



2021  
embedded  
**VISION**  
summit®  
VIRTUAL | MAY 25-28

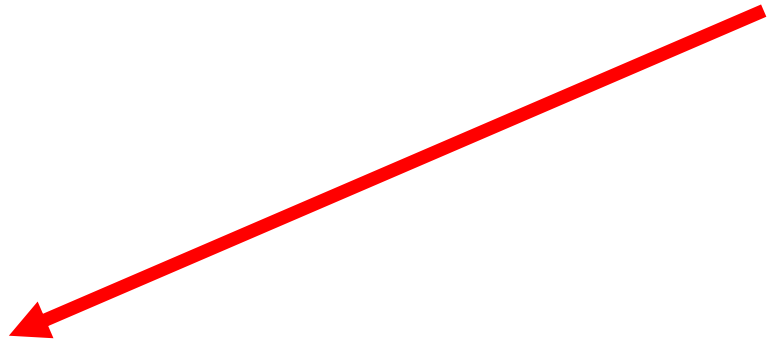
# Introduction To Simultaneous Localization and Mapping (SLAM)

Gareth Cross

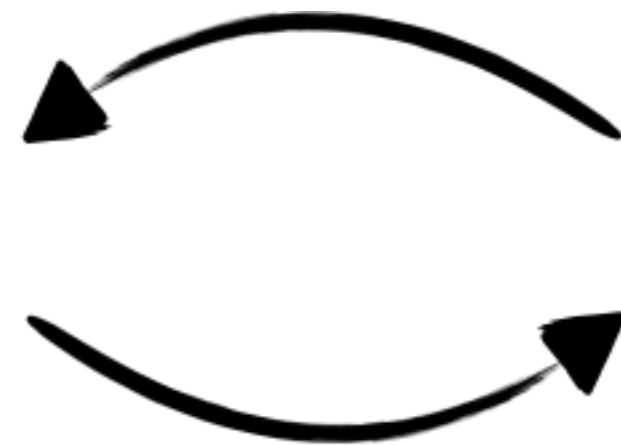
May, 2021



## Simultaneous *Localization and Mapping*



Recover state of a vehicle or sensor platform, usually over multiple time-steps.



Recover location of landmarks in some common reference frame.

*Simultaneous:* We must do these tasks at the same time, as both quantities are initially unknown.

# An age-old practice

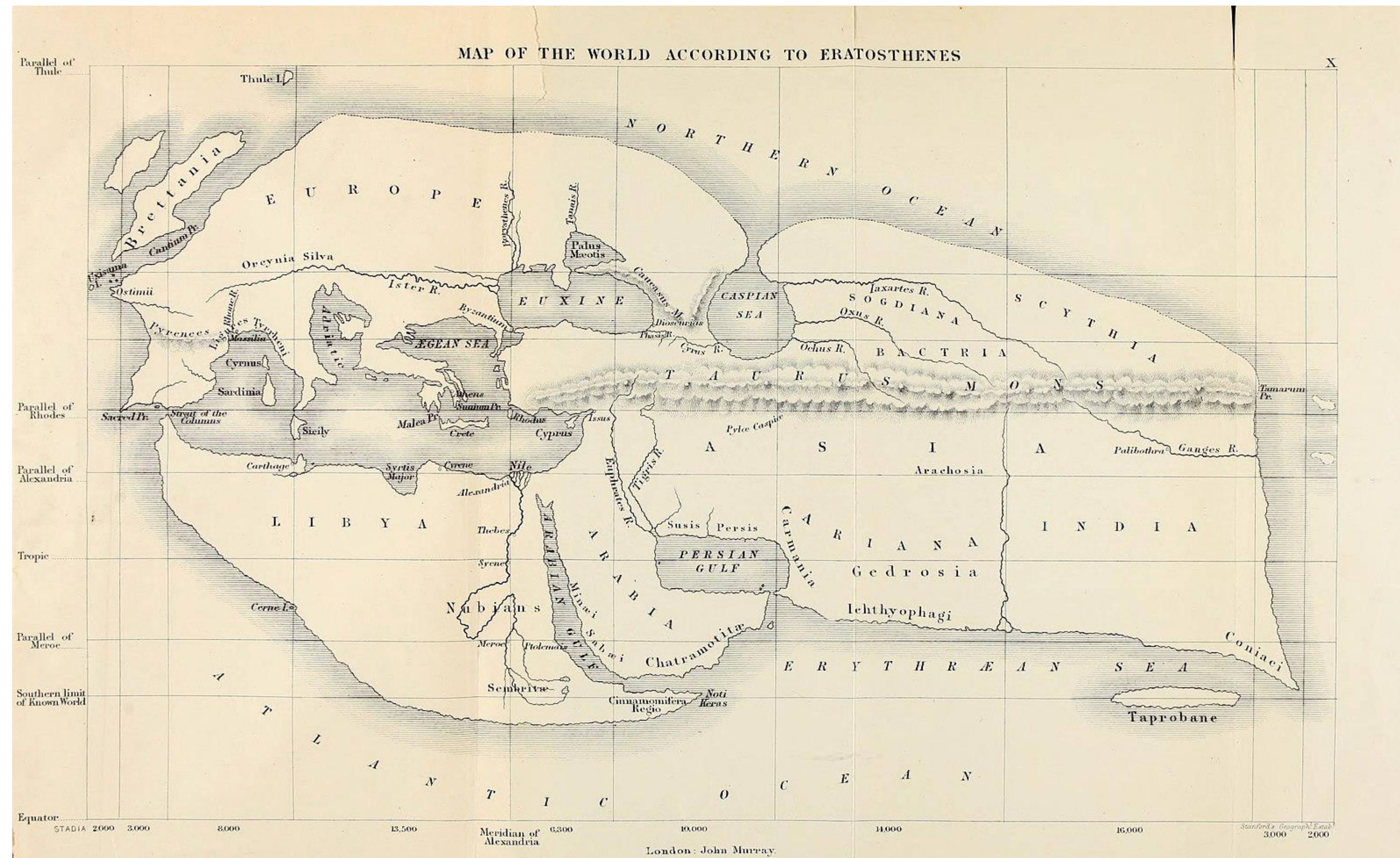
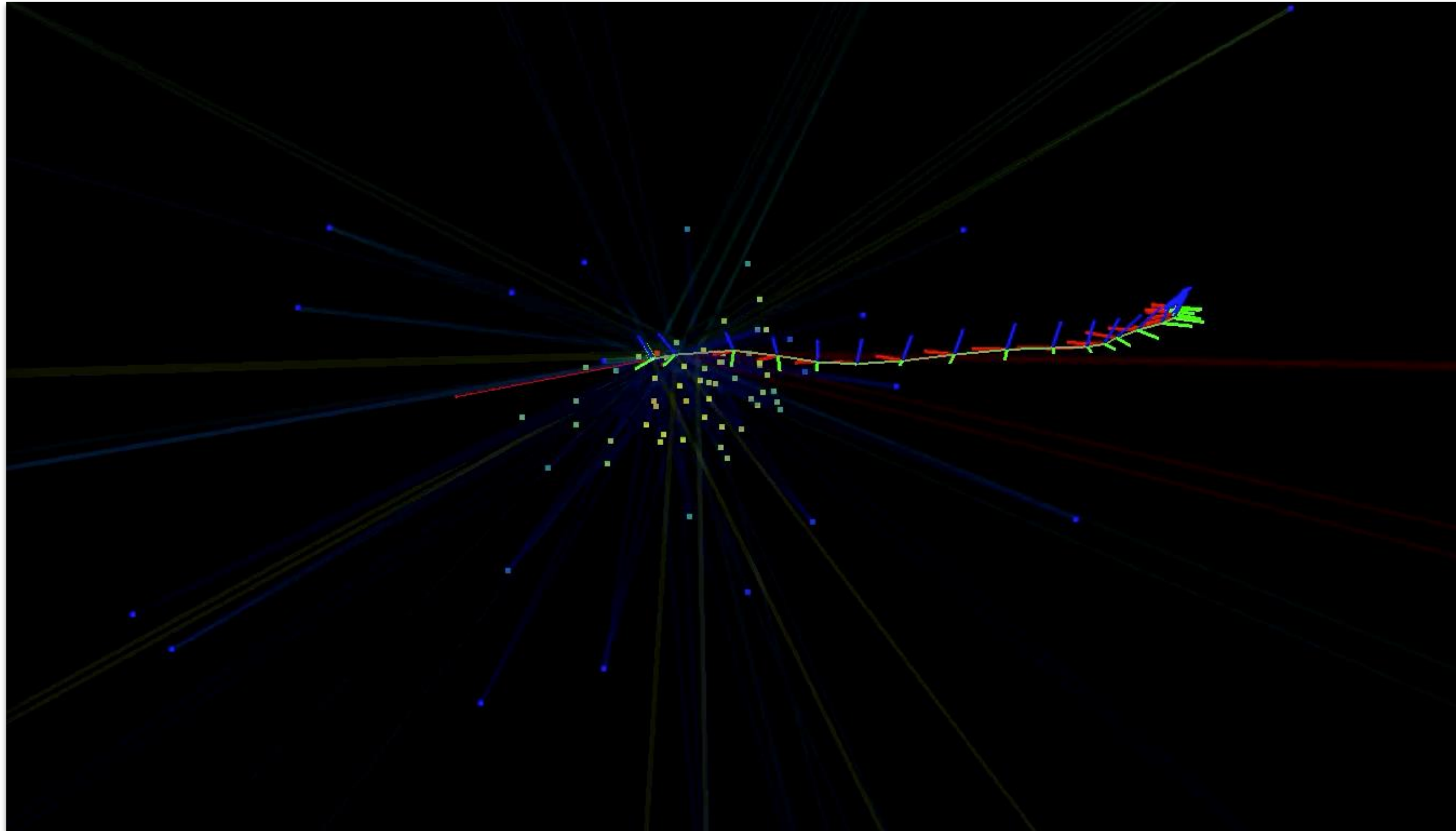


Image Source: [A History of Ancient Geography among the Greeks and Romans from the Earliest Ages till the Fall of the Roman Empire via Wikipedia](#)



Image Source: [COLMAP](#) / Schönberger, Johannes Lutz and Frahm, Jan-Michael, "Structure From Motion Revisited", CVPR 2016



Visual Inertial Odometry (VIO) on the Skydio drone, an embedded system.



