# From Paper to Progress: MASt3R-SLAM: Real-Time Dense SLAM with 3D Reconstruction Priors

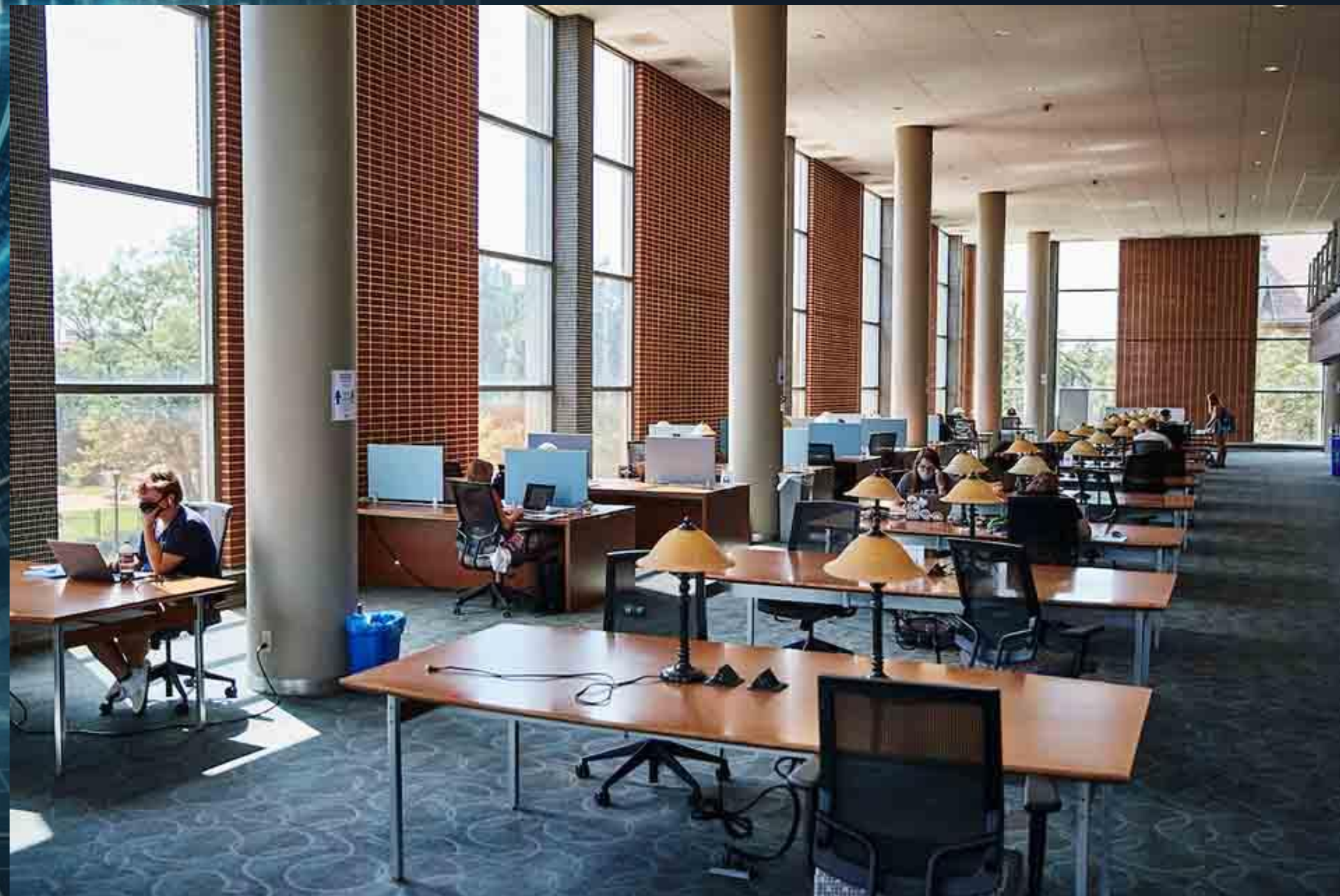**Guangping Liu, Ravali Gangineni, Rubeena Sultana**

# Overview

EaTemp

# Introduction

## What is SLAM?

Simultaneous Localization and Mapping (SLAM) is a method used for mobile robots that lets you build a map and localize your vehicle in that map at the same time.
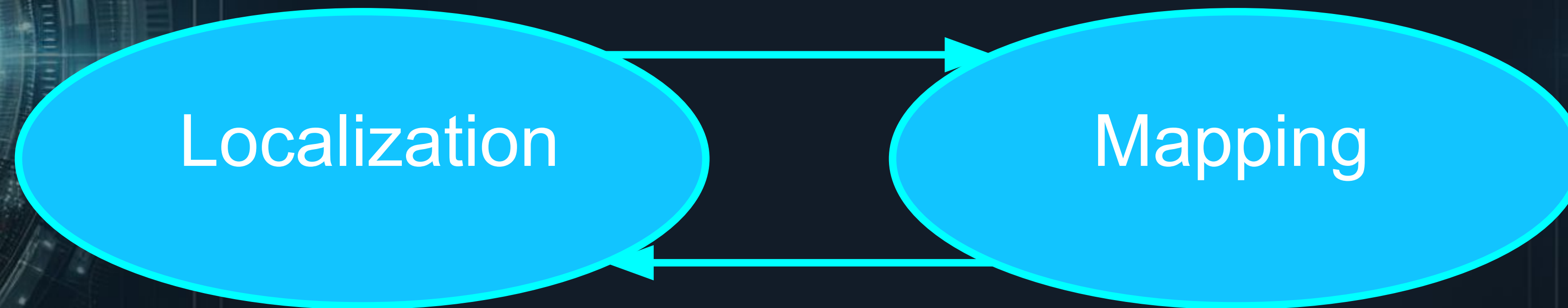


I am a new cleaner at SLU. How should I plan my work to clean the floor of the PLUS XII Library?

-I visited the library to plan my cleaning route.

EaTemp

# Introduction

What is SLAM?



