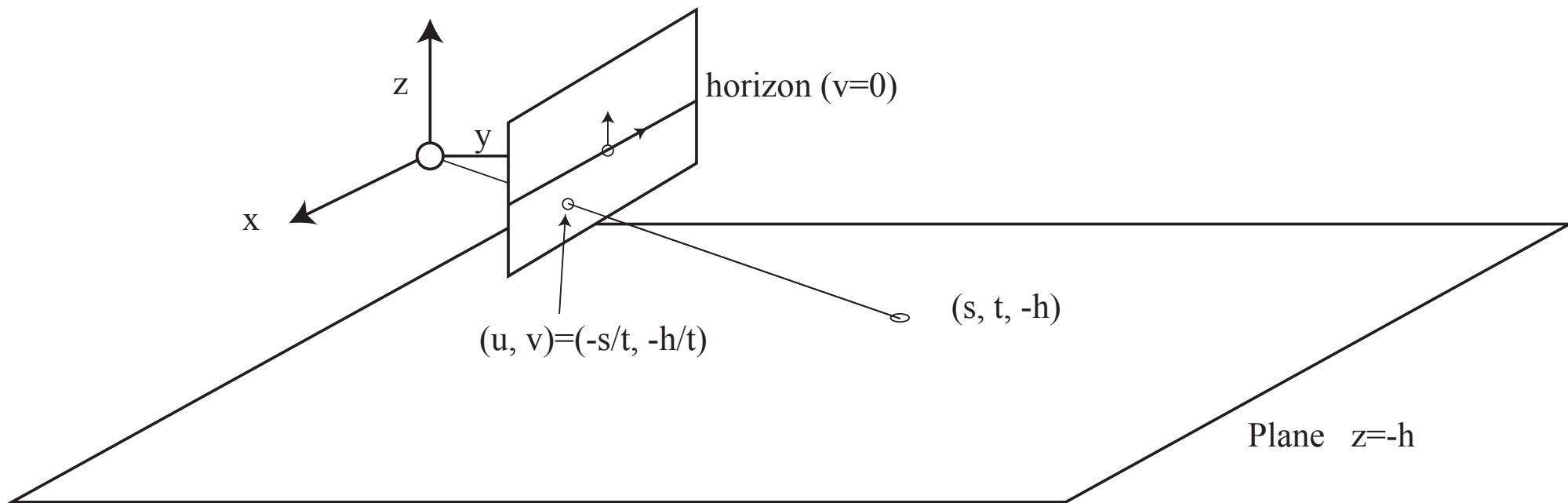


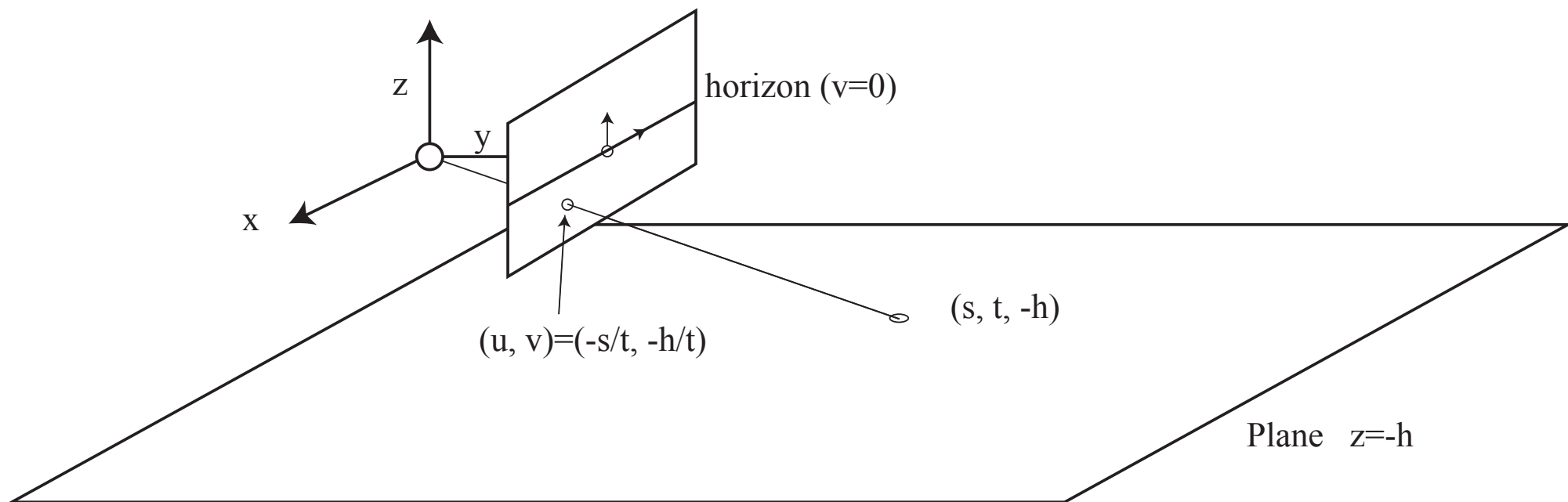
Direct Slam

D.A. Forsyth, UIUC

Basic direct method

- Imagine a camera at a fixed height
 - moving rigidly over a textured ground plane
 - bottom half of image is distorted ground plane texture
 - Q: when camera moves, how does distortion change?





$$\begin{pmatrix} U/W \\ V/W \\ 1 \end{pmatrix} \equiv \begin{pmatrix} U \\ V \\ W \end{pmatrix} = \begin{pmatrix} -1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} s \\ t \\ -h \end{pmatrix}$$

\mathcal{C}
 \downarrow

Image coordinates

Plane texture coords

How the image distorts

- Camera moves relative to ground plane

$$\begin{array}{c}
 \begin{pmatrix} U \\ V \\ W \end{pmatrix}_2 \\
 \text{New image coords}
 \end{array}
 =
 \begin{array}{c}
 \begin{pmatrix} -1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \\
 \text{Motion}
 \end{array}
 \begin{array}{c}
 \begin{pmatrix} \mathcal{R}_{2D} & \mathbf{T} \\ \mathbf{0} & 1 \end{pmatrix} \\
 \text{Motion}
 \end{array}
 \begin{array}{c}
 \begin{pmatrix} s \\ t \\ -h \end{pmatrix} \\
 \text{Original texture}
 \end{array}$$

