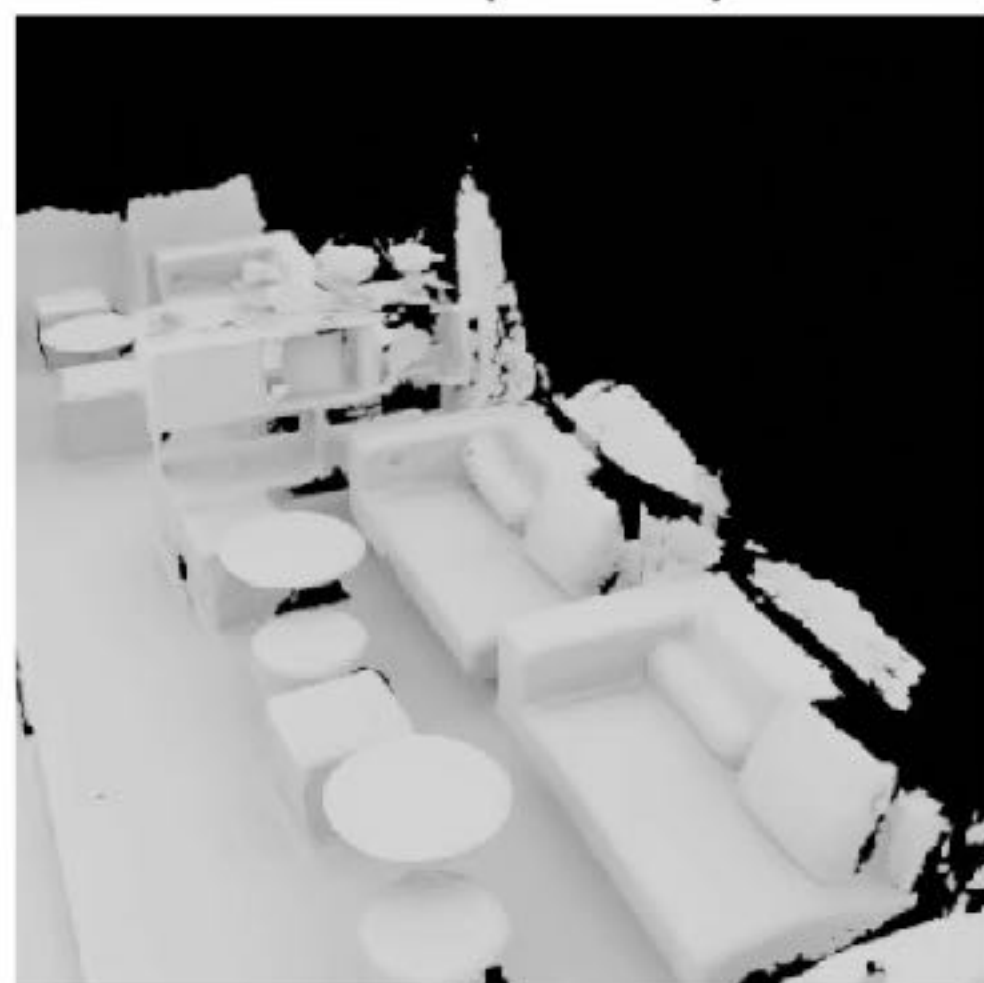
Ours (30ms)



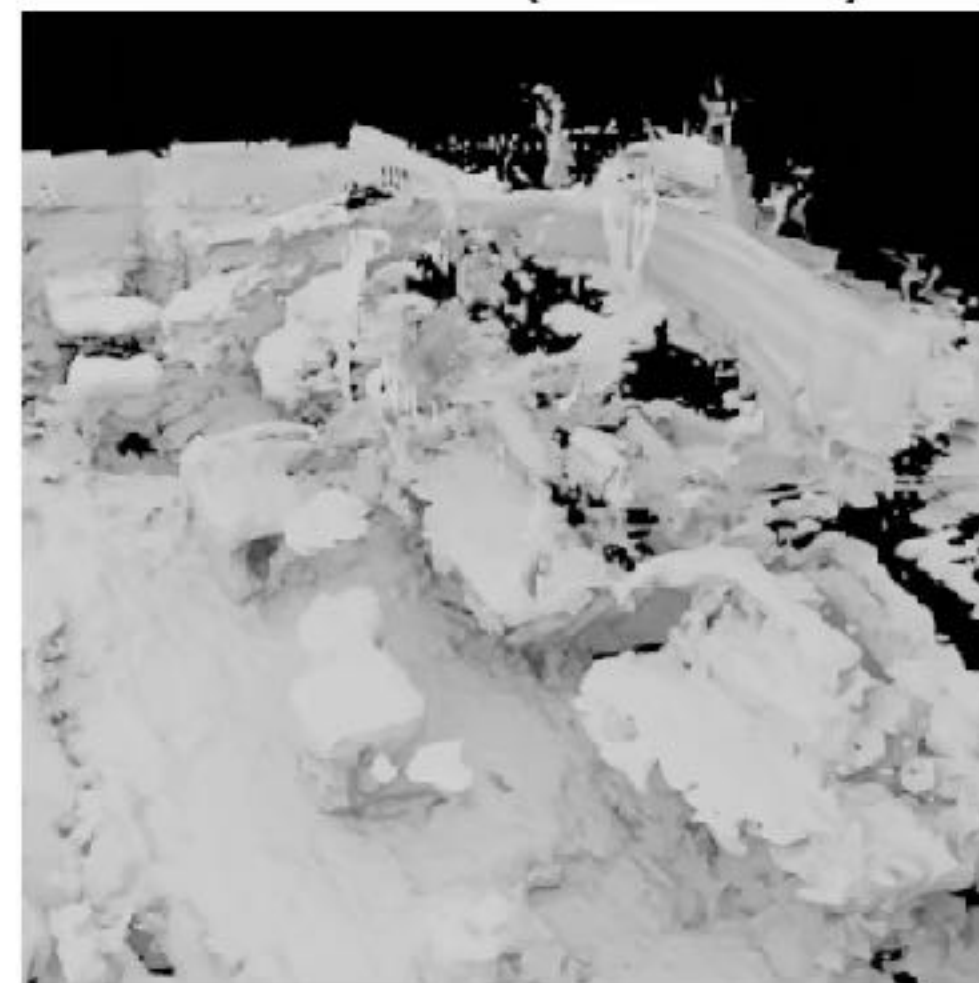
COLMAP (2076ms)