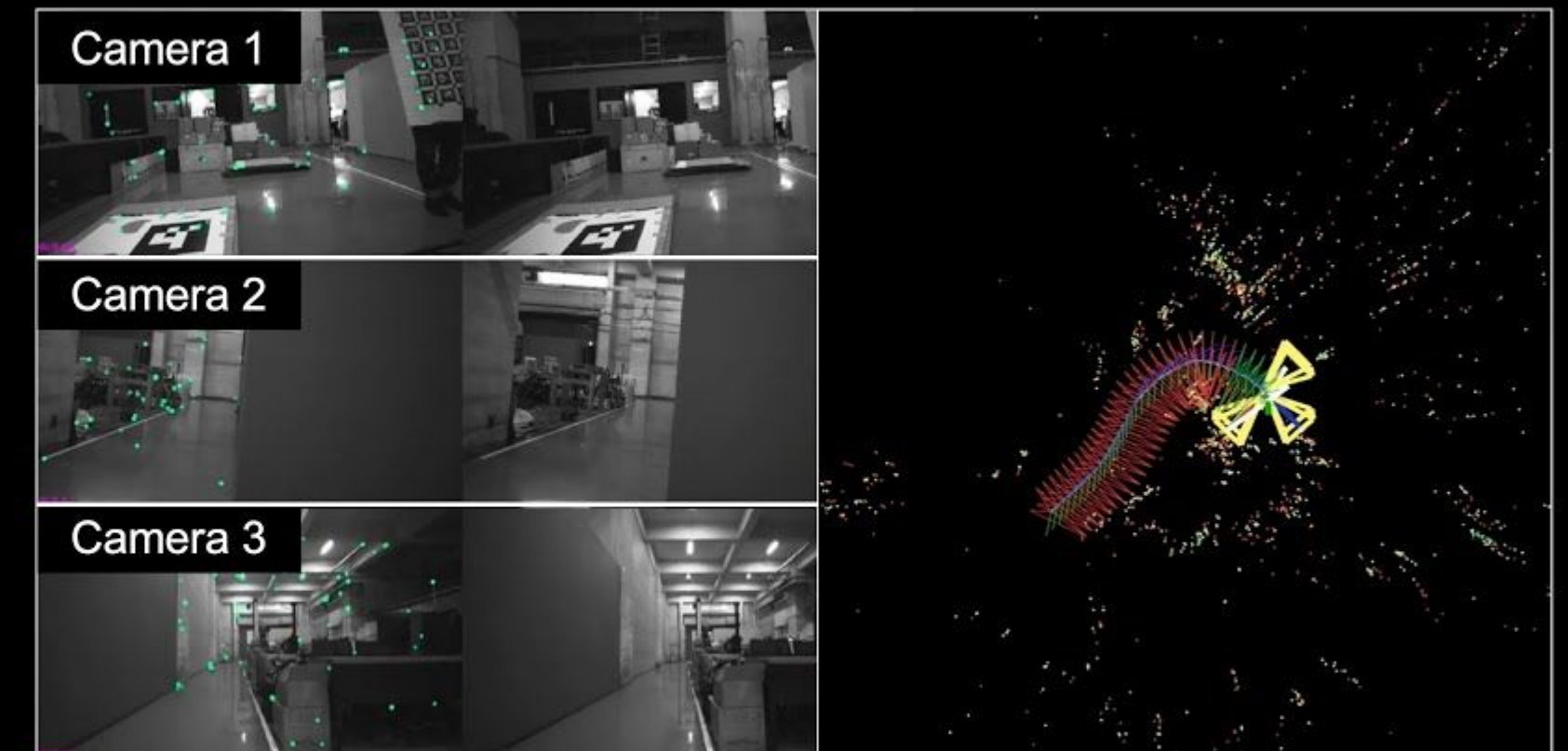| Localization | Mapping |
| --- | --- |
| Where I am? | What's around me? |

# Introduction
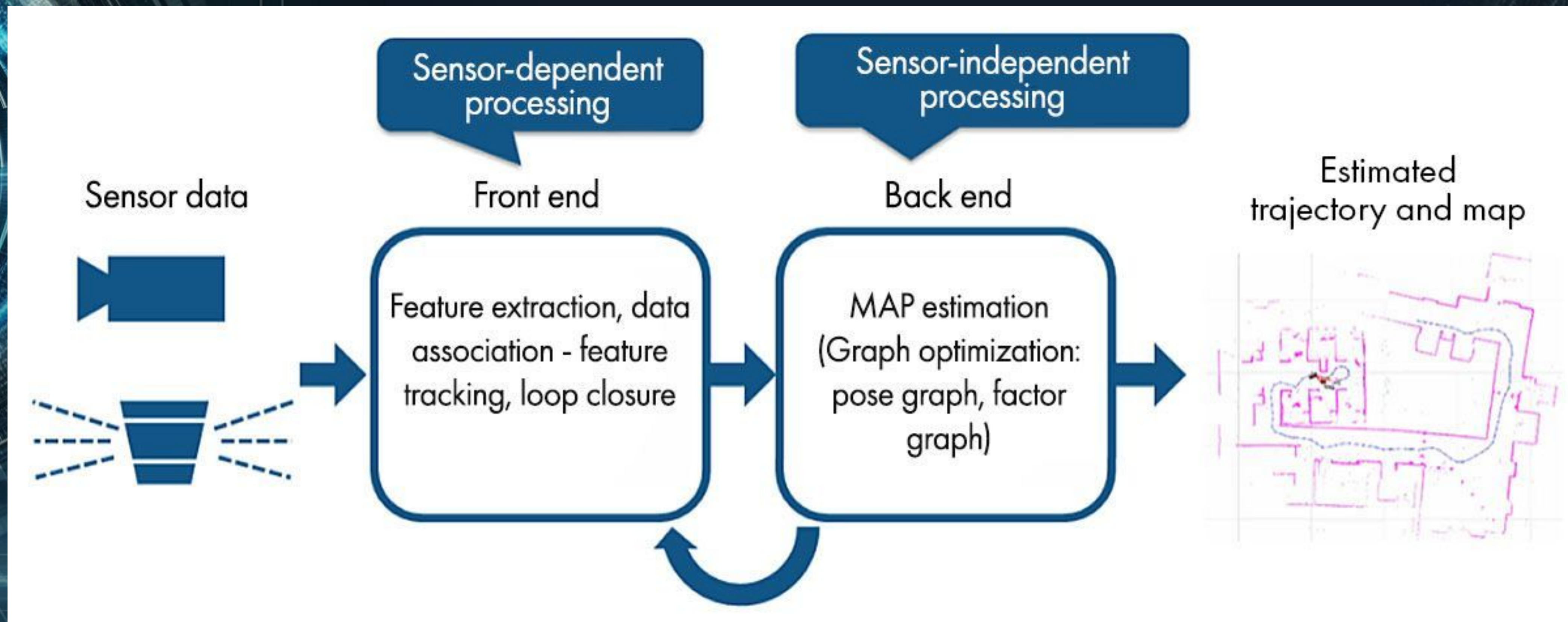
## What is SLAM?



Lidar-based SLAM



The 3 cameras on this sensor mount were used to run SLAM (Lidar at the top was not used for SLAM)

Visual SLAM

# Introduction

## What is SLAM?

# Introduction

## Gaps in Previous SLAM Systems

- **Multiple sensors:** Traditional SLAM relies on additional sensors such as IMU to estimate pose and trajectories.

- **Camera Requirements:** Needed Known Camera Settings (Camera Clibration), Fixed Camera Models: Geometry Prediction