- And we get an image transformation

$$\begin{array}{c}
 \begin{pmatrix} U \\ V \\ W \end{pmatrix}_2 \\
 \text{New image coords}
 \end{array}
 =
 \begin{array}{c}
 \begin{pmatrix} -1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \\
 \text{Motion}
 \end{array}
 \begin{array}{c}
 \begin{pmatrix} \mathcal{R}_{2D} & \mathbf{T} \\ \mathbf{0} & 1 \end{pmatrix} \\
 \text{Motion}
 \end{array}
 \left[\begin{array}{c}
 \begin{pmatrix} -1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \\
 \text{Motion}
 \end{array} \right]^{-1}
 \begin{array}{c}
 \begin{pmatrix} U \\ V \\ W \end{pmatrix}_1 \\
 \text{Original Image}
 \end{array}$$

Simplest approach

- Assume we *know* the camera height and intrinsics
- We can reconstruct the ground plane
 - apply the inverse of the given map
- Now if camera translates, ground plane translates
 - if camera rotates, ground plane rotates

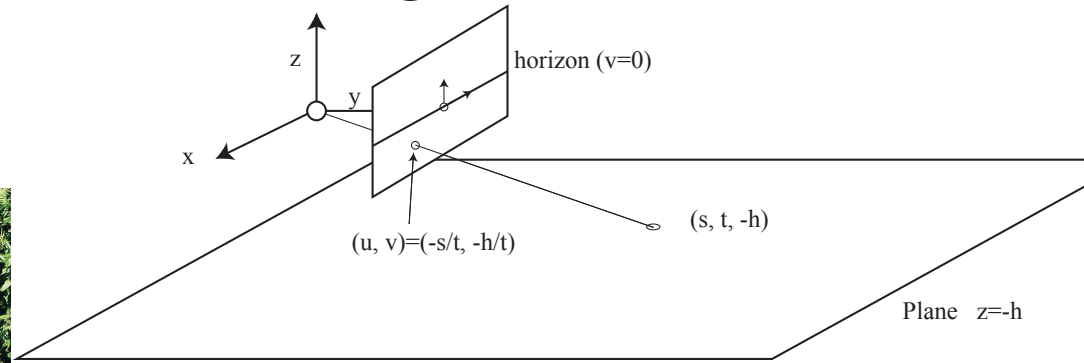
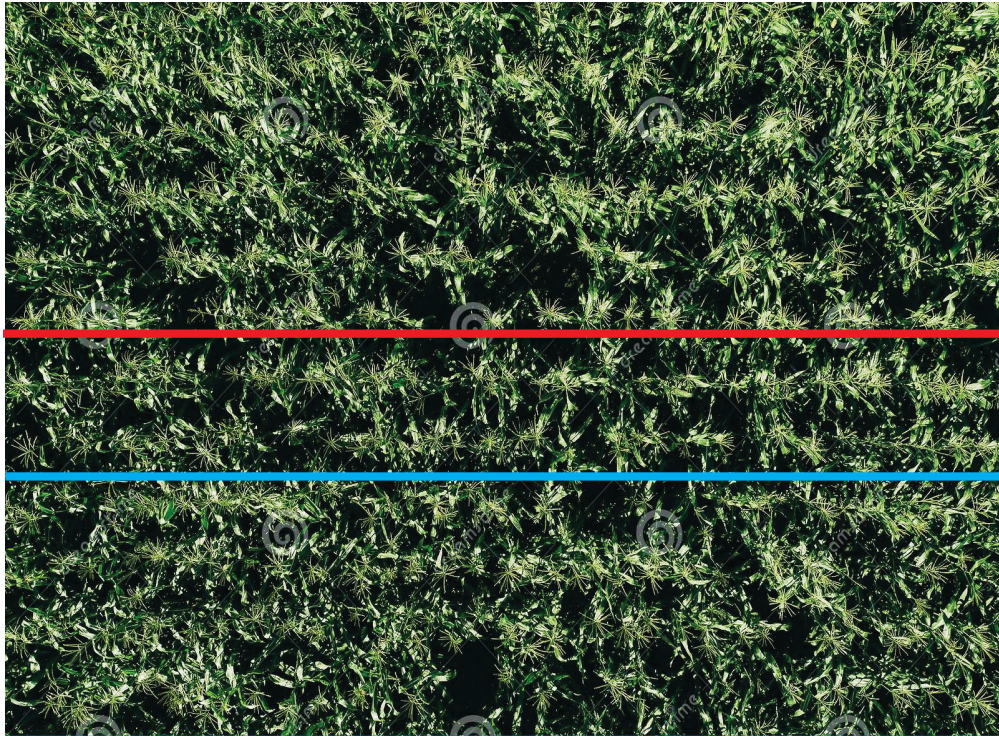


Camera image



Overhead view of ground plane

Ground plane for translating camera

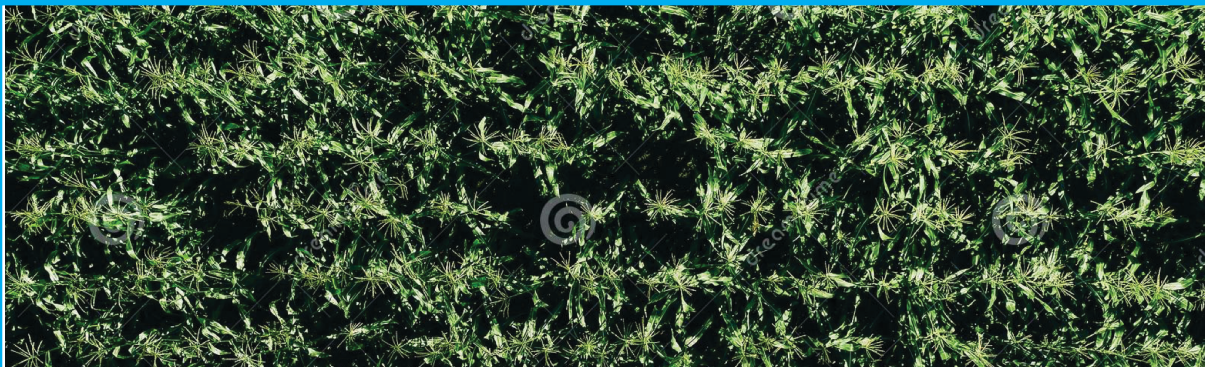


Red camera ($i+1$) sees forward of this



Blue camera (i) sees forward of this

The two ground planes...



can be aligned and this reveals translation



How to register?