# SLAM vs. Localization

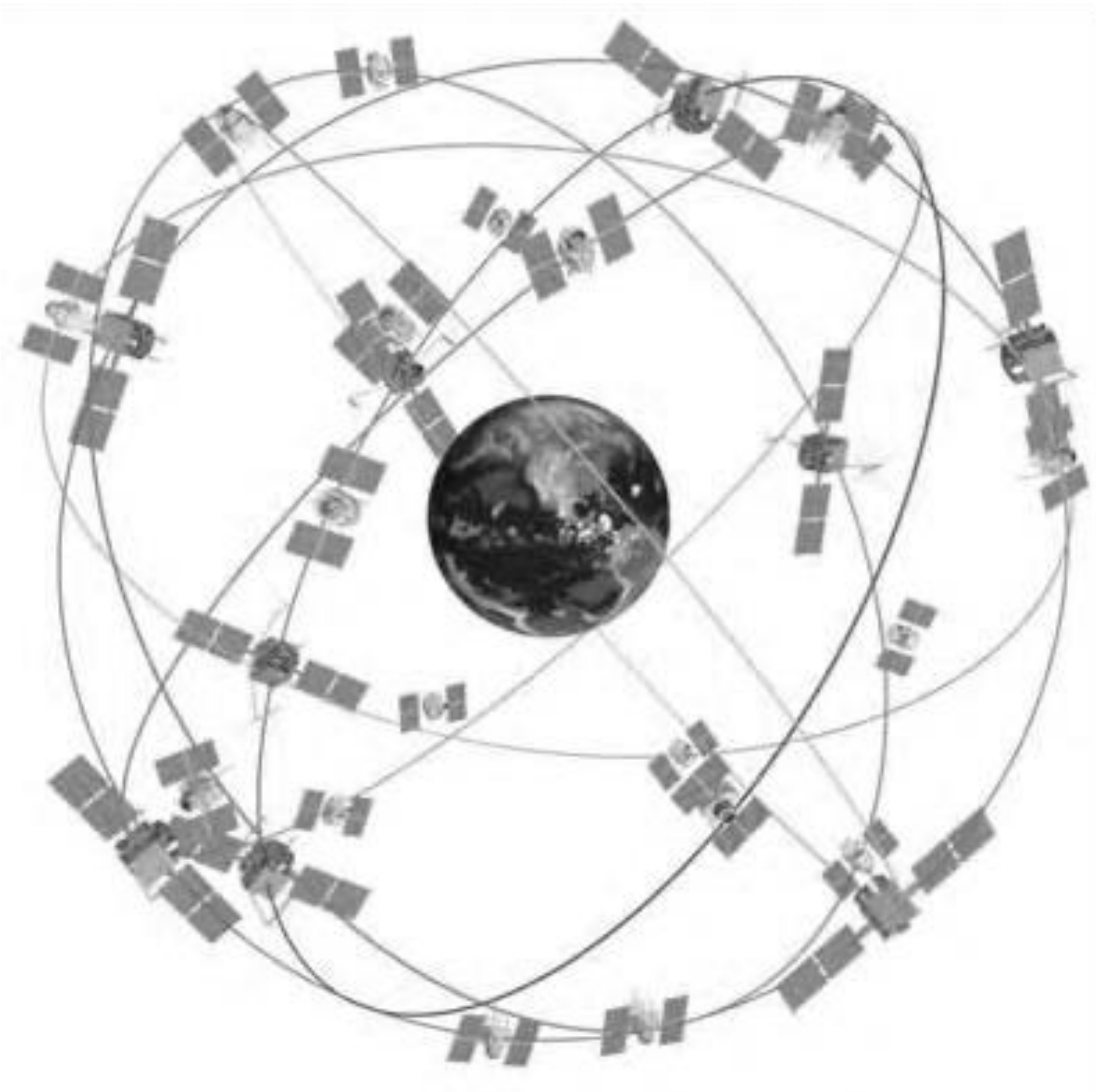
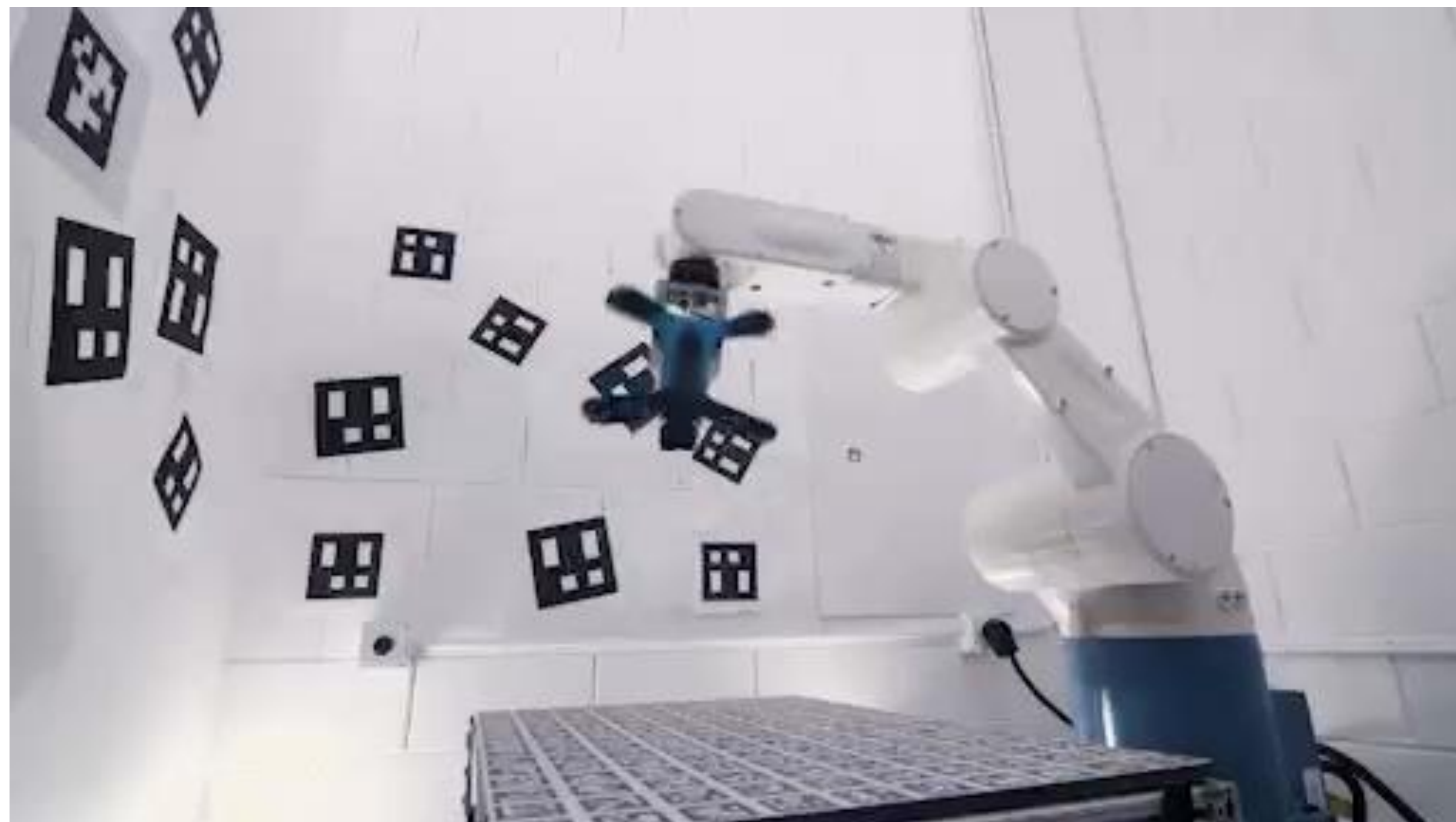


Image Source: E. Kaplan, C. Hergarty,  
*Understanding GPS Principles and Applications*,  
2005