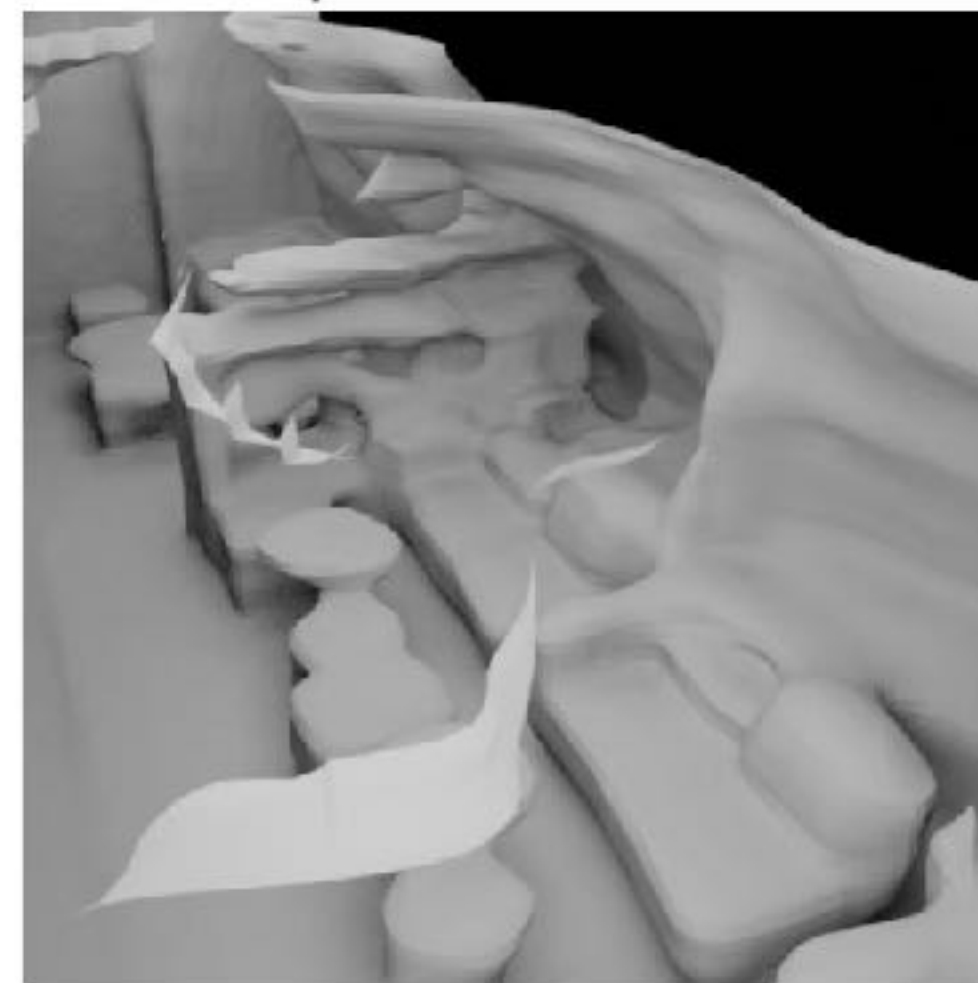
DeepV2D (347ms)