- Alternatives
 - find interest points; use registration methods, above
 - direct method (now)

Direct methods - I

- Notice that there is rotation, translation so that

$$G_{i+1}(\mathcal{R}\mathbf{x} + \mathbf{t}) \approx G_i(\mathbf{x})$$

- so we could try to minimize

$$\sum_{\mathbf{x} \in \text{relevant pixels}} [G_{i+1}(\mathcal{R}\mathbf{x} + \mathbf{t}) - G_i(\mathbf{x})]^2$$

- generally, this is a photometric consistency constraint

Direct methods - II

- How to minimize?
 - if motion is small, minimization is easier
 - eg translation

Cost is:

$$[G_{i+1}(x + \delta x, y + \delta y) - G_i(x, y)]^2$$

Approximation:

$$[G_{i+1}(x + \delta x, y + \delta y) - G_i(x, y)]^2 \approx \left[G_{i+1}(x, y) + (\nabla G_{i+1})^T \begin{bmatrix} \delta x \\ \delta y \end{bmatrix} - G_i(x, y) \right]^2$$

And can solve this by least squares

Issues

- Minimize how?
 - Typically, Newton's method or a variant
- What about robustness?
 - use an m-estimator, as in IRLS
- How to initialize?
 - you could use interest points...
 - but if movements are small you might not need to
- Why?
 - massively improved estimates of rotation and translation IF you can min
 - because very large numbers of points contribute
- Generalize
 - camera moving in 3D and viewing a plane - easy, from drawing
 - 3D world - more interesting

SLAM with depth measurements

- RGB-D or stereo (sketch)
 - align depth map in view 2 with that of view 1
 - using rotation, translation, m-estimator+IRLS
 - wrinkle - use intensity as well as depth
 - keep
 - aligned depths (prune redundancies) for mapping
 - transformation (for localization)
- Key question
 - How do we manage computational cost?

SLAM with depth measurements - II

- Strategies:
 - Find image/depth interest points and register those
 - advantage: we know how to do this, straightforward, fast
 - possible disadvantage: ignoring a lot of information
 - Direct method: minimize photometric/depth alignment cost at every point
 - advantage: we know how to do this, all points contribute, accurate
 - disadvantage: much more expensive

Semi-direct methods

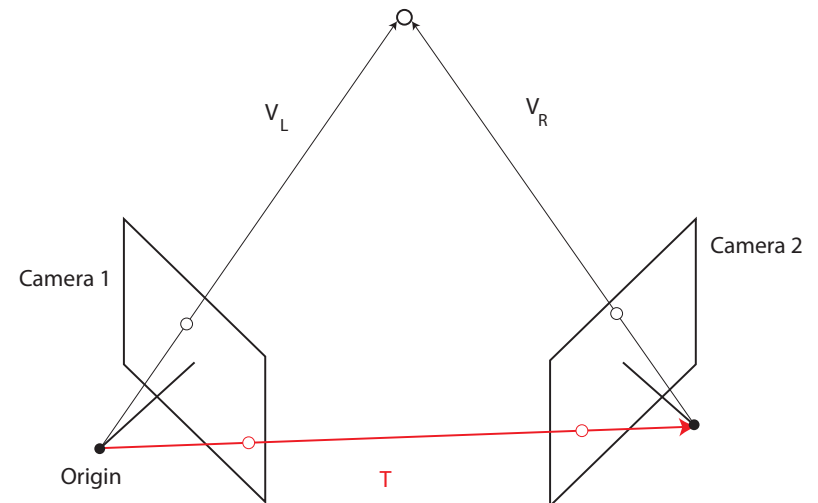
- Problem:
 - a direct method means you must touch many image/depth measurements
 - accurate, but likely slow
- Idea:
 - find interest points
 - use interest point AND its neighborhood
 - computing photometric consistency for neighborhood

Semi-direct visual odometry and mapping

- Monocular cameras
 - semi-direct
 - for the moment, consider interest points in the image
 - worry about neighborhoods in a moment
 - to predict image positions, we need depths
 - associate a depth with each interest point, and a depth estimate
 - for that point in each frame
 - stick a filter on this

SVO (semi-direct visual odometry)

- Must estimate camera motion from interest points
 - assume we have a depth associated with i 'th interest point
 - update when camera moves
 - start
 - depth from prior OR
 - stereo in first two frames OR
 - depth from single image estimate



SVO (semi-direct visual odometry)

- Steps:
 - estimate camera motion from correspondences using photometric error
 - assume constant depth at each patch
 - now move patches in 2D to improve photometric error
 - now adjust 3D configuration of points and camera motion