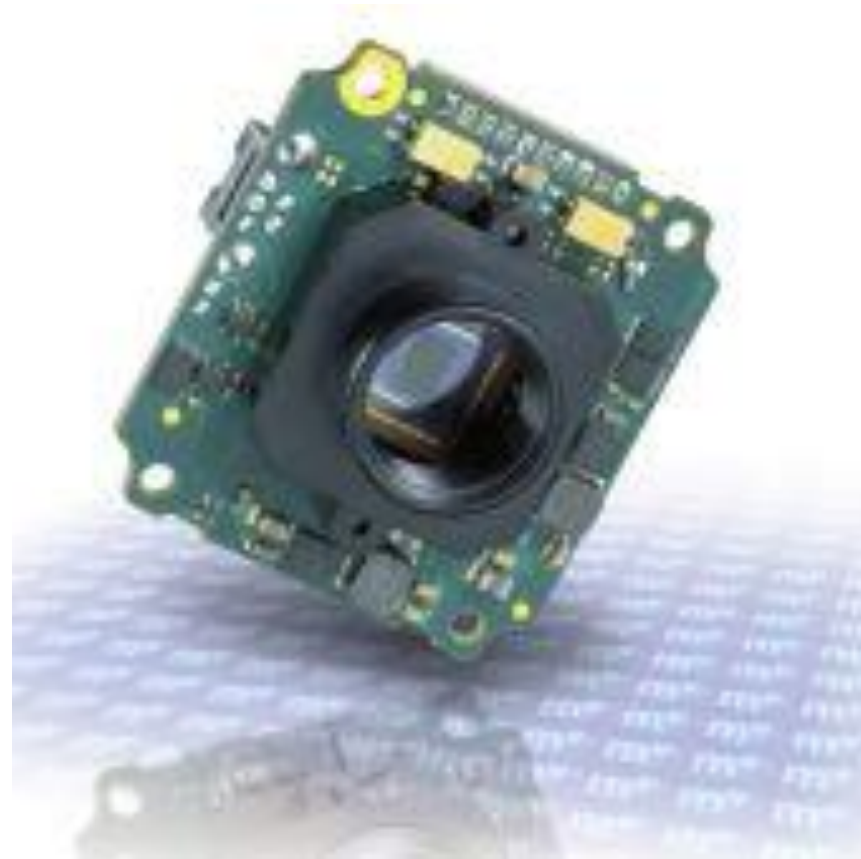


Video Source: Skydio

For every SLAM problem, we have two key ingredients:

1) One or more *sensors*:

Cameras



Source: [MatrixVision](#)

Inertial Measurement Unit



Source: [Lord MicroStrain](#)

LiDAR/Range-finders



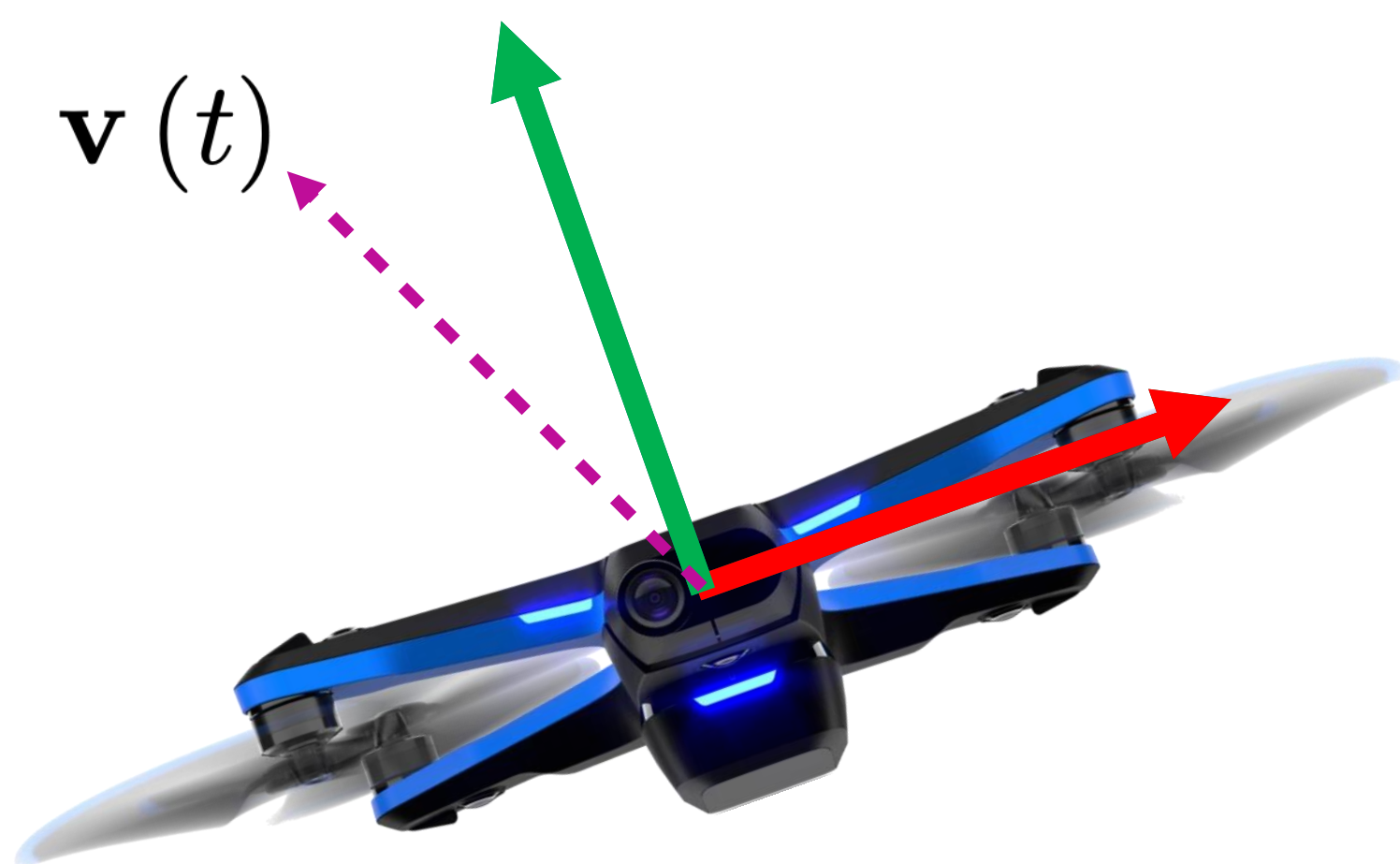
Source: [Velodyne](#)

RGB-D/Structured Light



Source: [Occipital](#)

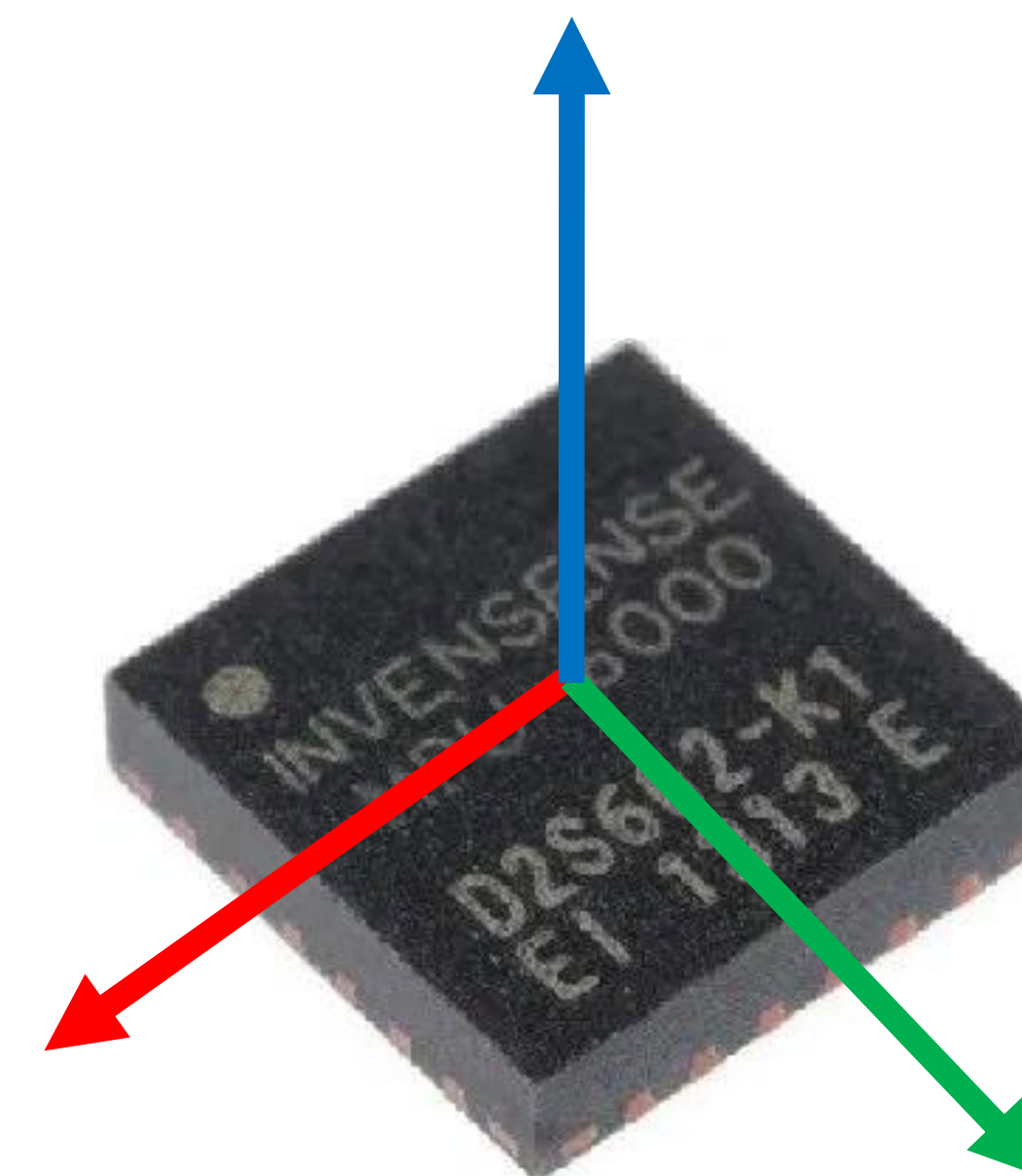
2) A set of *states* we wish to recover.