Ground Truth



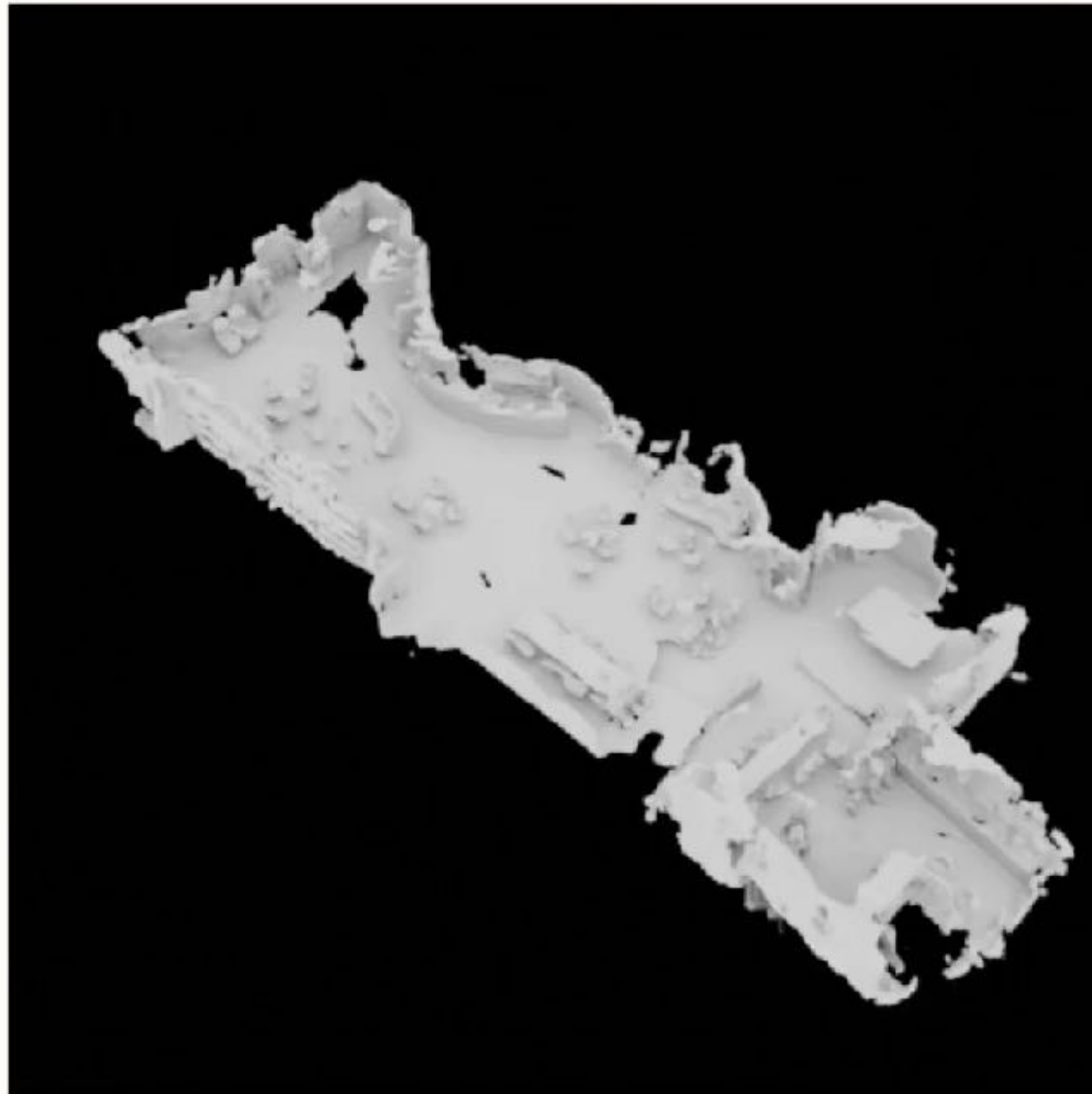
CNMNet (80ms)



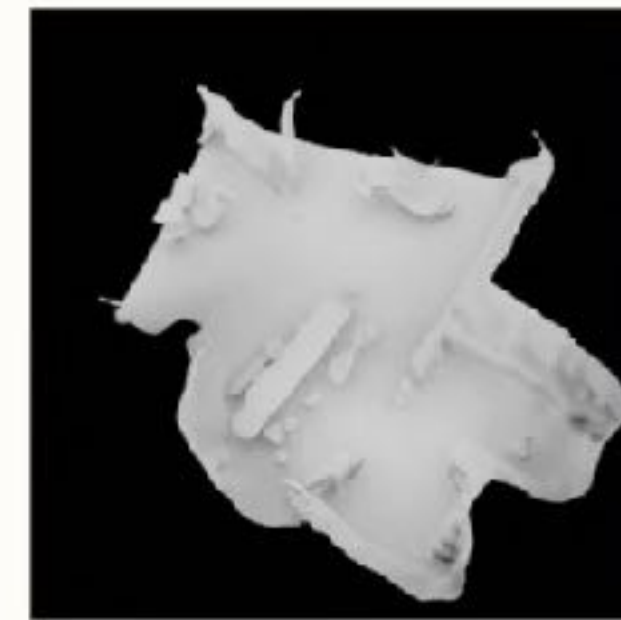
Atlas (292ms)

Experiments

Qualitative results: Comparison with Atlas on a **large** scene (30m x 10m)