- **Dense SLAM:** Many Systems Couldn't scale to real-time dense 3D mapping

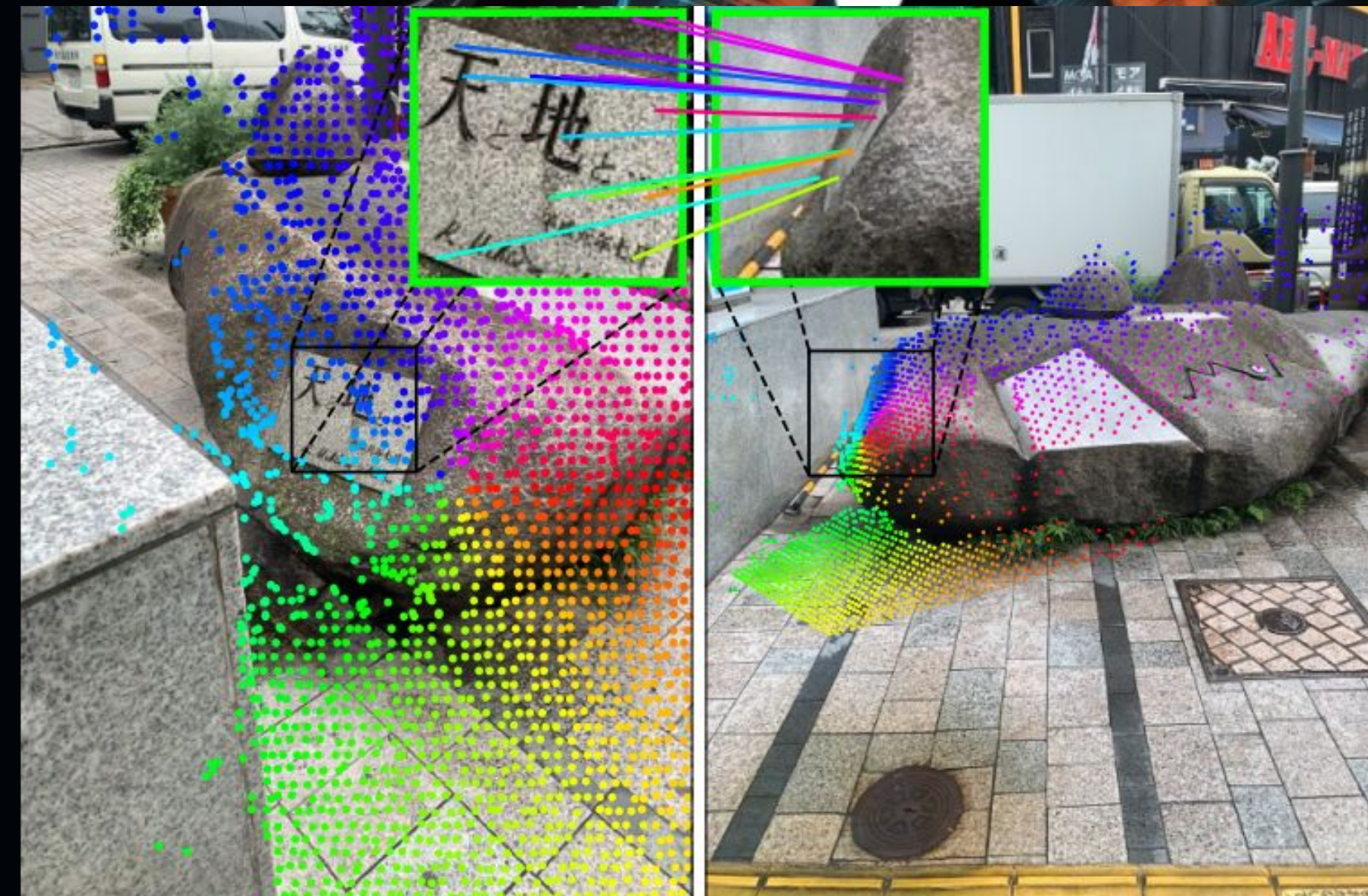- **Runtime:** learning based systems like MASt3r-sfm were offline and slow

EaTemp

# Introduction

## Motivation

- Instead of using traditional sub-steps of matching and trangulation, it uses a prior, MASt3R, for an end-to-end prediction of the pointmap.

**MASt3R**, Multi-view Aggregation for Structure-from-Motion with Transformer, is a deep learning model built on the earlier method DUSt3R.
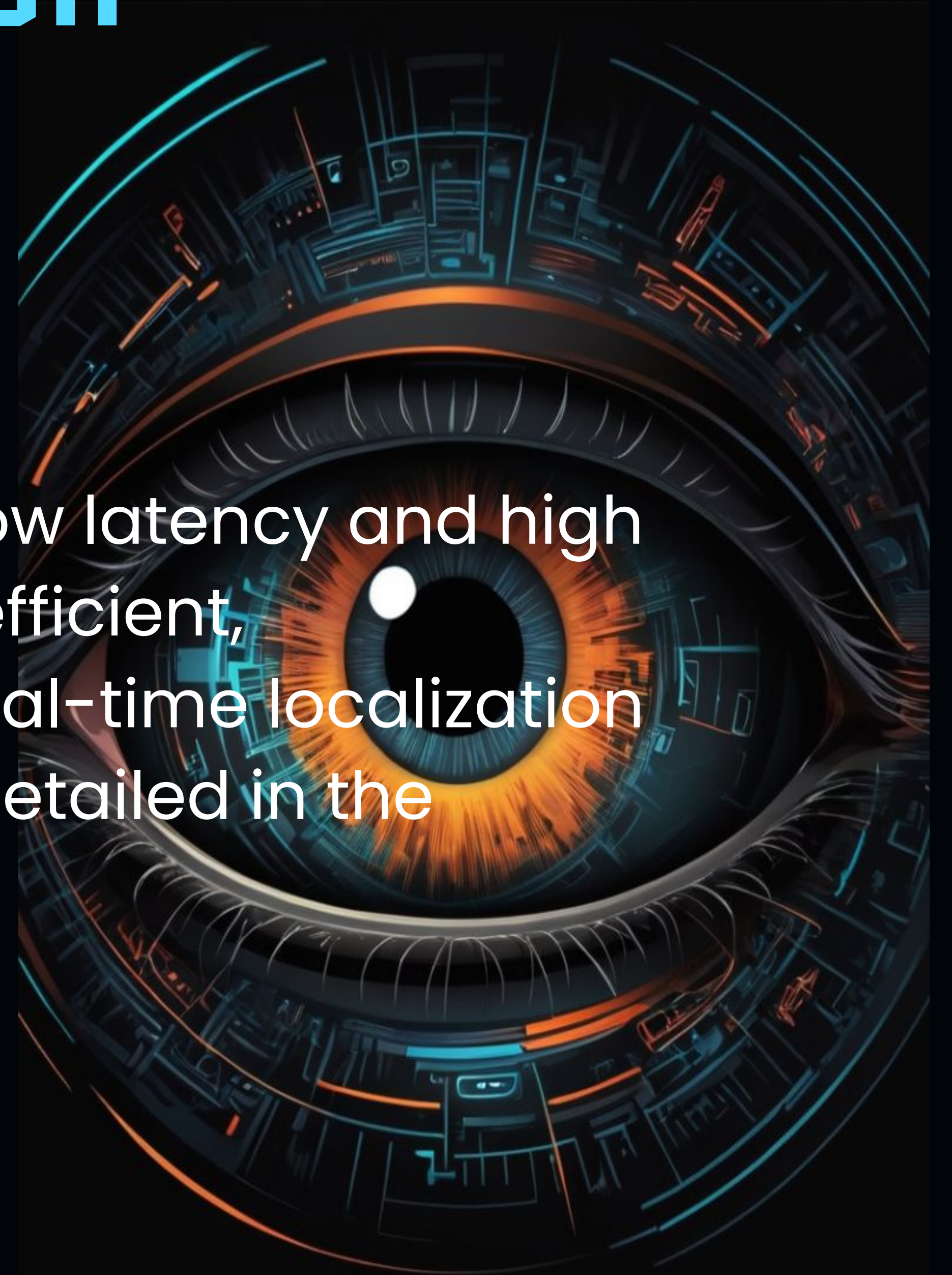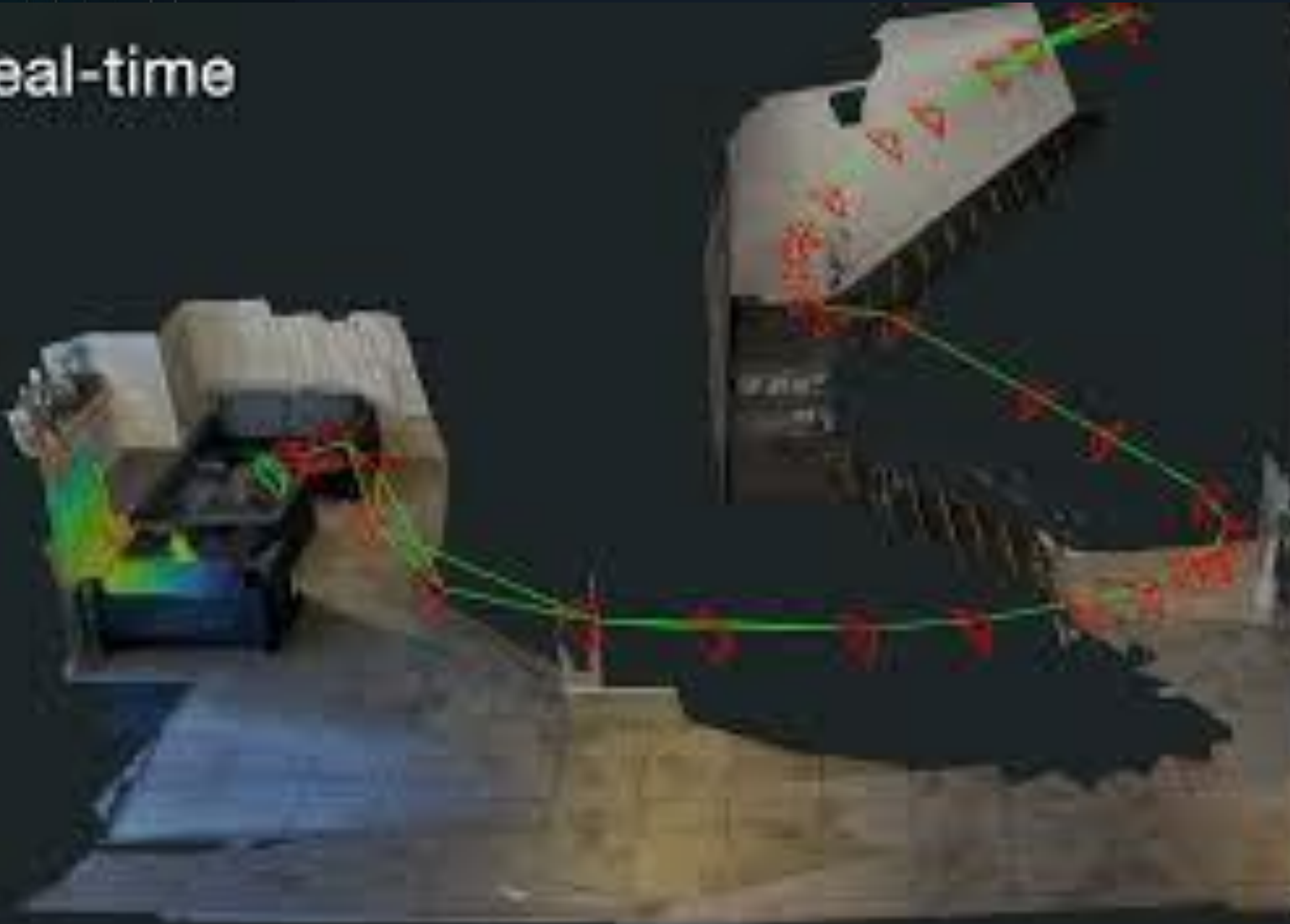
# Introduction