*Ego-motion*: Rotation, position, velocity



World Structure (Map)



Calibration Parameters

Image Source: [DroneTest](#)



- Choice of sensor will drive many downstream design considerations.
- Consider the *sensor measurement model*:

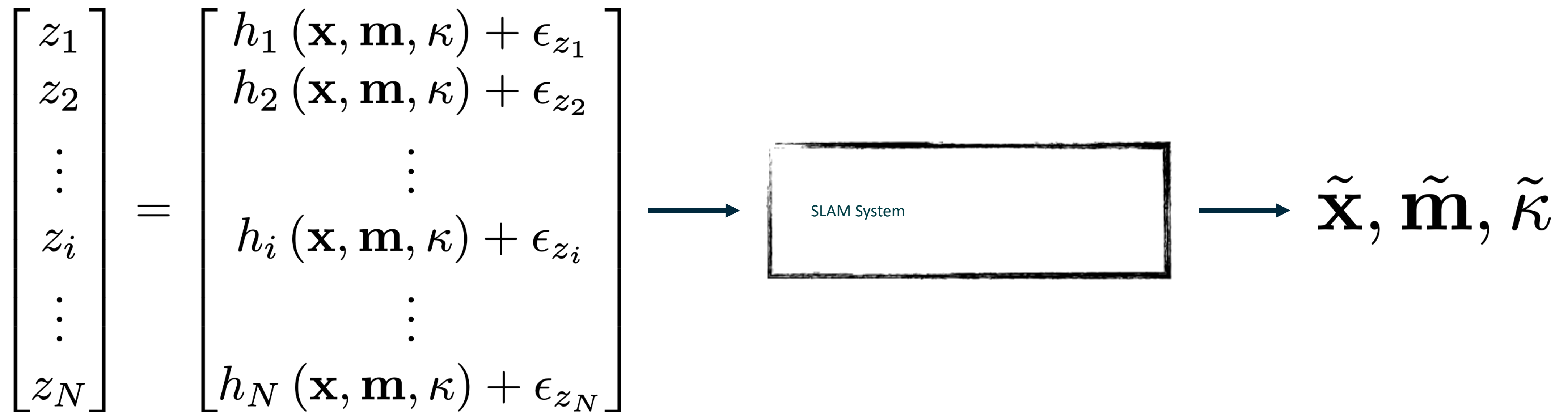
$$\mathbf{z} = h(\mathbf{x}, \mathbf{m}, \kappa) + \epsilon_z$$

Ego-motion                      Calibration parameters

Sensor output                      Map                      Noise

The diagram illustrates the sensor measurement model equation  $\mathbf{z} = h(\mathbf{x}, \mathbf{m}, \kappa) + \epsilon_z$ . It features three labels with arrows pointing to the corresponding variables in the equation: 'Ego-motion' points to  $\mathbf{x}$ , 'Calibration parameters' points to  $\kappa$ , and 'Map' points to  $\mathbf{m}$ . Below the equation, three labels are positioned: 'Sensor output' with an arrow pointing up to  $\mathbf{z}$ , 'Map' with an arrow pointing up to  $\mathbf{m}$ , and 'Noise' with an arrow pointing up to  $\epsilon_z$ .

Take many measurements (possibly from many sensors), and recover the ego-motion, map, and calibration parameters.



# Example — Calibration



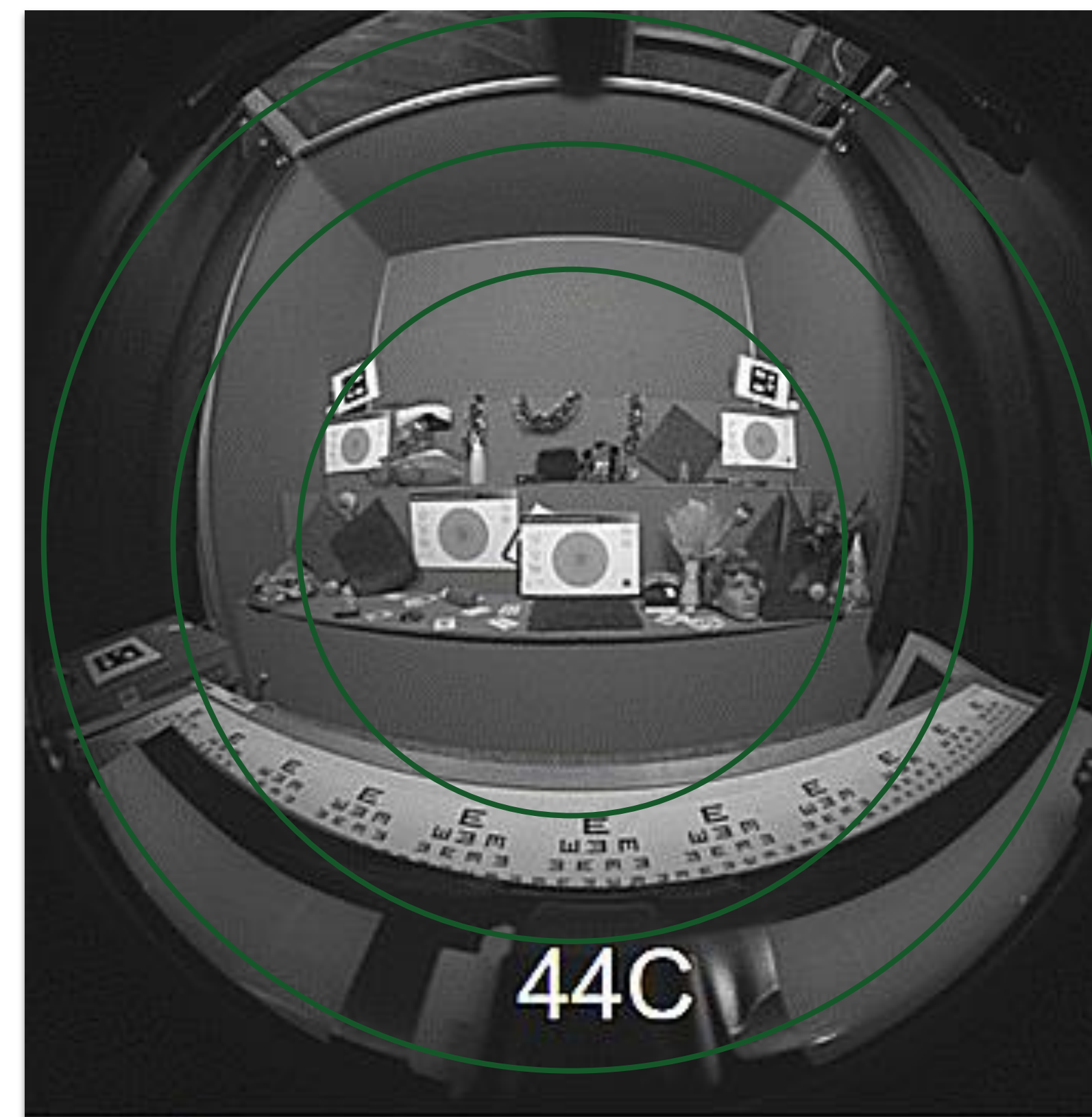
Tire inflation will affect the scale of wheel odometry, as could slippage between the tire and the road surface.

Image source: [MotorTrend.com](https://www.motortrend.com)



Un-modelled extrinsic rotation between IMU and camera may cause increased drift in a visual SLAM pipeline.

Image source: [MWee RF Microwave](https://www.mwrf.com)



Intrinsic temperature distortion may also introduce unexpected errors into vision estimates.

Image Source: Skydio

# Example — Rolling Shutter

When selecting a camera sensor for your platform, you have the choice of *global* or *rolling* shutter.