Estimate camera motion from correspondences using photometric error

Recall: we know how to make these patches

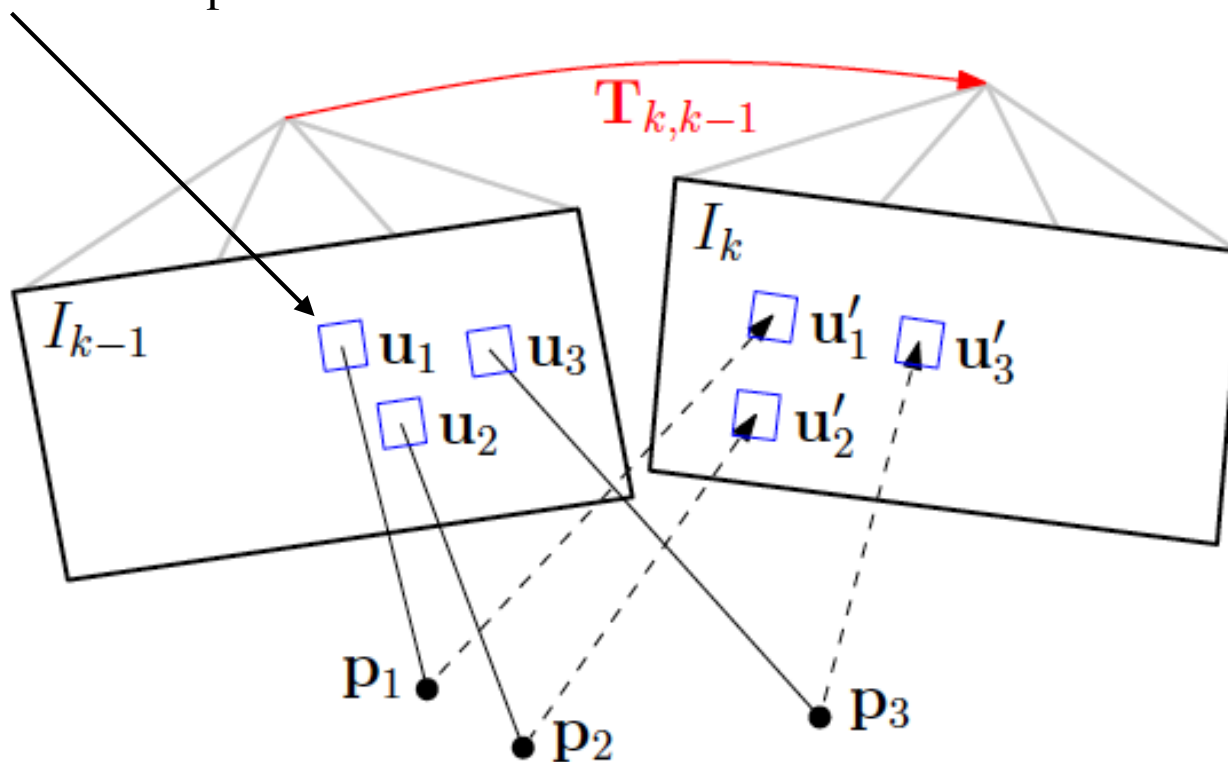


Fig. 2: Changing the relative pose $T_{k,k-1}$ between the current and the previous frame implicitly moves the position of the reprojected points in the new image u'_i . Sparse image alignment seeks to find $T_{k,k-1}$ that minimizes the photometric difference between image patches corresponding to the same 3D point (blue squares). Note, in all figures, the parameters to optimize are drawn in red and the optimization cost is highlighted in blue.

now move patches in 2D to improve photometric error

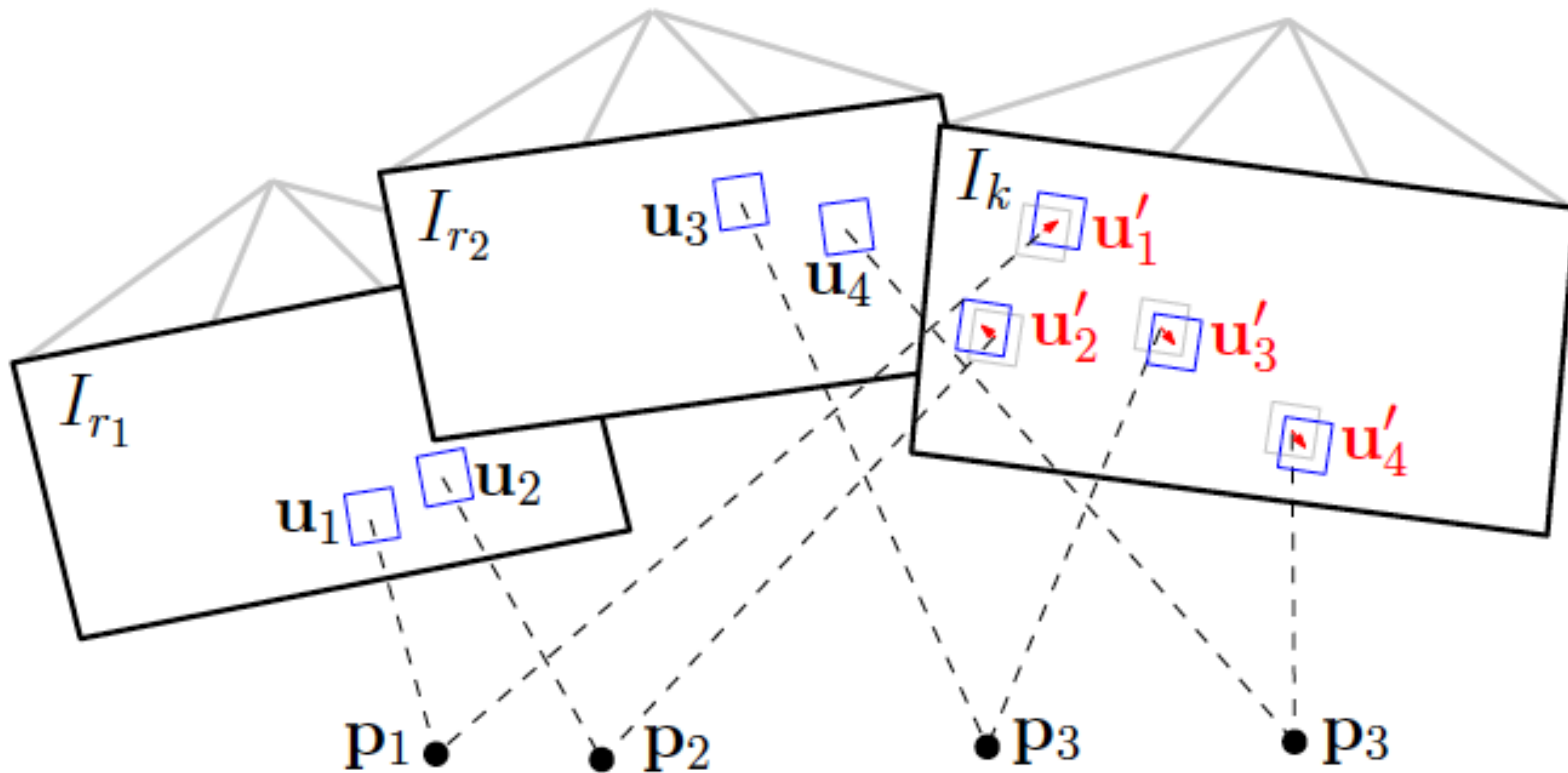


Fig. 3: Due to inaccuracies in the 3D point and camera pose estimation, the photometric error between corresponding patches (blue squares) in the current frame and previous keyframes r_i can further be minimised by optimising the 2D position of each patch individually.