## Motivation

Considering that SLAM requires low latency and high accuracy, the paper introduces efficient, camera-agnostic methods for real-time localization and map estimation, which are detailed in the Methodology section.
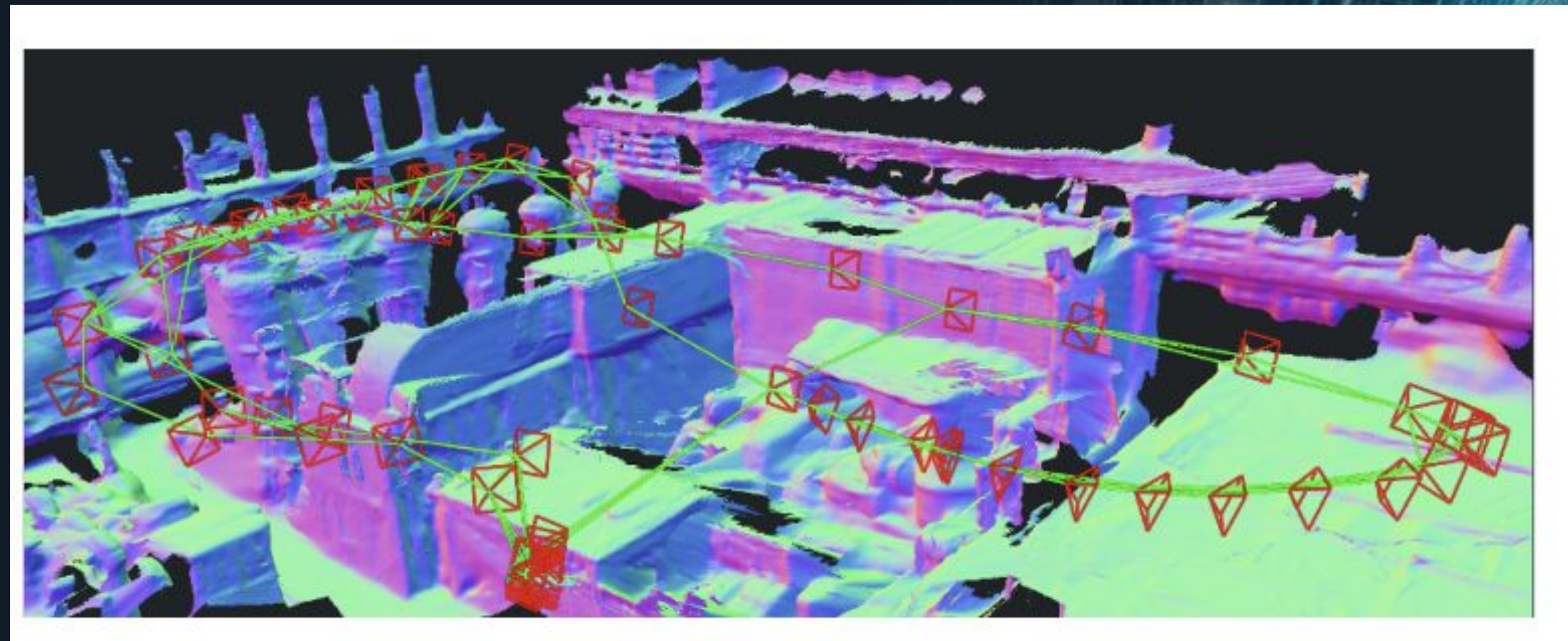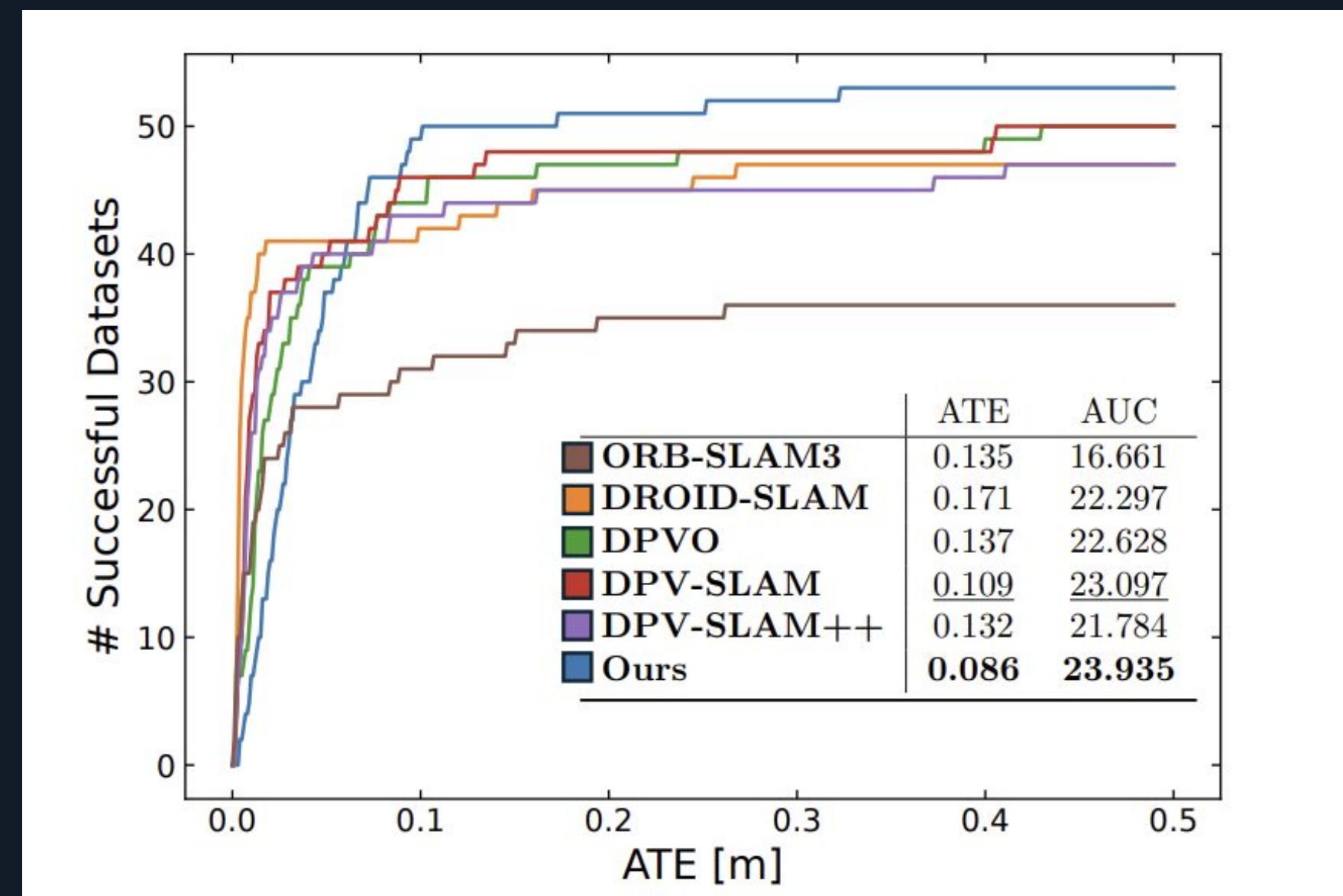
EaTemp