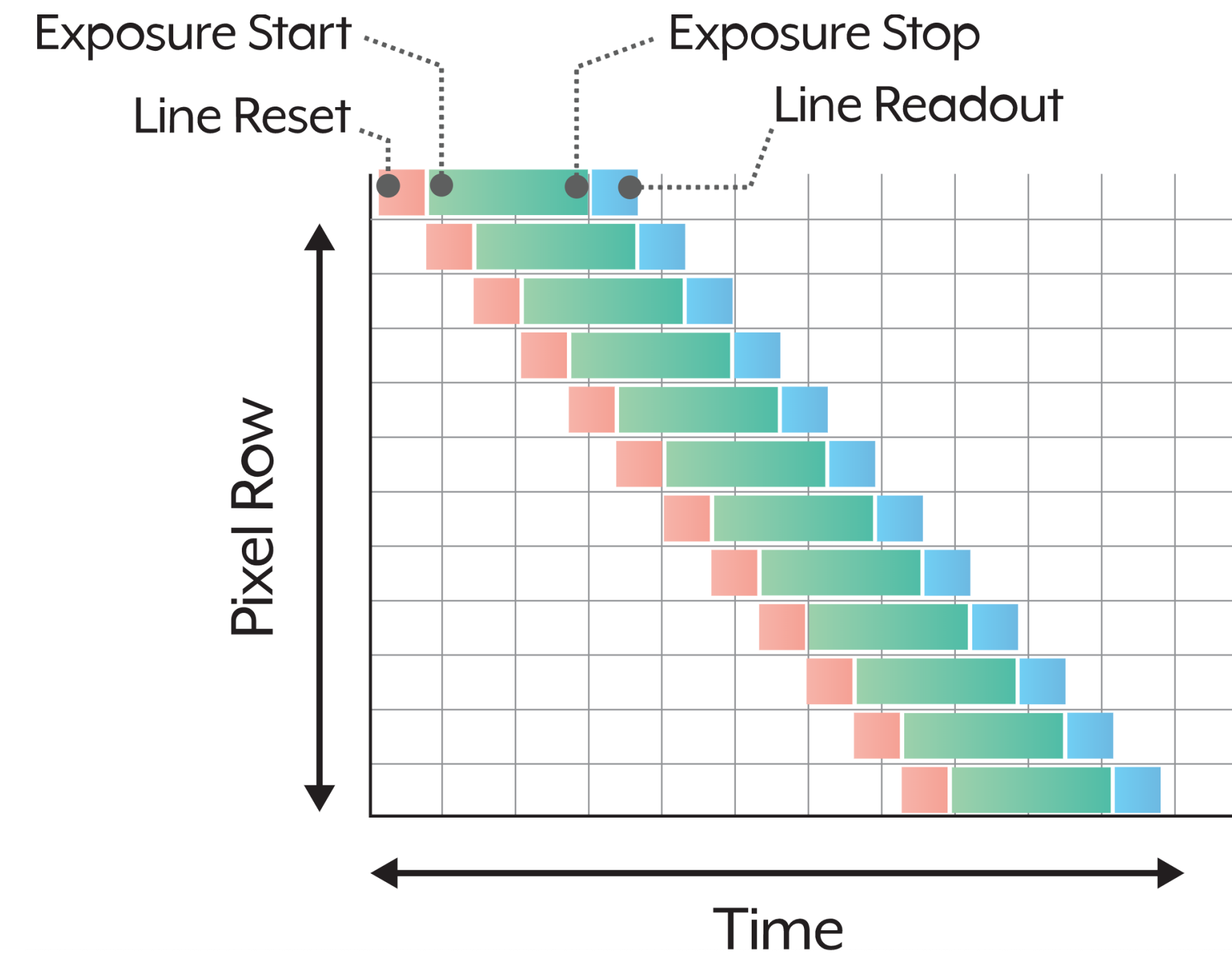
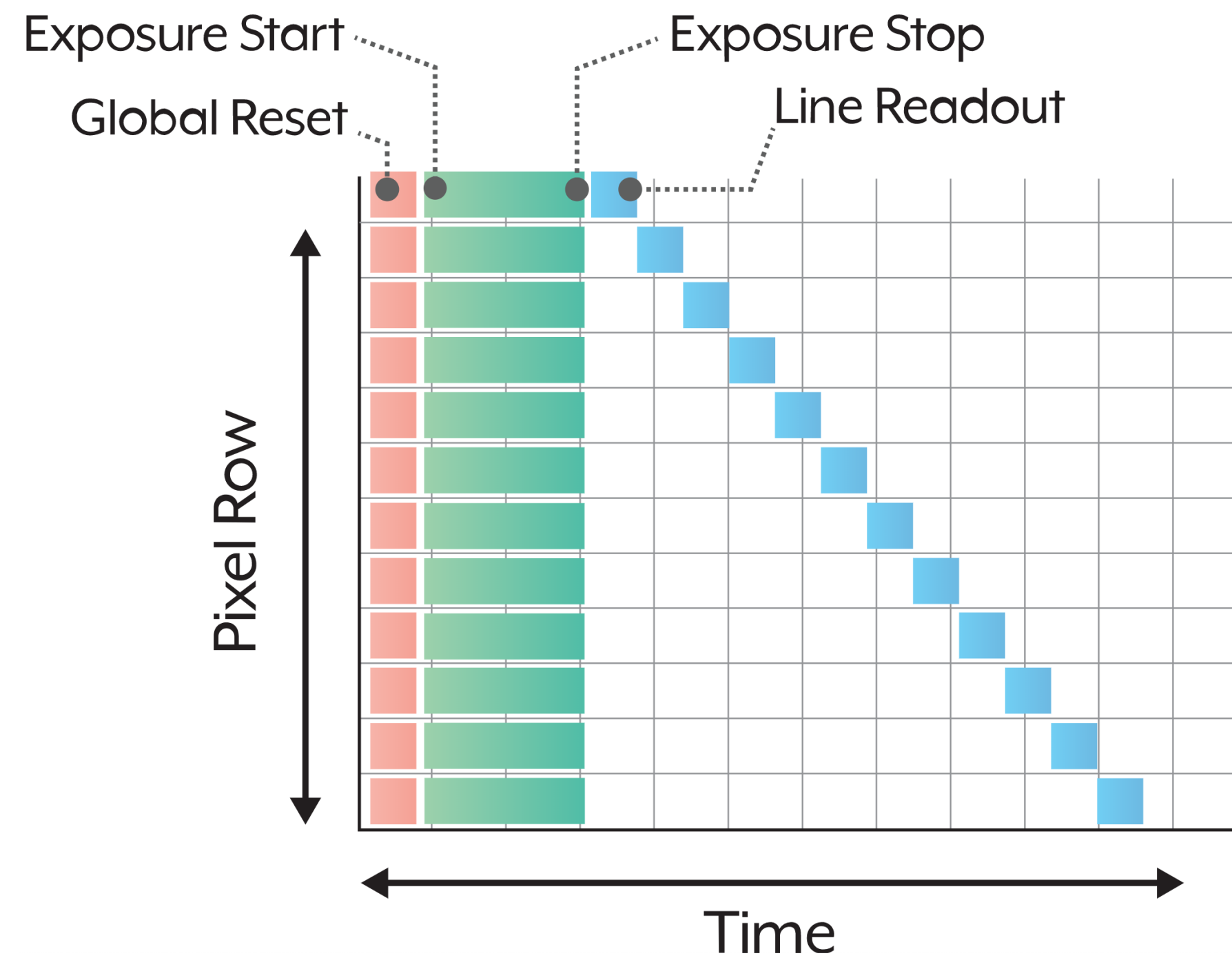


Image credit: [LucidVision](#)

# Example — Rolling Shutter



Rolling shutter deforms rigid objects like the horizon line and the vehicle itself.

Global-shutter model:

$$\mathbf{z} = \pi \left( T_w^c(t_0) \oplus p_w \right)$$

Diagram illustrating the Global-shutter model equation  $\mathbf{z} = \pi \left( T_w^c(t_0) \oplus p_w \right)$ . Annotations include:

- Camera projection (pointing to  $\pi$ )
- Measured location in image (pointing to  $\mathbf{z}$ )
- Pose of the camera at time  $t_0$  (pointing to  $T_w^c(t_0)$ )
- Point in the world (pointing to  $p_w$ )
- Transformation of world points into sensor frame (pointing to  $\oplus$ )

Rolling-shutter model:

$$\mathbf{z} = \pi \left( T_w^c(t_0) \int_{t_0}^{\tau} \dot{T}_w^c(t) dt \oplus p_w \right)$$

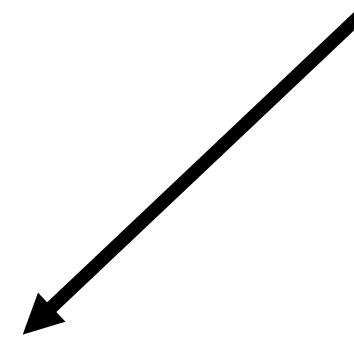
Diagram illustrating the Rolling-shutter model equation  $\mathbf{z} = \pi \left( T_w^c(t_0) \int_{t_0}^{\tau} \dot{T}_w^c(t) dt \oplus p_w \right)$ . An annotation includes:

- Inclusion of higher-order derivatives in the measurement model increases computational cost (pointing to  $\dot{T}_w^c(t)$ )

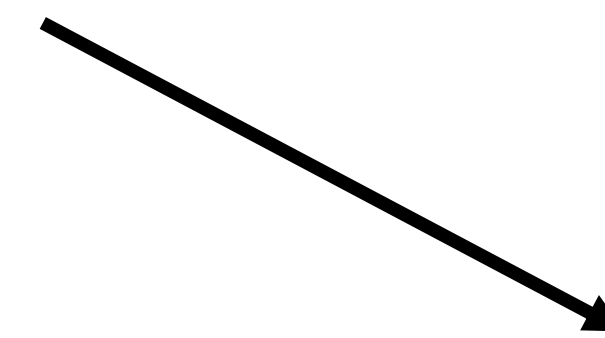
Inclusion of higher-order derivatives in the measurement model increases computational cost.

$$\mathbf{z} = h(\mathbf{x}, \mathbf{m}, \kappa) + \epsilon_z$$

- All sensors exhibit some minimum amount of *noise*.
- We distinguish between *noise* and *model error*.



A random error that can only be modeled via statistical means.  
Example: thermal electrical noise.



Errors resulting from a limitation in our sensor model.  
Example: failure to include a calibration parameter.

- Owing to noise in the sensor inputs, SLAM is an inherently uncertain process.
- We can never recover the “true” states, only uncertain estimates of them.
  - More measurements *usually* means reduced uncertainty...
  - ... But it also means increased computational cost.