now adjust 3D configuration of points and camera motion

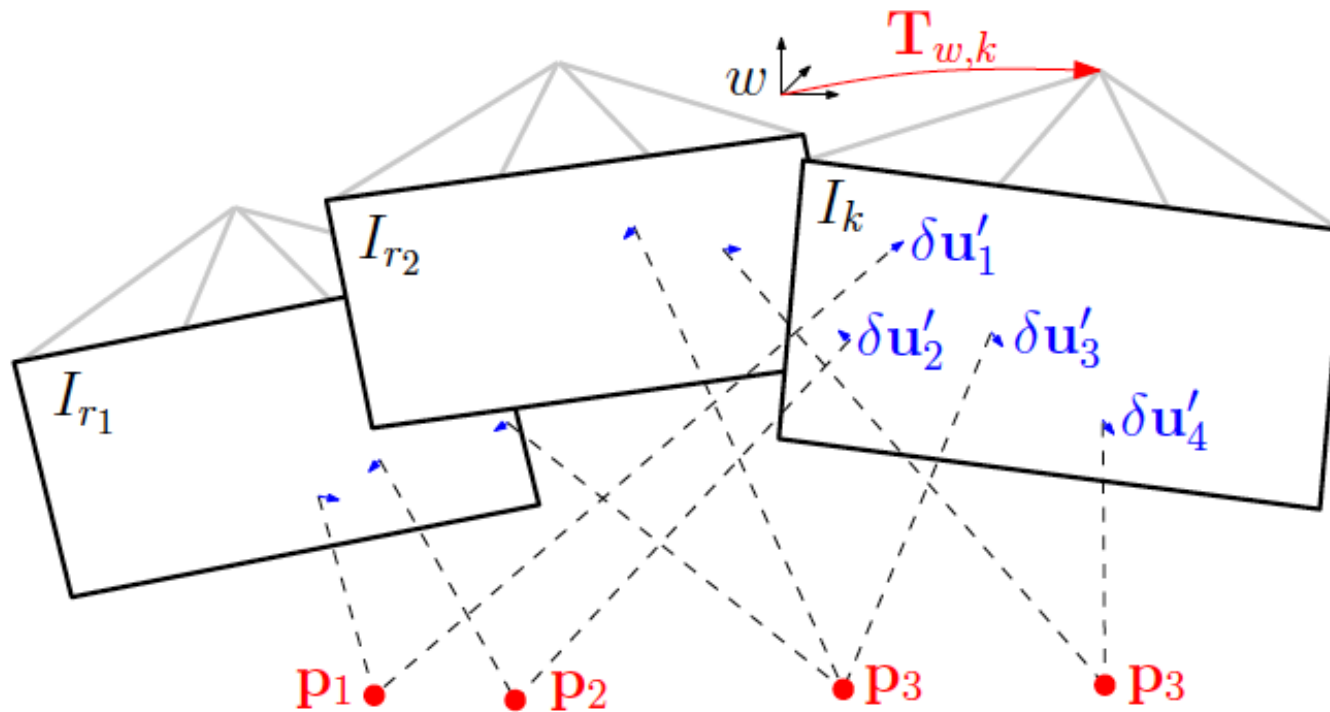


Fig. 4: In the last motion estimation step, the camera pose and the structure (3D points) are optimized to minimize the reprojection error that has been established during the previous feature-alignment step.