# Results

# Results

## Evaluation metrics for results

- **ATE(Absolute Trajectory Error):** Measures how accurate the estimated camera path is compared to ground truth
- **Chamfer Distance:** Measures how close the reconstructed 3D geometry is to the real world scene
- **Accuracy, completion:** how accurately and completely the scene is reconstructed

## What did they Evaluate?

They tested MASt3R- SLAM on 4 real world datasets TUM RGB-D, 7-Scenes, EuRoC MAV, ETH3D-SLAM



| | ATE | AUC |
|---|---|---|
| ORB-SLAM3 | 0.135 | 16.661 |
| DROID-SLAM | 0.171 | 22.297 |
| DPVO | 0.137 | 22.628 |
| DPV-SLAM | 0.109 | 23.097 |
| DPV-SLAM++ | 0.132 | 21.784 |
| Ours | **0.086** | **23.935** |



This 3D image(EuRoC machine hall) shows the reconstructed point cloud (colored surface) and camera poses (little red pyramids and green lines).

# Results

Shows two consecutive frames (just 1 second apart) with extreme zoom change.
Despite no camera calibration, MASt3R-SLAM reconstructs a clean 3D model of the building.
Colored pyramids show camera positions over time.



Figure 7. Dense uncalibrated SLAM with extreme zoom changes shown by two consecutive keyframes for an outdoor scene.

# Methodology:Preliminaries

$\mathcal{I}^i$

$\mathcal{I}^j$

MASt3R

$$\mathbf{X}_i^i, \mathbf{X}_i^j \in \mathbb{R}^{H \times W \times 3}$$

$$\mathbf{C}_i^i, \mathbf{C}_i^j \in \mathbb{R}^{H \times W \times 1}$$

$$\mathbf{D}_i^i, \mathbf{D}_i^j \in \mathbb{R}^{H \times W \times d}$$

**The forward pass of MASt3R is Fm[Ii, Ij]**

EaTemp

# Methodology:Preliminaries

# Methodology:PointMap Matching

$$m_{i,j} = \mathcal{M}(\mathbf{X}_i^i, \mathbf{X}_i^j, \mathbf{D}_i^i, \mathbf{D}_i^j)$$



Cost Function:

$$\mathbf{p}^* = \arg \min_{\mathbf{p}} \left\| \psi\left([\mathbf{X}_i^i]_{\mathbf{p}}\right) - \psi(\mathbf{x}) \right\|^2$$

There is no closed-form projection.
Assumption: **Each frame has a unique camera center**

EaTemp

# Methodology:PointMap Matching



Since it is not a linear equation, how can we minimize the error? The paper uses Levenberg-Marquardt which is for non-linear least square.

## Cost Function:

$$\mathbf{p}^* = \arg\min_{\mathbf{p}} \left\| \psi\left([\mathbf{X}_i^i]_{\mathbf{p}}\right) - \psi(\mathbf{x}) \right\|^2$$

How to represent the error? The difference of two rays can be converted into the theta by the following equation:

$$\|\psi_1 - \psi_2\|^2 = 2(1 - \cos\theta), \quad \cos\theta = \psi_1^T \psi_2.$$

EaTemp

# Methodology:PointMap Matching

## What is Levenberg-Marquardt?

The primary application of the Levenberg–Marquardt algorithm is in the least-squares curve fitting problem: given a set of $m$ empirical pairs $(x_i, y_i)$ of independent and dependent variables, find the parameters $\beta$ of the model curve $f(x, \beta)$ so that the sum of the squares of the deviations $S(\beta)$ is minimized:

$$\hat{\beta} \in \operatorname{argmin}_\beta S(\beta) \equiv \operatorname{argmin}_\beta \sum_{i=1}^{m} [y_i - f(x_i, \beta)]^2, \text{ which is assumed to be non-empty.}$$