- Choice of sensor may also influence map parameterization.

Collection of  
photographs?

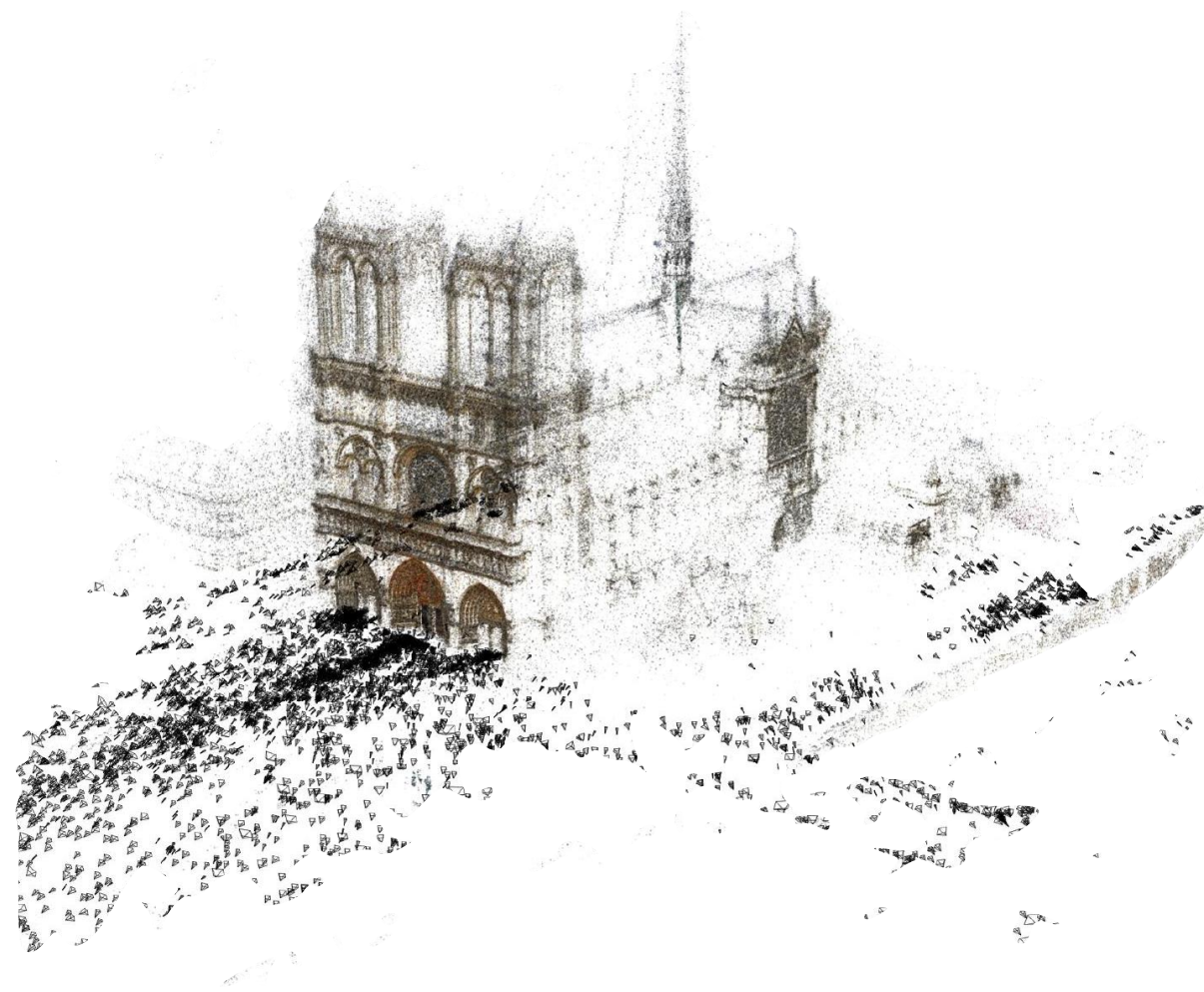


Image Source: [Noah Snavely](#)

2D LiDAR scans?

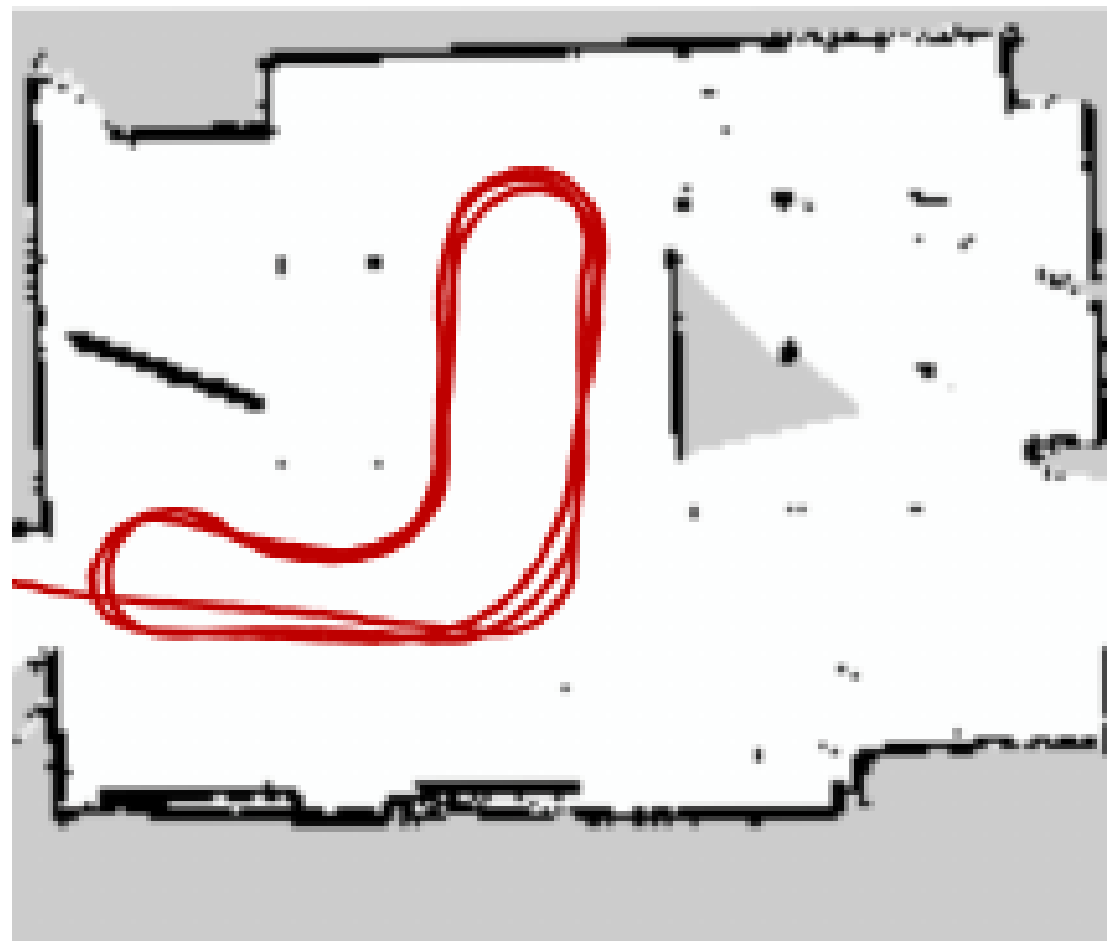


Image Source: B. Bellekens et al., [A Benchmark Survey of Rigid 3D Point Cloud Registration Algorithms](#), 2015

3D range images?