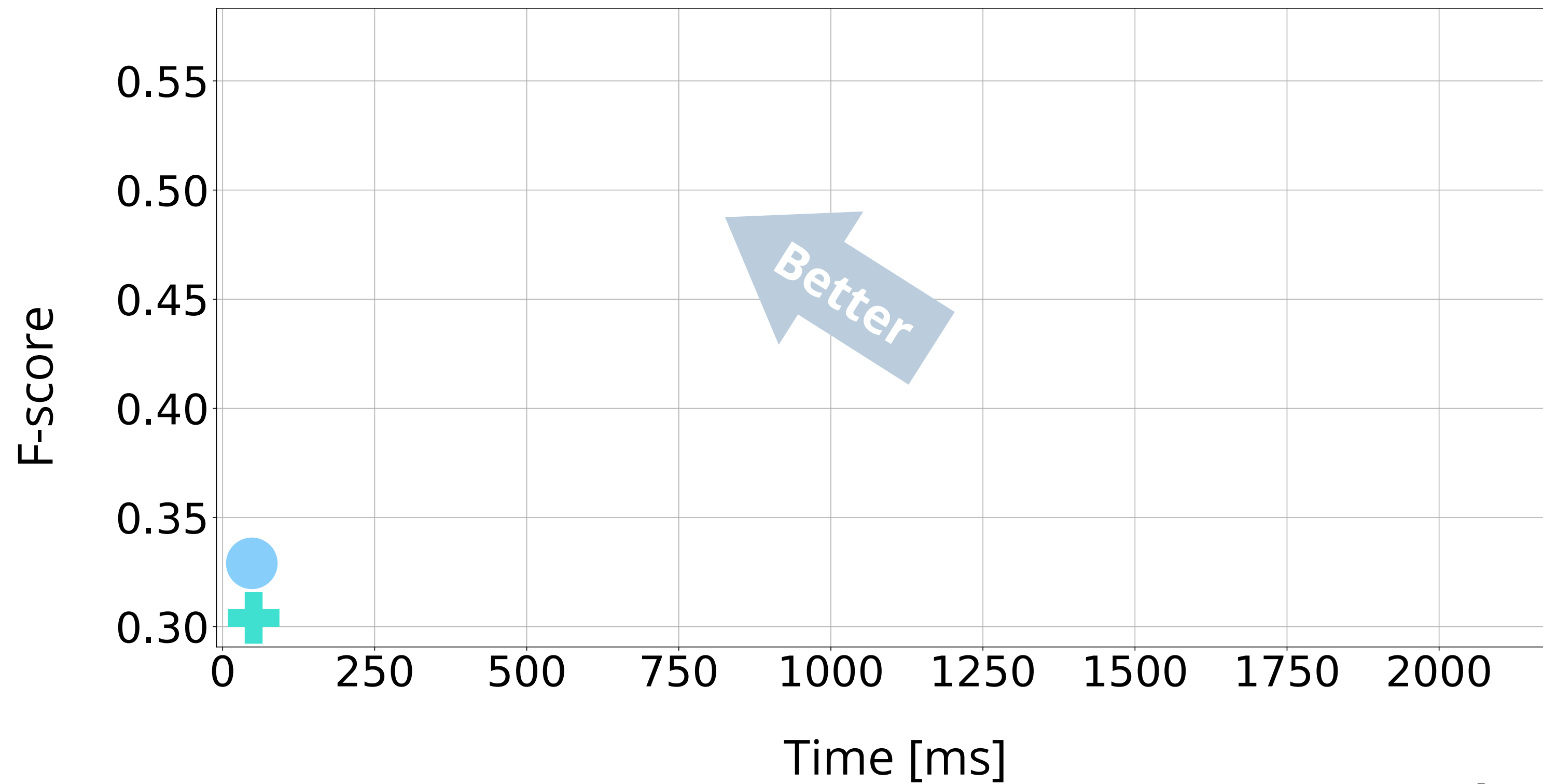
Ours
Max GPU Memory: 3.29GB



Atlas
Max GPU Memory: >24GB (OOM)
The reconstruction is incomplete
due to out of memory (OOM) error on the remaining sequence.

Experiments

Quantitative results



Real-time methods: ● MVDepthNet + GPMVS

Speed

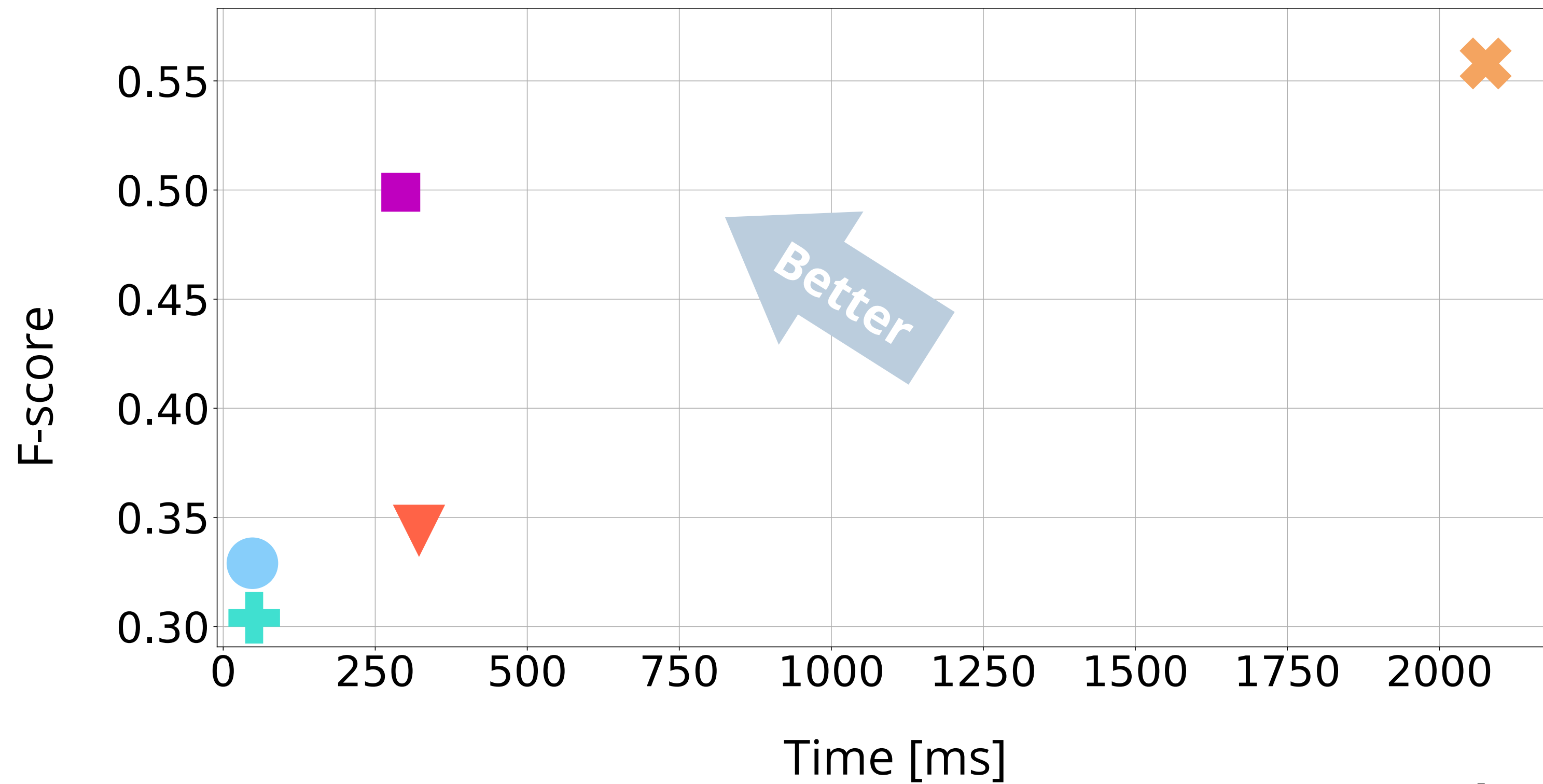


Accuracy



Experiments

Quantitative results



Real-time methods:

● MVDepthNet

+ GPMVS

Speed



Accuracy



Multiple View Stereo methods:

▼ DPSNet

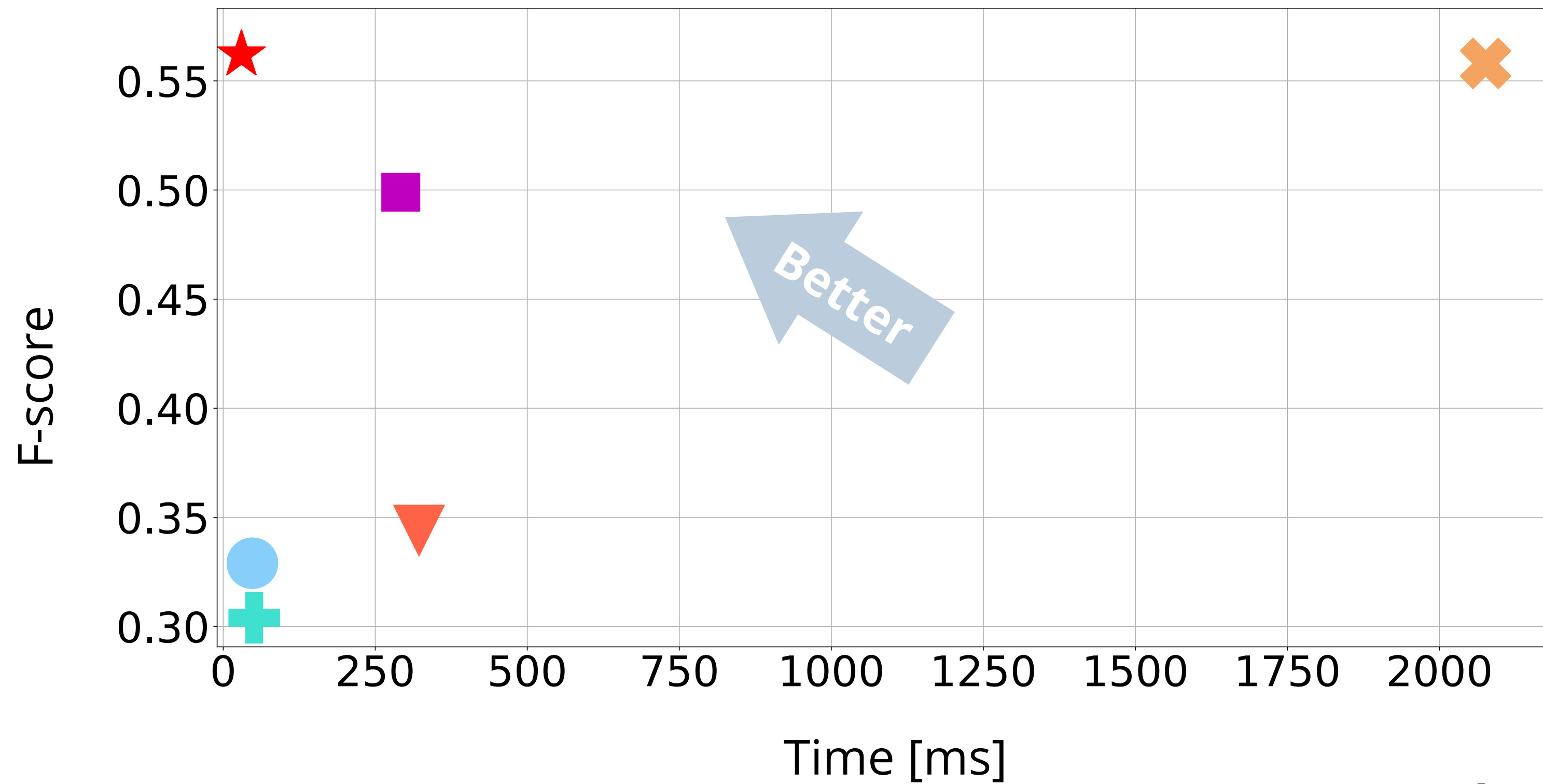
× COLMAP

■ Atlas



Experiments

Quantitative results



Real-time methods:

● MVDepthNet

+ GPMVS

Speed



Accuracy



Multiple View Stereo methods:

▼ DPSNet

× COLMAP

■ Atlas

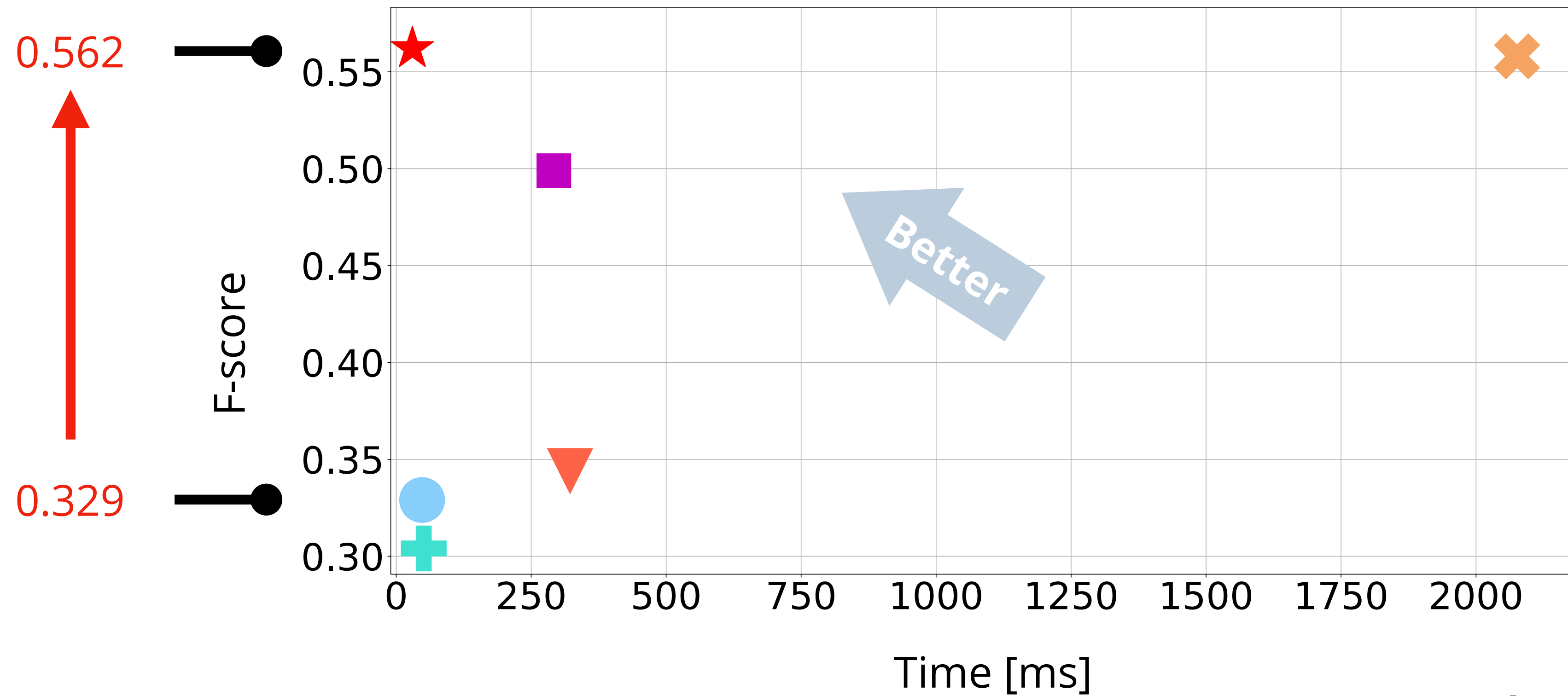


★ Ours



Experiments

Quantitative results



Real-time methods:

● MVDepthNet

+ GPMVS

Speed



Accuracy



Multiple View Stereo methods:

▼ DPSNet

× COLMAP

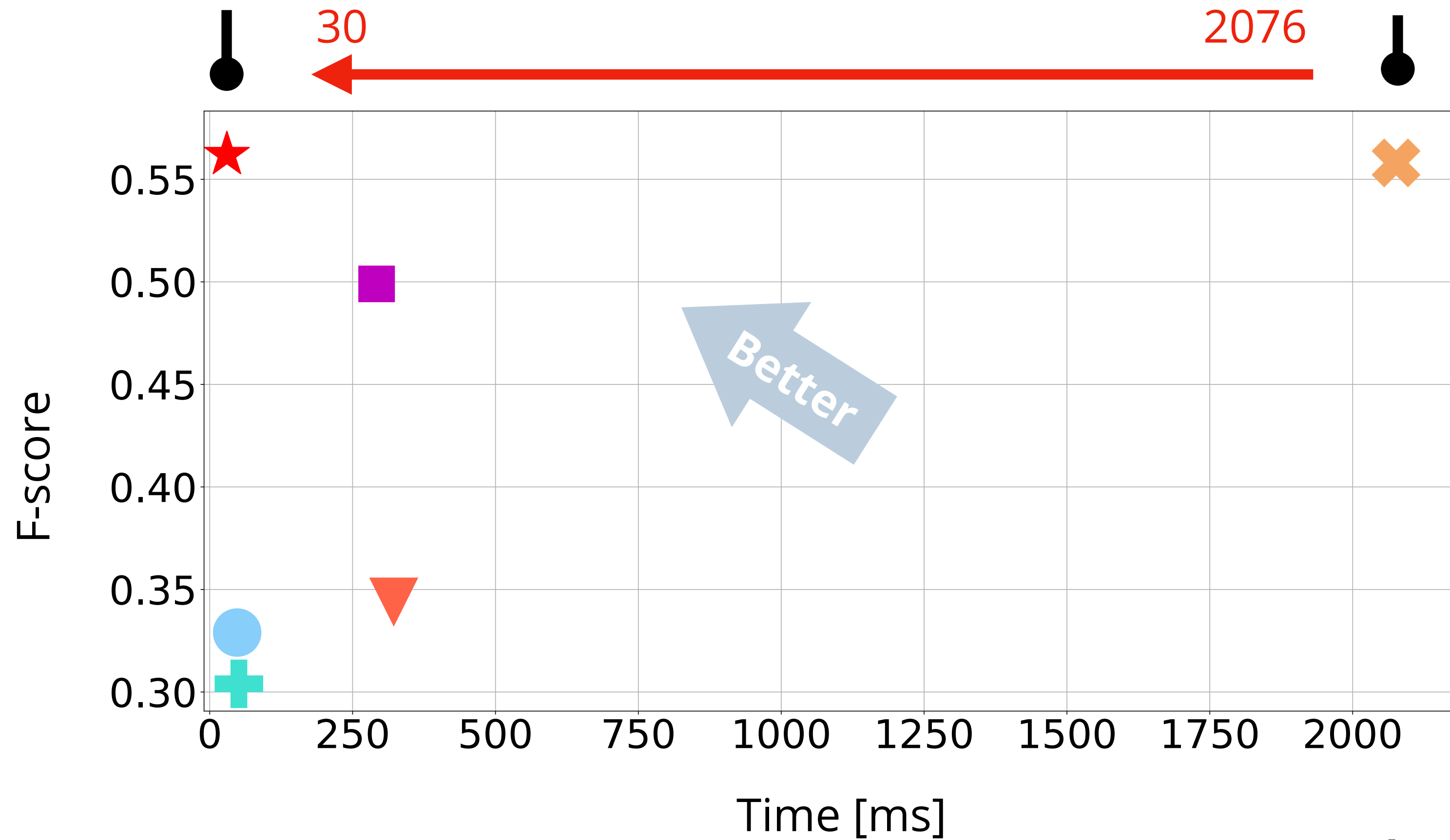
■ Atlas



★ Ours



Experiments



			Speed	Accuracy
<i>Real-time methods:</i>	● MVDDepthNet	+ GPMVS	😊	😞
<i>Multiple View Stereo methods:</i>	▼ DPSNet	× COLMAP	😞	😊
	★ Ours	■ Atlas	😊	😊