### Add a little increment

$$f(x_i, \beta + \delta) \approx f(x_i, \beta) + \mathbf{J}_i \delta$$

$$\mathbf{J}_i = \frac{\partial f(x_i, \beta)}{\partial \beta}$$

### Calculate the cost function with the increment

Finally , calculate the derivative

$$
\begin{aligned}
S(\beta + \delta) &\approx \|\mathbf{y} - \mathbf{f}(\beta) - \mathbf{J}\delta\|^2 \\
&= [\mathbf{y} - \mathbf{f}(\beta) - \mathbf{J}\delta]^T [\mathbf{y} - \mathbf{f}(\beta) - \mathbf{J}\delta] \\
&= [\mathbf{y} - \mathbf{f}(\beta)]^T [\mathbf{y} - \mathbf{f}(\beta)] - [\mathbf{y} - \mathbf{f}(\beta)]^T \mathbf{J}\delta - (\mathbf{J}\delta)^T [\mathbf{y} - \mathbf{f}(\beta)] + \delta^T \mathbf{J}^T \mathbf{J}\delta \\
&= [\mathbf{y} - \mathbf{f}(\beta)]^T [\mathbf{y} - \mathbf{f}(\beta)] - 2[\mathbf{y} - \mathbf{f}(\beta)]^T \mathbf{J}\delta + \delta^T \mathbf{J}^T \mathbf{J}\delta.
\end{aligned}
$$

$$(\mathbf{J}^T \mathbf{J}) \delta = \mathbf{J}^T [\mathbf{y} - \mathbf{f}(\beta)]$$

EaTemp

# Methodology:PointMap Matching

With a cost function:

$$\min_{\mathbf{p}} \frac{1}{2} \|\mathbf{r}(\mathbf{p})\|^2$$

$$\mathbf{r}(\mathbf{p}^{(k)} + \delta\mathbf{p}) \approx \mathbf{r}(\mathbf{p}^{(k)}) + \mathbf{J}(\mathbf{p}^{(k)})\,\delta\mathbf{p}$$

$$\nabla_{\delta\mathbf{p}} \|\mathbf{r}^{(k)} + \mathbf{J}\,\delta\mathbf{p}\|^2 = 2\mathbf{J}^{\mathsf{T}}\left(\mathbf{r}^{(k)} + \mathbf{J}\,\delta\mathbf{p}\right) = \mathbf{0}$$

$$\delta\mathbf{p} = -\left(\mathbf{J}^{\mathsf{T}}\mathbf{J}\right)^{-1}\mathbf{J}^{\mathsf{T}}\mathbf{r}^{(k)}$$

EaTemp

# Tracking the camera pose

Track the camera's position (pose) with respect to a reference keyframe (the last keyframe) in real-time(relative transformation).

The system relies on the point maps generated from two images to track the camera's movement.
Optimization
Equation for 3D Point Error:

$$E_p = \sum_{m,n \in \mathbf{m}_{f,k}} \left\| \frac{\tilde{\mathbf{X}}_{k,n}^{k} - \mathbf{T}_{kf}\mathbf{X}_{f,m}^{f}}{w(\mathbf{q}_{m,n}, \sigma_p^2)} \right\|_\rho,$$

The error between the predicted 3D point location and the actual position of the camera is minimized to estimate the pose.
The error is weighted by confidence values from MASt3R and robust norms like Huber norm are used to make the system more resilient to outliers.

# Pixel Wise alignment Vs Explicit matching

The paper uses explicit matching of points between the current frame and the keyframe to improve accuracy, especially when there is a larger baseline (greater camera movement between frames).

**Fusing Predictions:**
The system averages multiple pointmap predictions into a single pointmap.

**Error accumulation:**
Tracking errors due to depth inaccuracies accumulate during fusion.

**Impact on Keyframe's Pointmap:**
These errors degrade the keyframe's pointmap.

The degraded keyframe affects the backend map and overall pose estimation, reducing accuracy.

EaTemp

# Ray based error for improved accuracy

Instead of directly using depth values, the system converts 3D points to rays.

Under the assumption that the camera is central, the system treats each point as a ray originating from the camera center.

equation for ray based error:

$$E_r = \sum_{m,n \in \mathbf{m}_{f,k}} \left\| \frac{\psi\left(\tilde{\mathbf{X}}_{k,n}^{k}\right) - \psi\left(\mathbf{T}_{kf}\mathbf{X}_{f,m}^{f}\right)}{w(\mathbf{q}_{m,n}, \sigma_r^2)} \right\|_{\rho}$$

Advantage: Less sensitive to depth inaccuracies: Angular error between rays is more robust to depth discrepancies than 3D point error.

# Pose Update via Gauss-Newton Optimization

**The system computes the update to the pose by solving the linear system**

$$(\mathbf{J}^T \mathbf{W} \mathbf{J}) \, \boldsymbol{\tau} = -\mathbf{J}^T \mathbf{W} \mathbf{r}$$

**Once the update is computed, the pose is updated as**

$$\mathbf{T}_{kf} \leftarrow \boldsymbol{\tau} \oplus \mathbf{T}_{kf}.$$

**After updating the pose, the canonical pointmap is updated to improve accuracy**

**New points merged via confidence-weighted average**

$$\tilde{\mathbf{X}}_k^k \leftarrow \frac{\tilde{\mathbf{C}}_k^k \tilde{\mathbf{X}}_k^k + \mathbf{C}_f^k \left( \mathbf{T}_{kf} \mathbf{X}_f^k \right)}{\tilde{\mathbf{C}}_k^k + \mathbf{C}_f^k}, \tilde{\mathbf{C}}_k^k \leftarrow \tilde{\mathbf{C}}_k^k + \mathbf{C}_f^k.$$
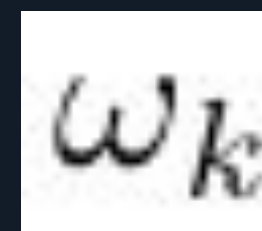
# Methodology:Graph Construction and Loop Closure