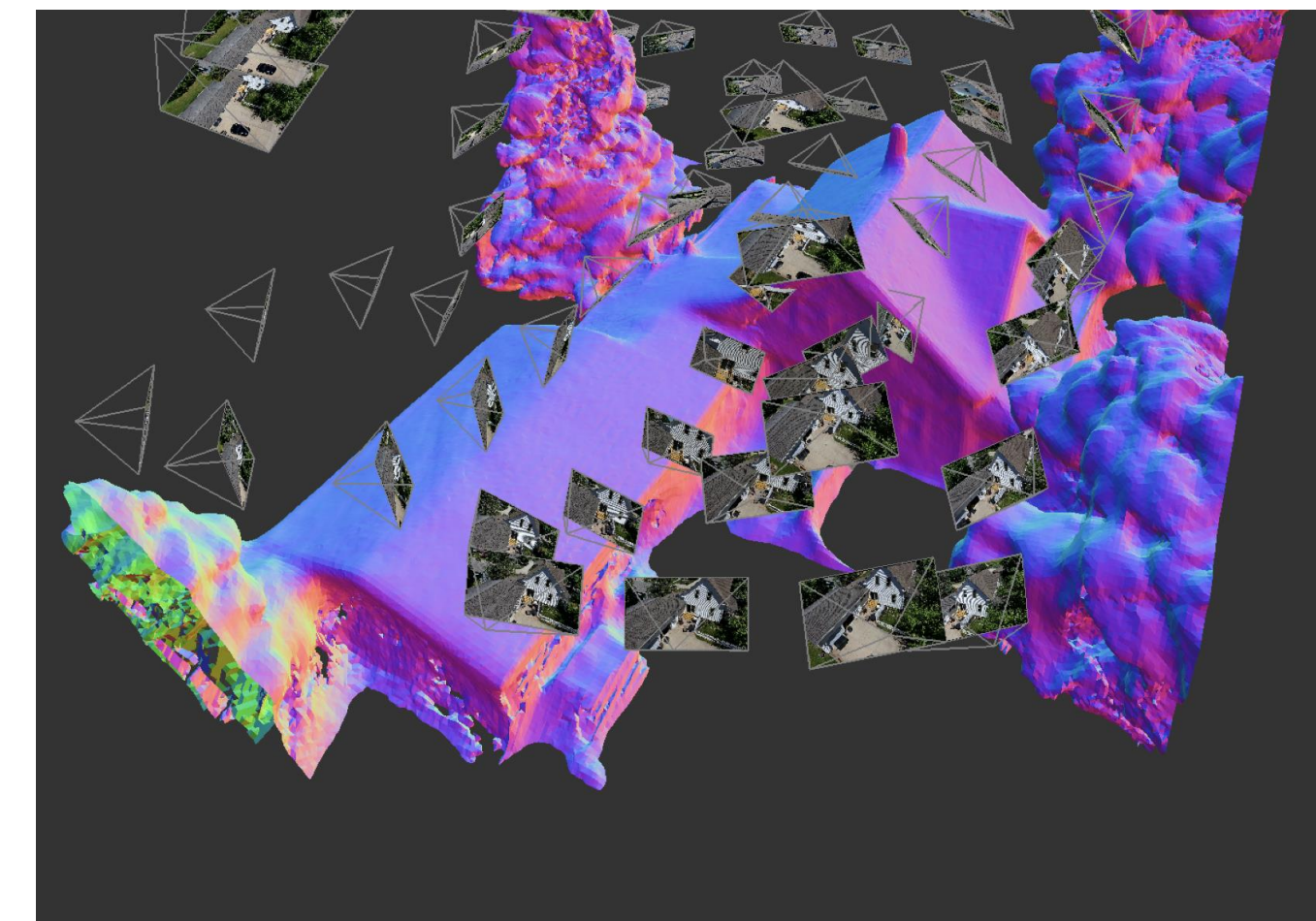
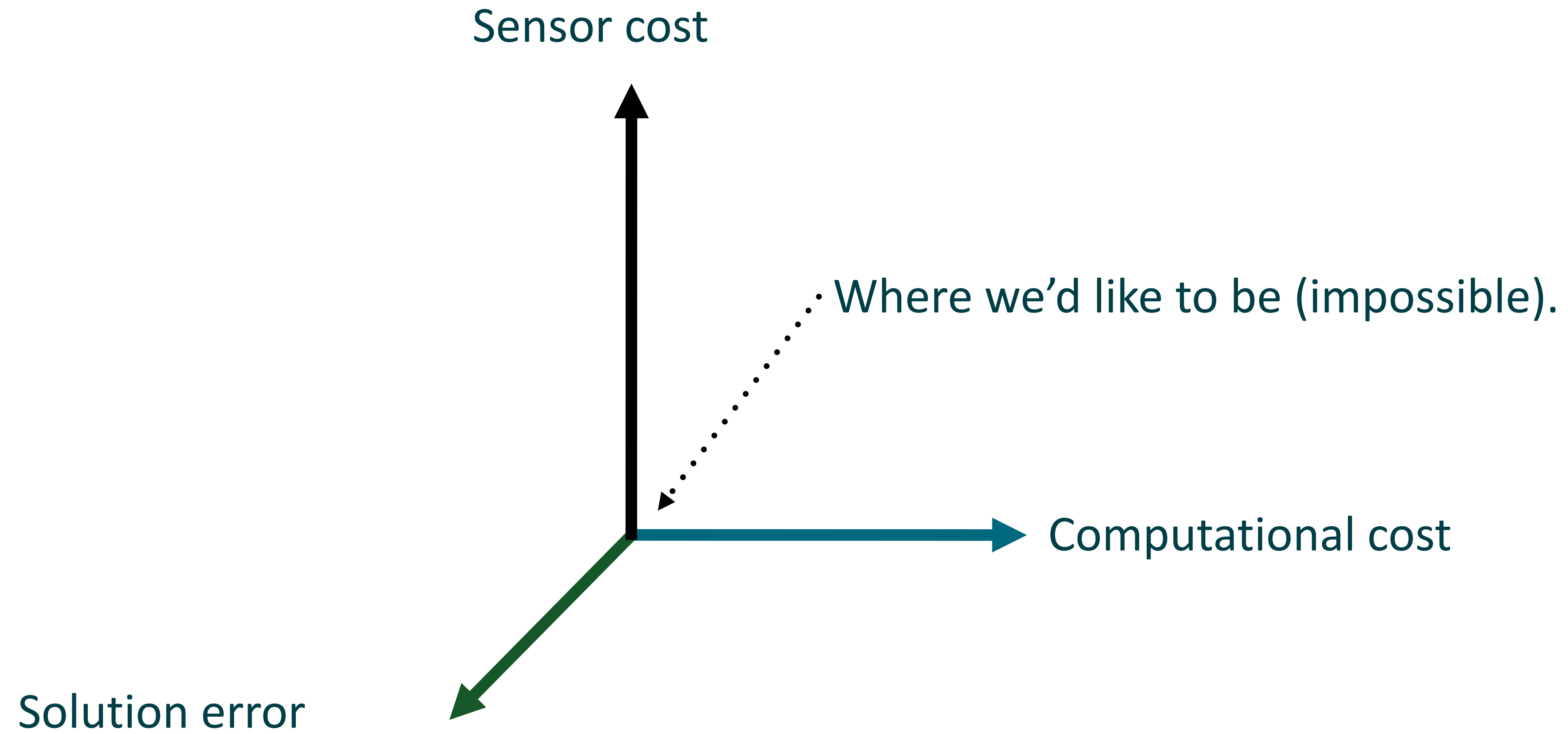
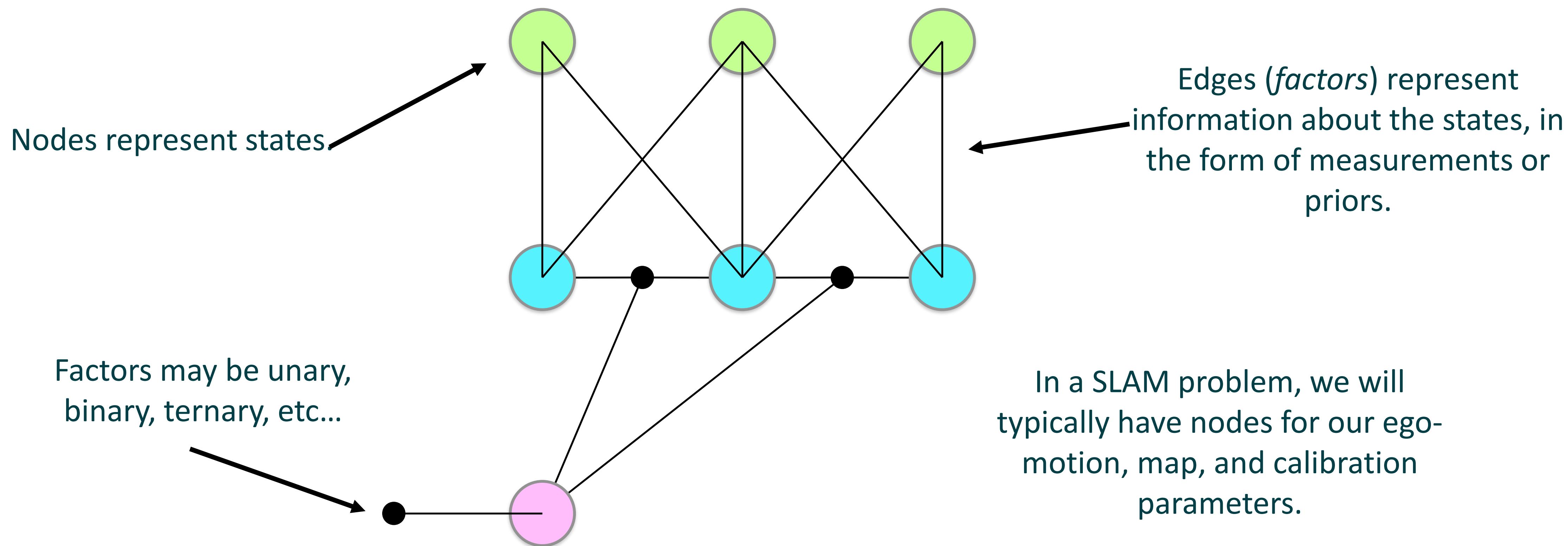