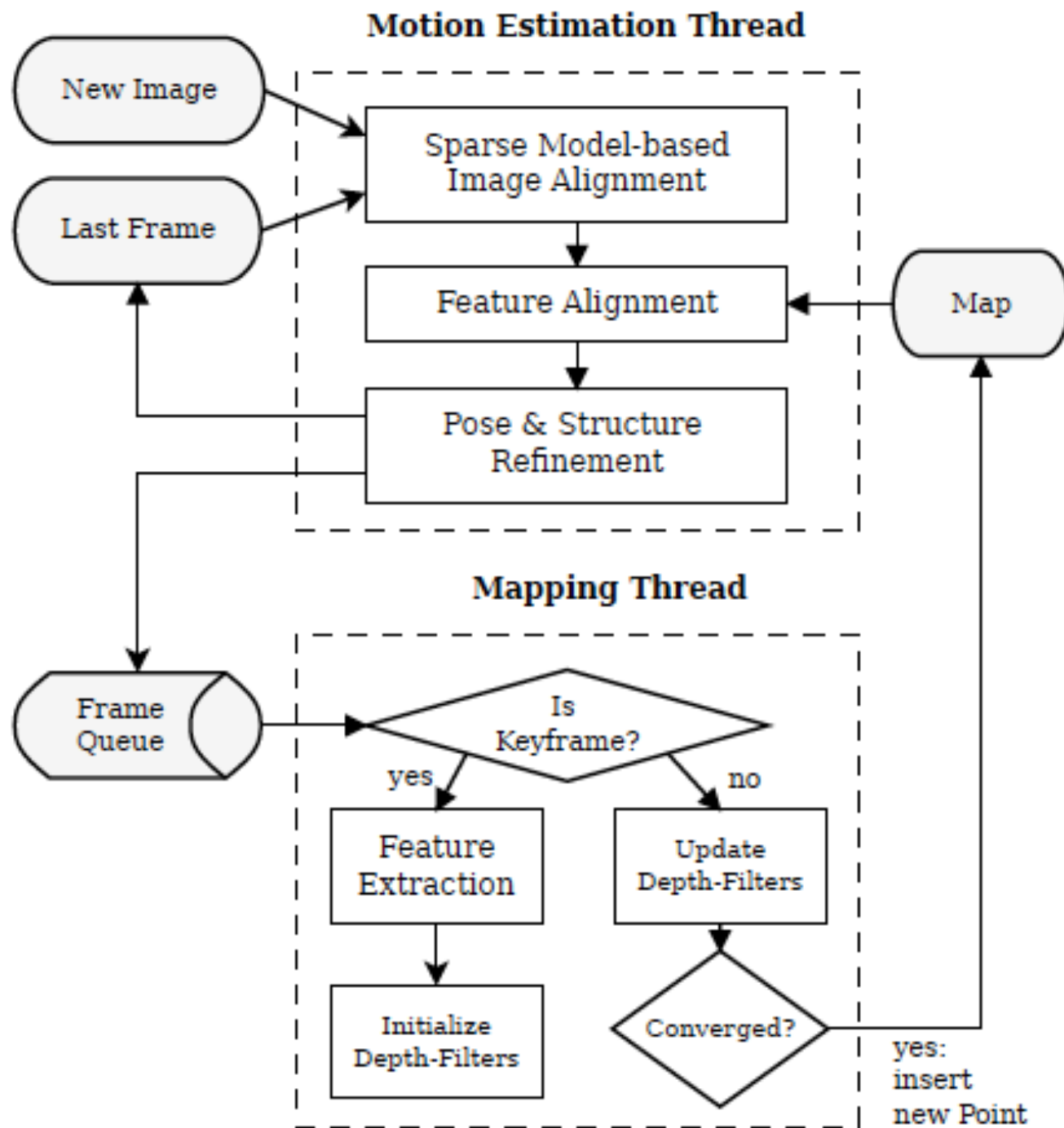


Fig. 1: Tracking and mapping pipeline

Quick and efficient

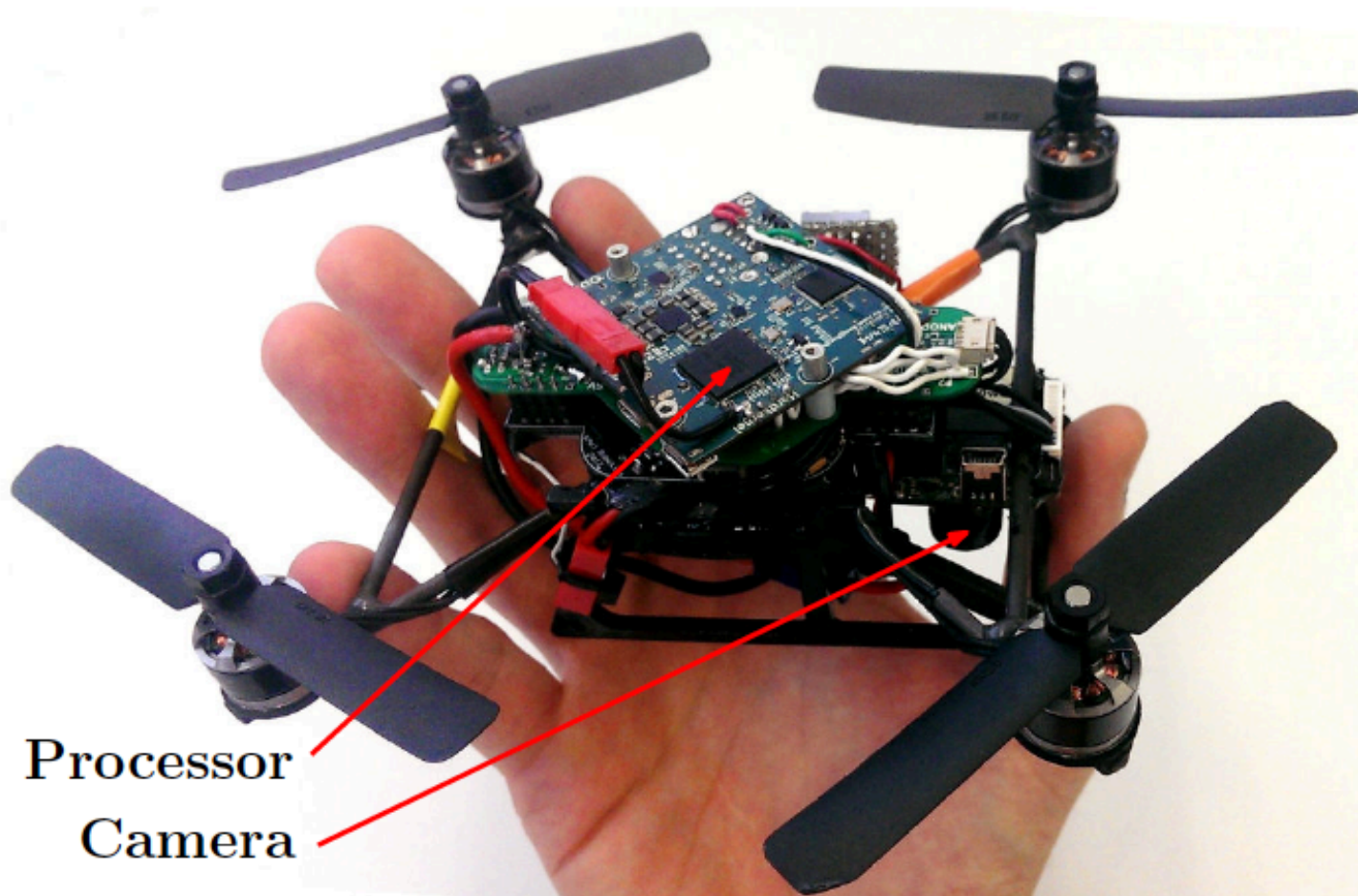


Fig. 17: “Nano+” by KMeI Robotics, customized with embedded processor and downward-looking camera. SVO runs at 55 frames per second on the platform and is used for stabilization and control.