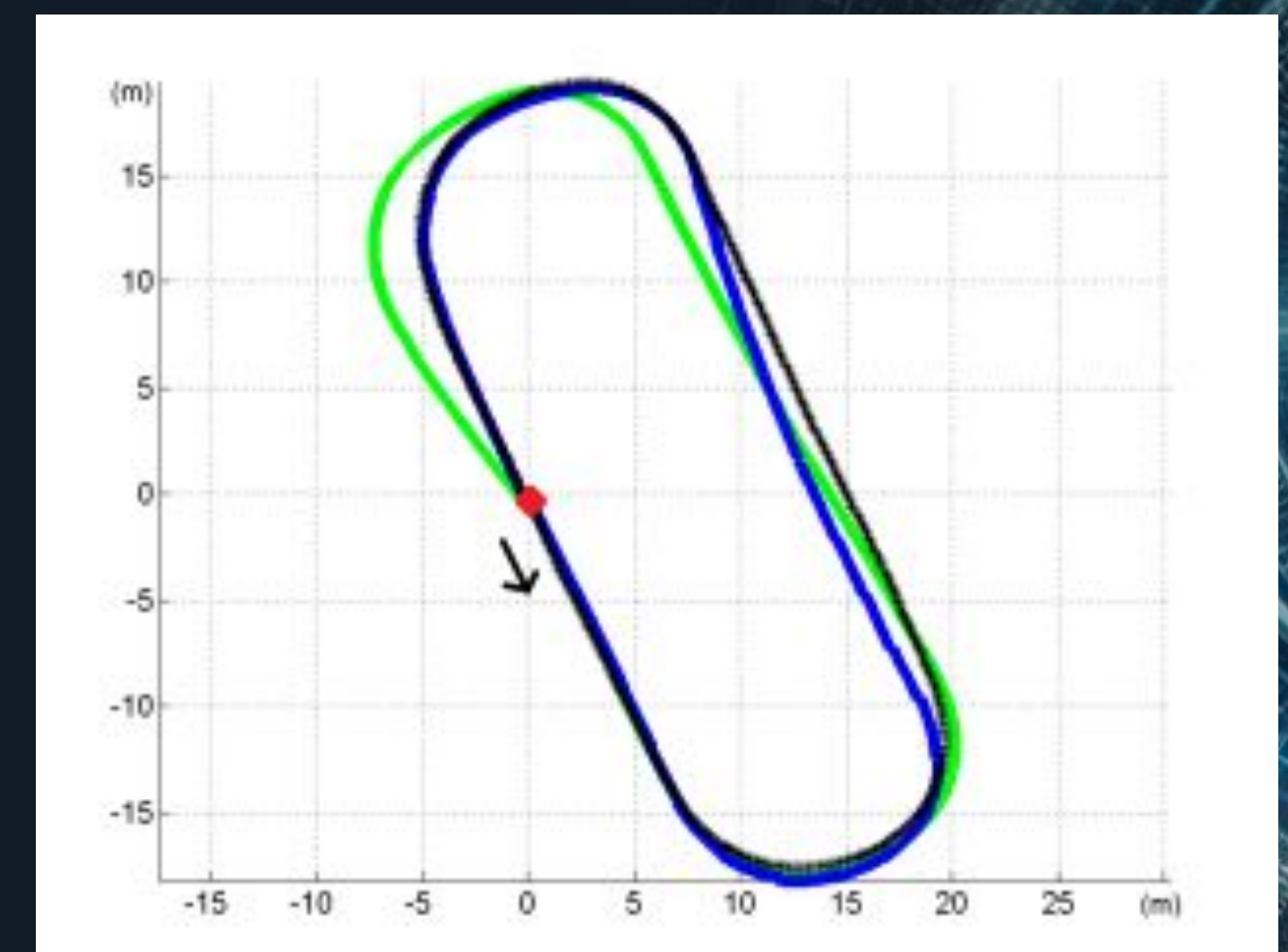**When should a new keyframe be added?**

$\mathcal{K}_{i-1}$

Corresponding(such as feature matches)

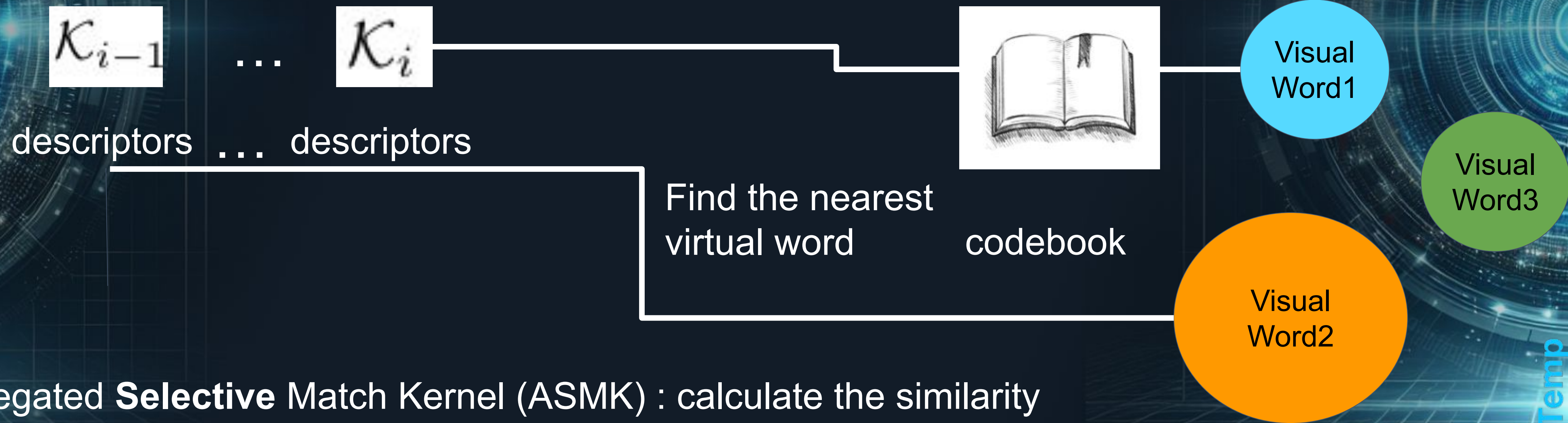$\mathcal{K}_i$

reference frame

$\omega_k$

New frame

It's time to add a new keyframe!



Drifting

# Methodology:Graph Construction and Loop Closure

**A loop closure is used to decrease drifting. But how to know when the loop is closure?**

$$\mathcal{K}_{i-1} \quad \ldots \quad \mathcal{K}_i$$

descriptors ... descriptors

Find the nearest virtual word          codebook

Visual Word1
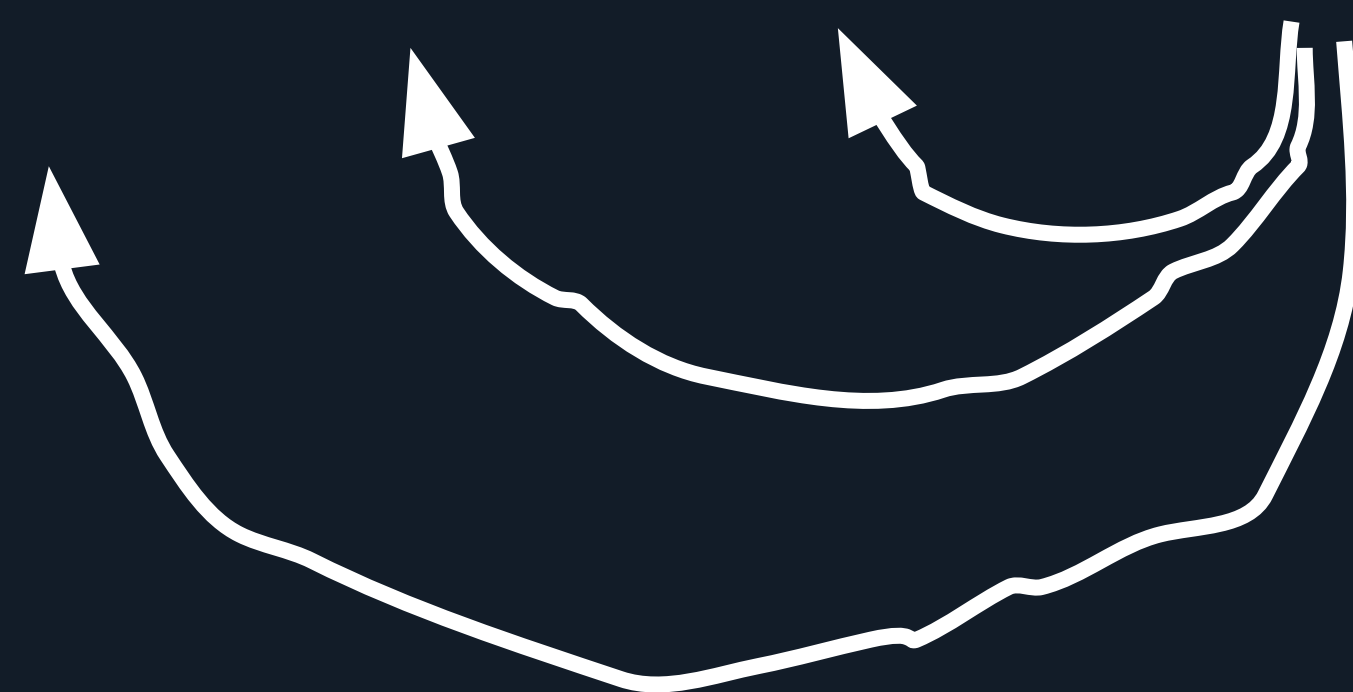
Visual Word3

Visual Word2

Aggregated **Selective** Match Kernel (ASMK) : calculate the similarity score by focusing on local descriptors that share a common visual word and have sufficiently similar orientation