Demo

Indoor scene at our office



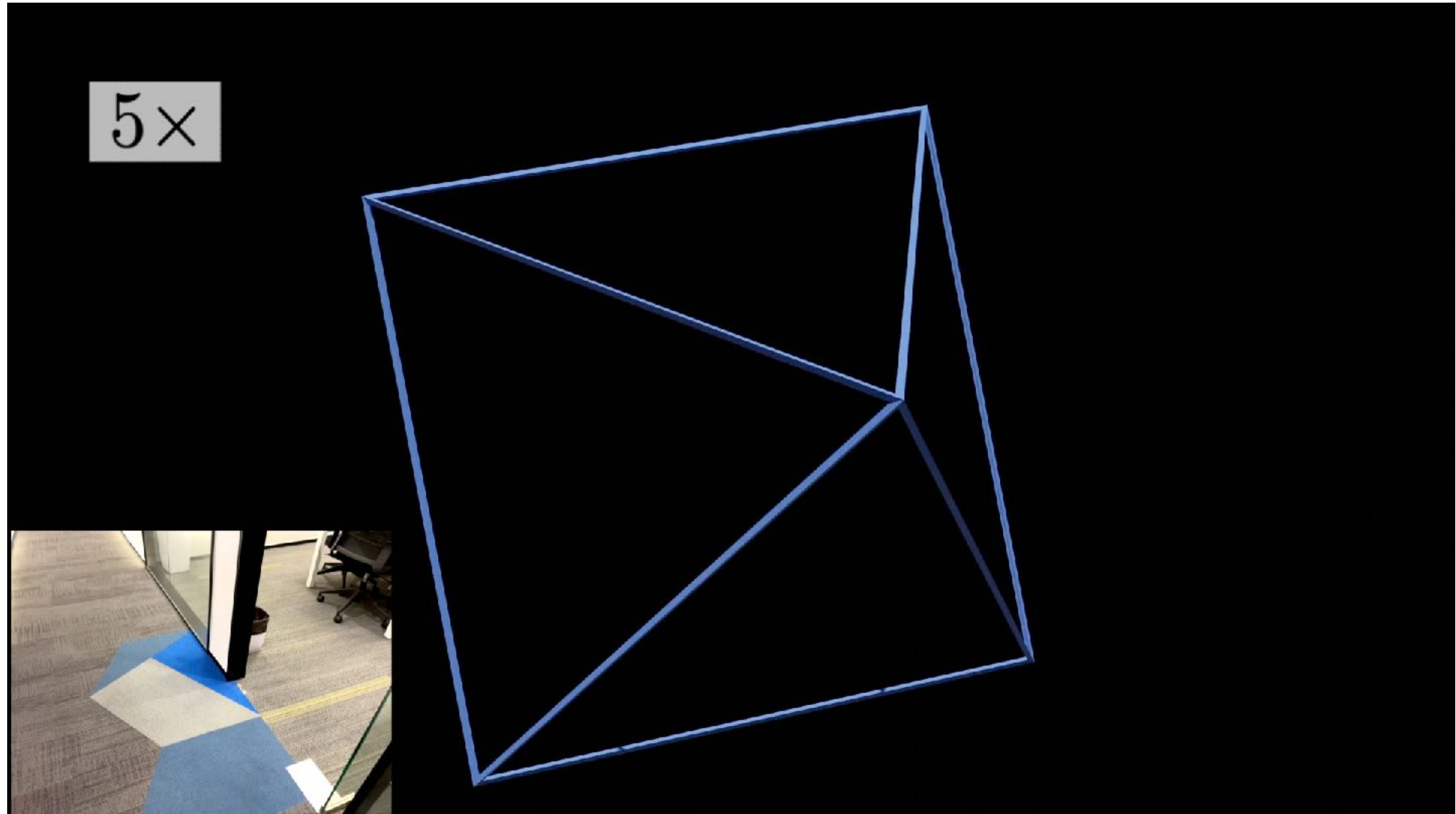
Input video with camera poses



3D reconstruction

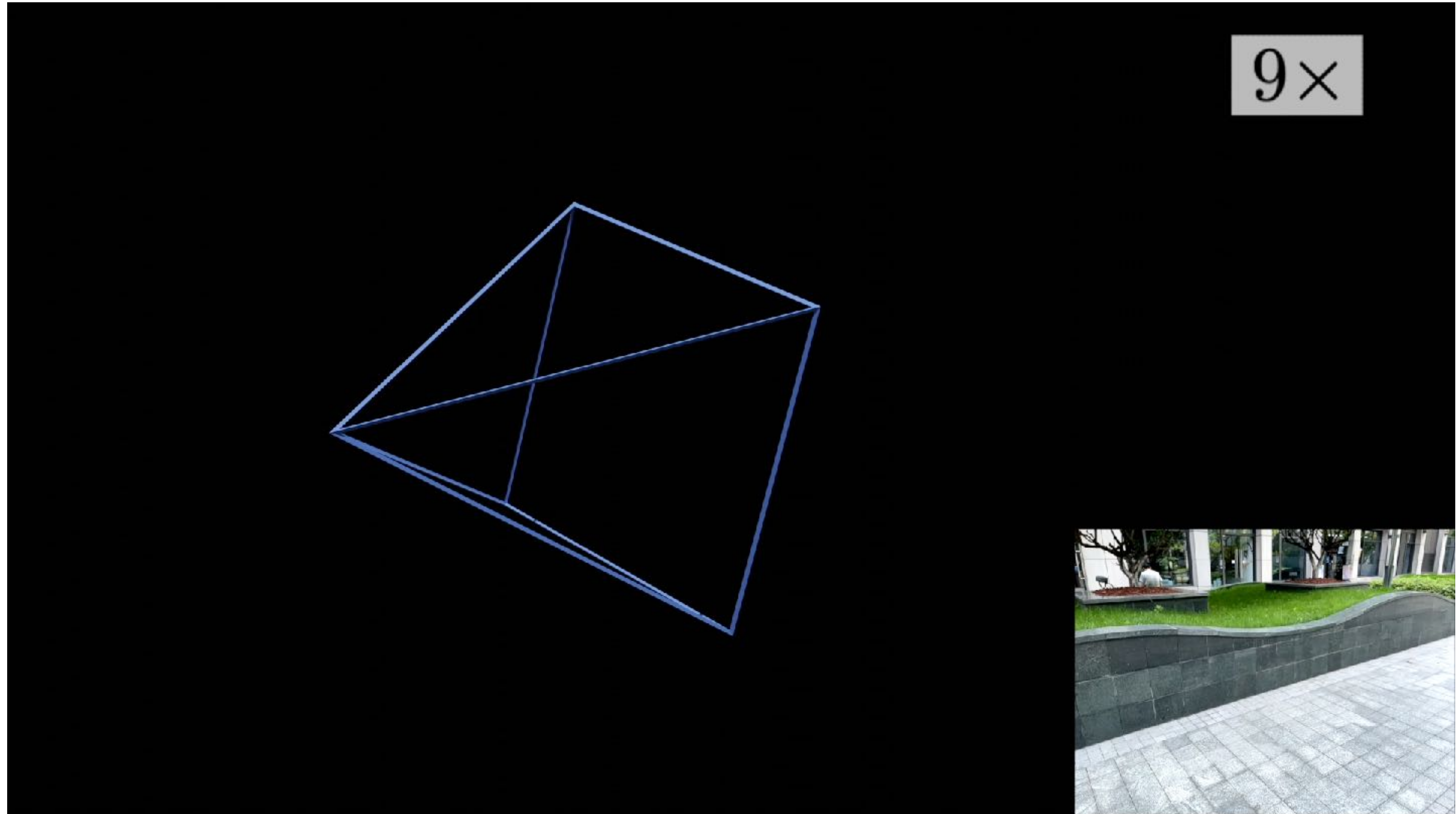
Demo

Indoor scene with extremely low texture



Demo

Generalization to outdoor scenes



Demo

AR Demo



NeuralRecon: Real-Time Coherent 3D Reconstruction from Monocular Video

Jiaming Sun* Yiming Xie* Linghao Chen Xiaowei Zhou Hujun Bao
CVPR 2021 (Oral)

Project page: <https://zju3dv.github.io/neuralrecon/>

Code: <https://zju3dv.github.io/NeuralRecon/>

Paper link: <https://arxiv.org/pdf/2104.00681.pdf>

Conclusion and Discussion

- **Learning-based 3D reconstruction and localization**
 - Determine what's really necessary to be learned
(e.g. features for matching, surface prior for reconstruction and fusion. BA, PnP, RANSAC are fine and should be kept.)
 - Design end-to-end learnable systems with differentiable modules
(e.g. GRU-based TSDF reconstruction and fusion in NeuralRecon, differentiable matching in LoFTR, differentiable and learnable optimization, etc.)

Conclusion and Discussion

- **Learning-based 3D reconstruction and localization**
 - Determine what's really necessary to be learned
(e.g. features for matching, surface prior for reconstruction and fusion. BA, PnP, RANSAC are fine and should be kept.)
 - Design end-to-end learnable systems with differentiable modules
(e.g. GRU-based TSDF reconstruction and fusion in NeuralRecon, differentiable matching in LoFTR, differentiable and learnable optimization, etc.)
- **Future 3D perception systems?**
 - Learnable and embodied (e.g. semantics, context-dependent)
 - Self-supervised from observations and interactions
(e.g. with differentiable rendering, simulators)
 - Hardcoded with well-established knowledges as prior
(e.g. multi-view geometry, laws of physics)

Acknowledgements



Xiaowei Zhou



Yiming Xie



Zehong Shen



Yu'ang Wang



Linghao Chen

Thanks for watching!
Q&A