LSD-SLAM

- Essential point:
 - careful use of keyframes speeds things up a lot
 - Monocular cameras
 - direct
 - to predict image positions, we need depths
 - Recall:
 - keyframes are fine
 - Discovery:
 - keyframe depths are enough
 - pose graph:
 - key frames (nodes) linked by transforms (edges)

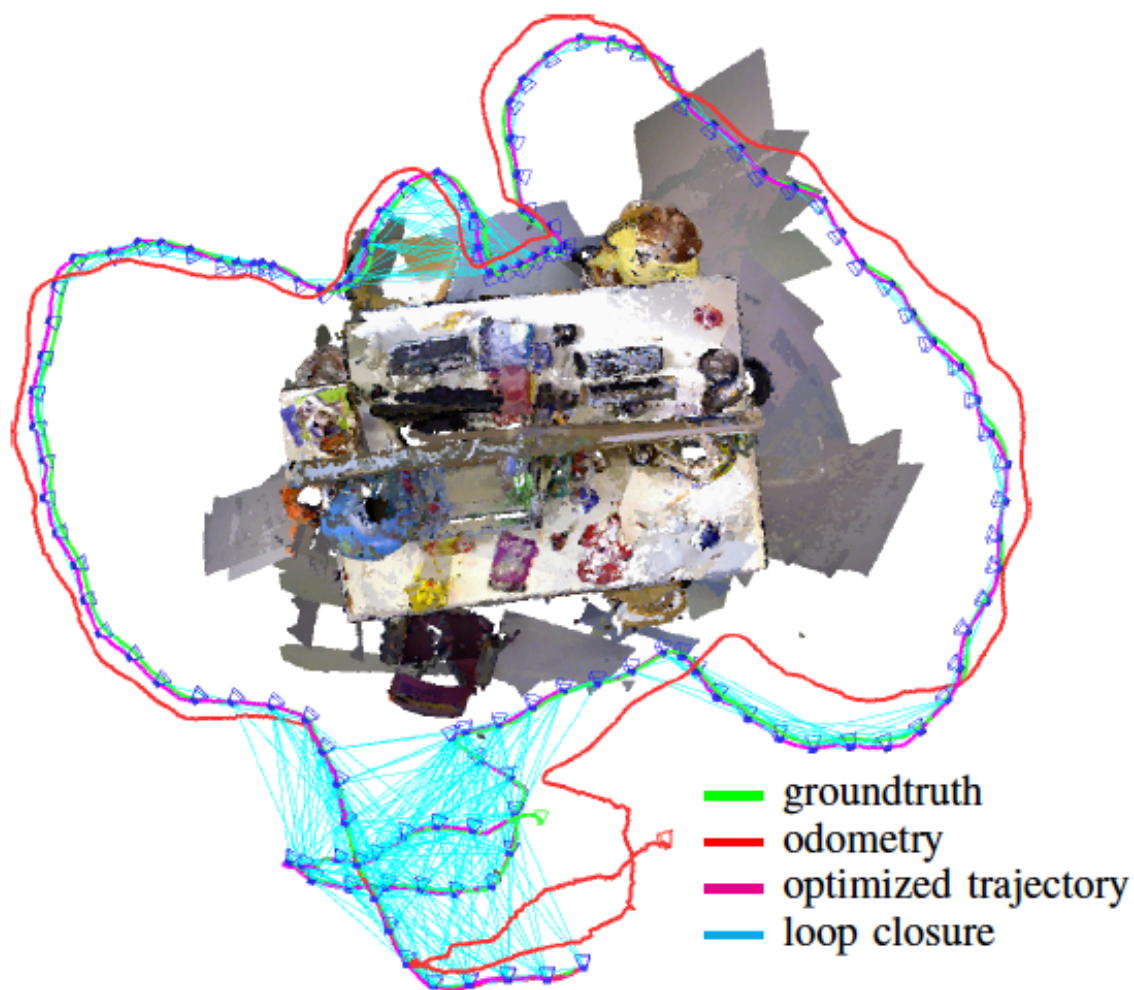


Fig. 1: We propose a dense SLAM method for RGB-D cameras that uses keyframes and an entropy-based loop closure detection to eliminate drift. The figure shows the groundtruth, frame-to-keyframe odometry, and the optimized trajectory for the fr3/office dataset.

Essential steps

- Compare new frame to keyframe
 - which has known depths
 - compute R, T, from
 - photometric error and depth error
 - recall:
 - depth \rightarrow image match \rightarrow photometric error
 - (could adjust depths in keyframe using R,T, photometric error)
- Keyframe selection
 - even sampling OR
 - entropy of R,T from last keyframe to current frame $j = H_j$
 - look at H_j/H_1
 - ie am I getting bad at computing motion?
 -

Essential steps

- Loop Closure
 - match keyframes
- Map
 - pose graph
- Start
 - how do I get depth for first keyframe?
 - doesn't seem to matter - random initialization
 - as long as you refine the depths
 - (roughly) stereo matching yields depth