# Methodology:Graph Construction and Loop Closure
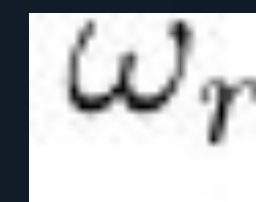
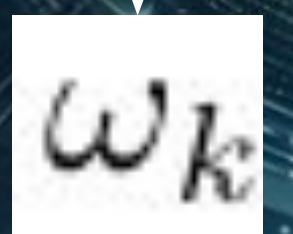**A loop closure is used to decrease drifting. But how to know when the loop is closure?**



*Query*

old ——————————→ latest

MASt3R

ASMK scores:
If the score of a pair is higher than the threshold

$w_r$

$w_k$

**Loop Closed**

# Backend Optimisation

## RAY-BASED BACKEND OPTIMIZATION

**1. BUILD POSE GRAPH**

matches

keyframes

**2. FORMULATE ERROR ON RAYS**

$$\frac{\left|\psi\bar{X})-(\bar{X}_j)\right|}{w(q_{mn},\sigma_r^2)}$$

**3. FIX THE GAUGE**

$$-\frac{\overset{\curvearrowright}{S^x}}{pose}$$

**3. FIX THE GAUGE**

**4. BUILD THE HESSIAN**

$$\begin{bmatrix} & 0 & \cdots & \cdots & 0 \\ 14 & & 0 & \cdots & 0 \\ \vdots & 0 & & \vdots & \vdots \\ \vdots & \vdots & \cdots & & 0 \end{bmatrix}$$

$$7N \times 7N$$

**5. SOLVE WITH GAUSS-NEWTON**

$$Hx = -b$$

**6. UPDATE ALL POSES**

**7. FIX DRIFT**

---

Minimize global consistency across all poses and geometry

Used Gauss-Newton optimization(a second-order method)
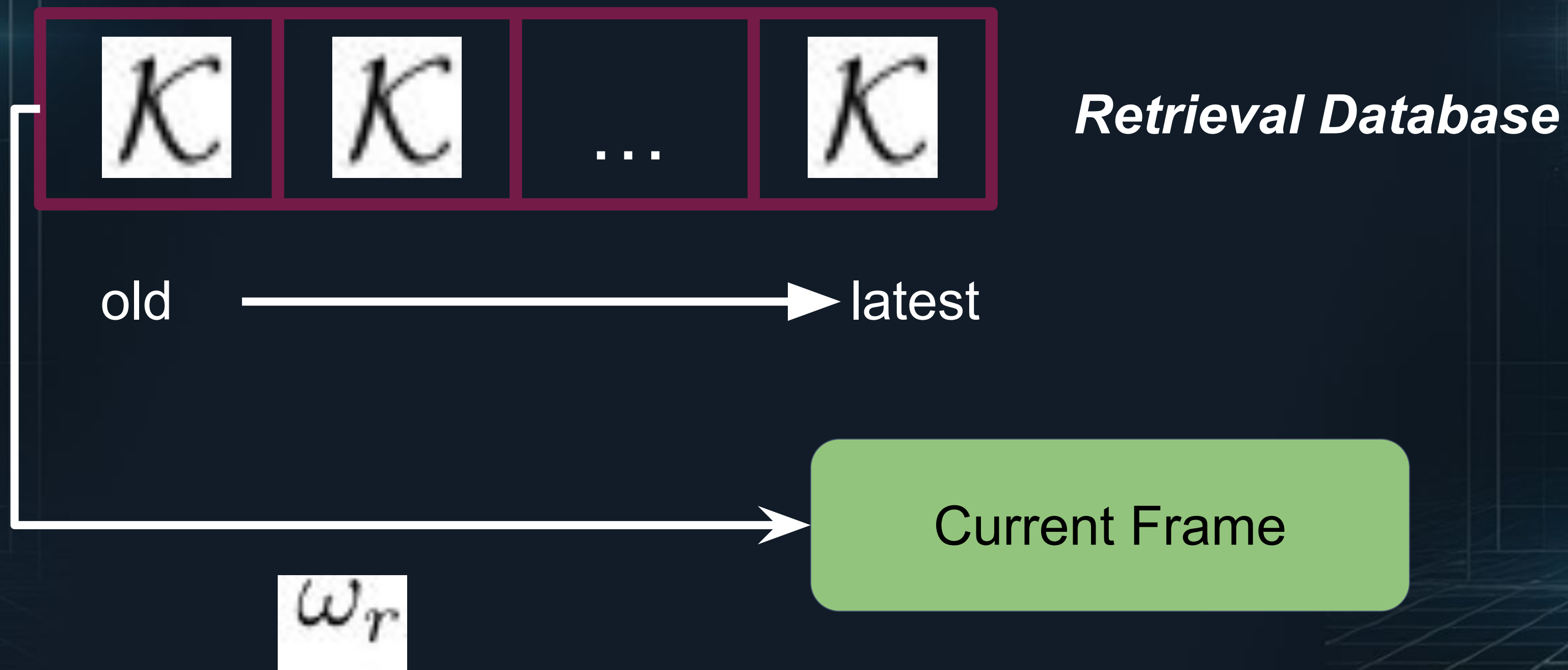
Why

First-order requires rescaling after every iteration

Gauge Freedom - First 7 DOF

3 for Translation
3 for Rotation
1 for scaling

# Methodology:Relocalization

**When the system loses tracking due to an insufficient number of matches, relocalisation mode is triggered.**



*Retrieval Database*

old ⟶ latest

Current Frame

$\omega_r$

# Methodology:Known Calibration

**1.Constraining the Pointmap with Calibration: Calculate ray with known intrinsic parameters**

**2.Pixel-Space Residuals Instead of Ray-Space**

Instead of calculating in ray space, the residual can be calculated in pixel space

$$E_\Pi = \sum_{i,j \in \mathcal{E}} \sum_{m,n \in \mathbf{m}_{i,j}} \left\| \frac{\mathbf{p}_{i,m}^i - \Pi\left(\mathbf{T}_{ij}\tilde{\mathbf{X}}_{j,n}^j\right)}{w(\mathbf{q}_{m,n}, \sigma_\Pi^2)} \right\|_\rho ,$$

Instead of using confidence, we can use other advance methodology to pointmap fusion such as Kalman Filter)

Train MASt3R with more datasets

# Improvement

Thank you for watching

EaTemp