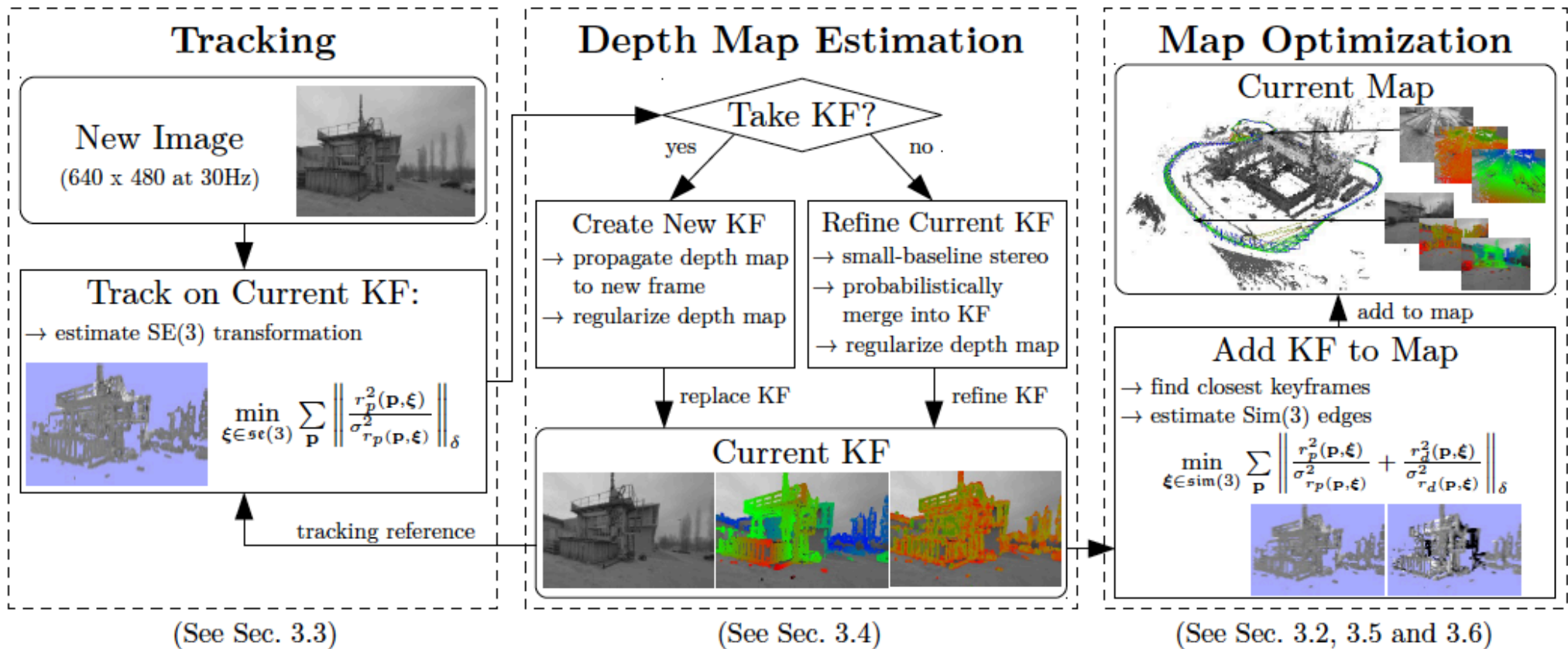


Fig. 3: Overview over the complete LSD-SLAM algorithm.

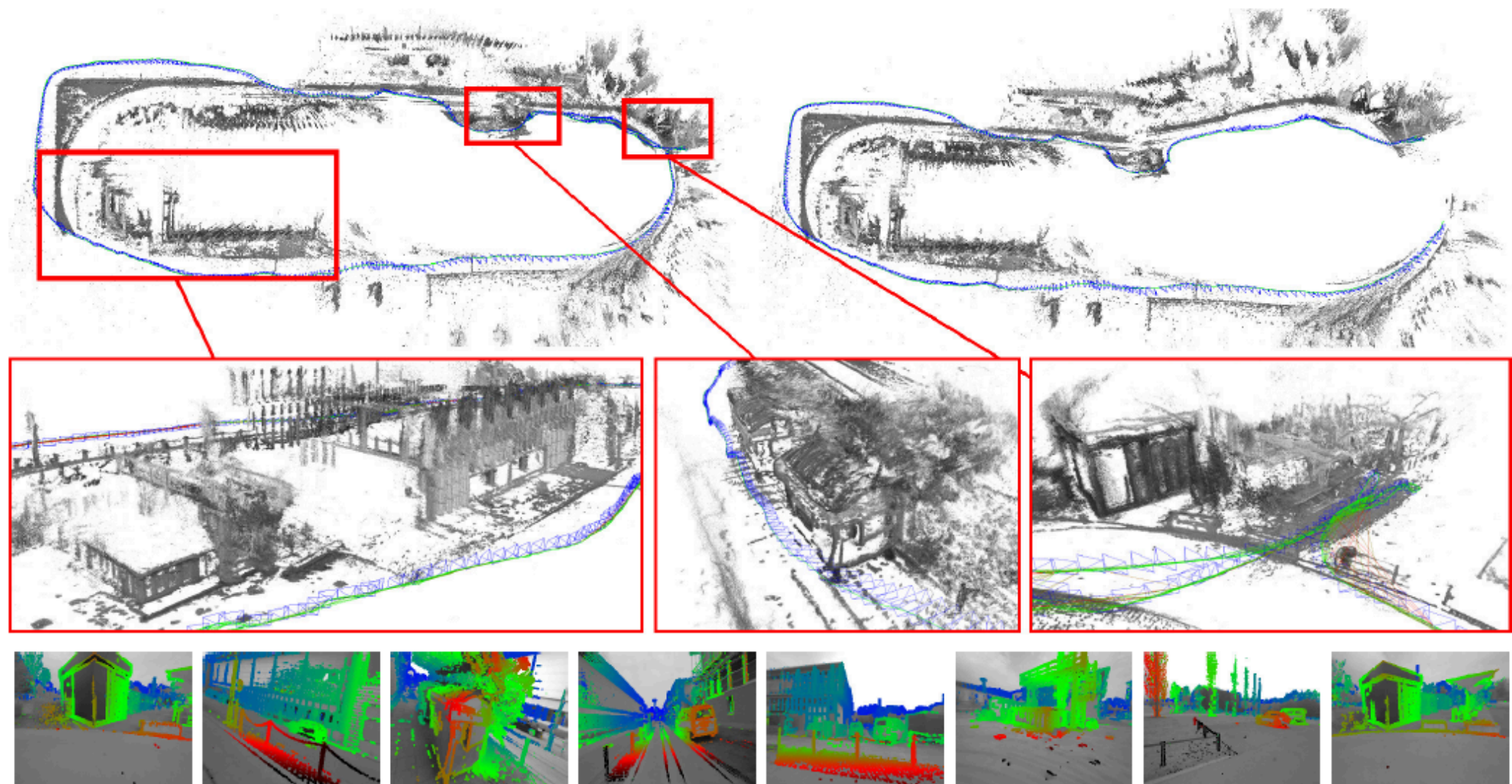


Fig. 7: Loop closure for a long and challenging outdoor trajectory (after the loop closure on the left, before on the right). Also shown are three selected close-ups of the generated pointcloud, and semi-dense depth maps for selected keyframes.

More...

<https://vision.in.tum.de/research/vslam/lslam>

References on web page