Image Source: Skydio

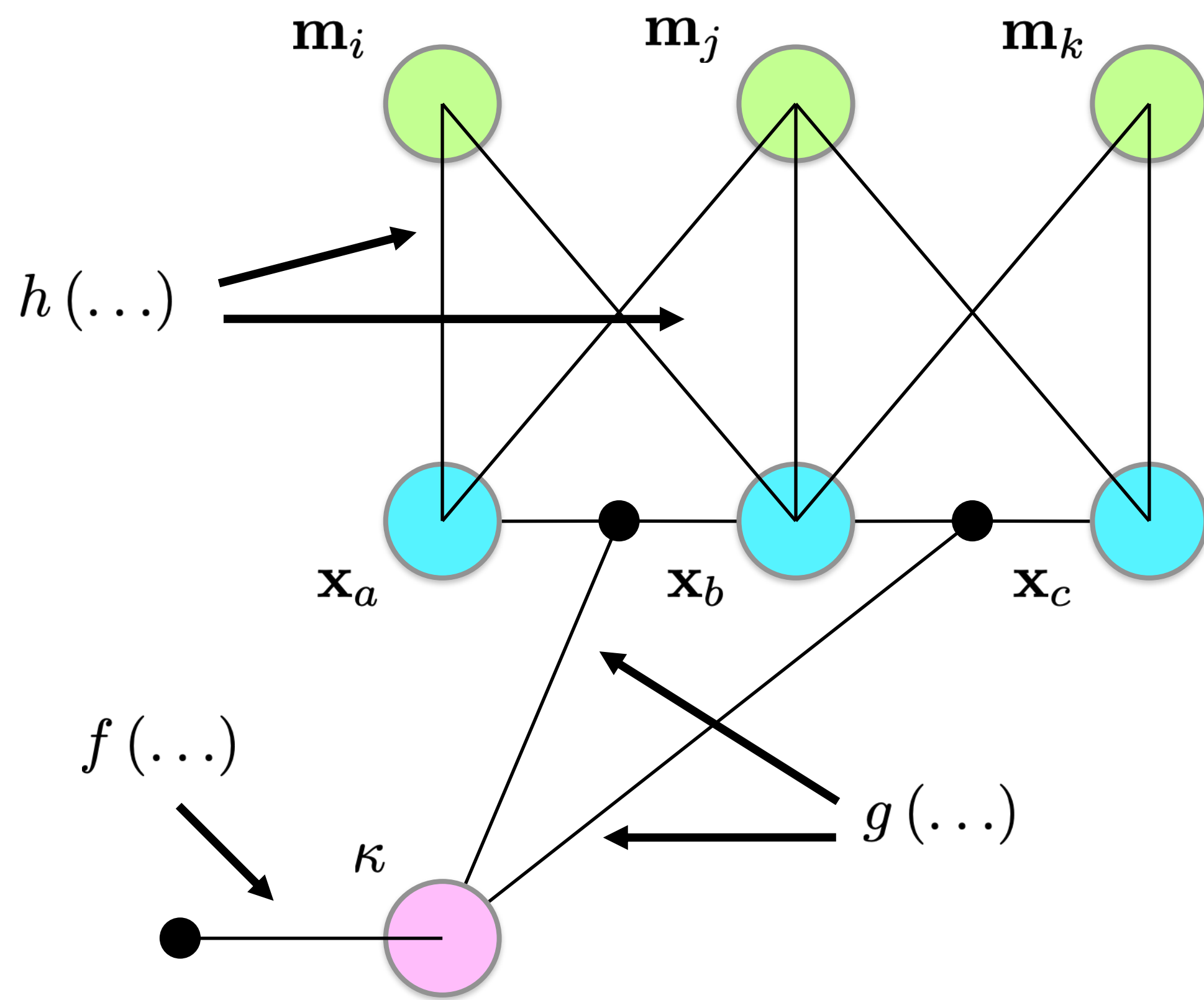




- Factor Graphs are a convenient method of graphically representing a SLAM problem.



There is a mapping from the factor graph to our sensor measurements:



$$\begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \\ z_5 \\ z_6 \\ z_7 \\ z_8 \\ z_9 \\ z_{10} \end{bmatrix} = \begin{bmatrix} h(\mathbf{x}_a, \mathbf{m}_i) + \epsilon_{z_1} \\ h(\mathbf{x}_a, \mathbf{m}_j) + \epsilon_{z_2} \\ h(\mathbf{x}_b, \mathbf{m}_i) + \epsilon_{z_3} \\ h(\mathbf{x}_b, \mathbf{m}_j) + \epsilon_{z_4} \\ h(\mathbf{x}_b, \mathbf{m}_k) + \epsilon_{z_5} \\ h(\mathbf{x}_c, \mathbf{m}_j) + \epsilon_{z_6} \\ h(\mathbf{x}_c, \mathbf{m}_k) + \epsilon_{z_7} \\ g(\mathbf{x}_a, \mathbf{x}_b, \kappa) + \epsilon_{z_8} \\ g(\mathbf{x}_b, \mathbf{x}_c, \kappa) + \epsilon_{z_9} \\ f(\kappa) + \epsilon_{z_{10}} \end{bmatrix}$$

# Real Example: Bundle Adjustment (BA)

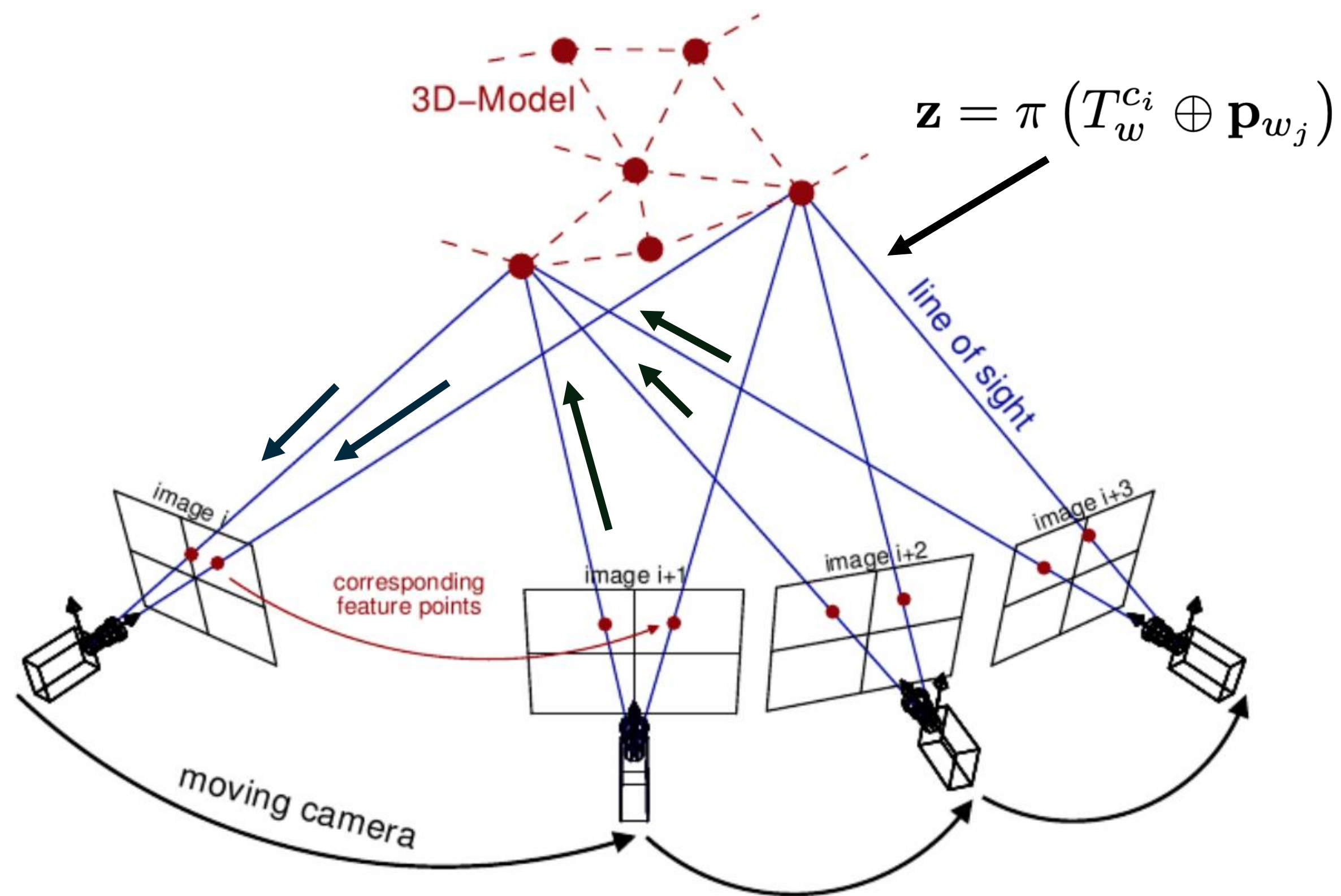


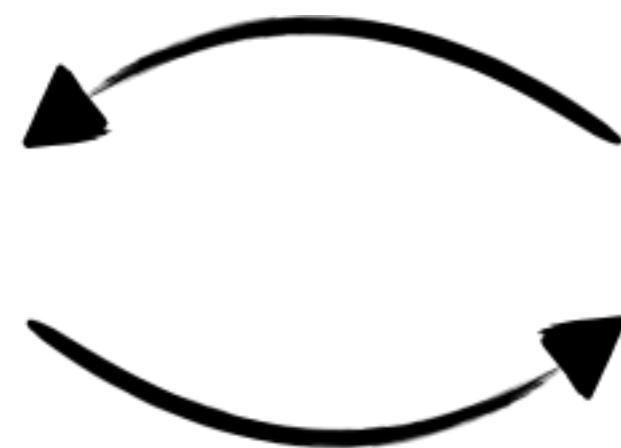
Image source: [Theia SFM](#)

- A form of *Structure from Motion (SFM)*.
- Leverage *projective geometry* to recover 3D landmarks and poses from 2D feature associations.
- Highly scalable and can be quite accurate.
- Using *marginalization* the compute cost can be bounded.

For more details on SFM, see [Richard Szeliski's book](#) as a jumping off point.

## Simultaneous *Localization and Mapping*

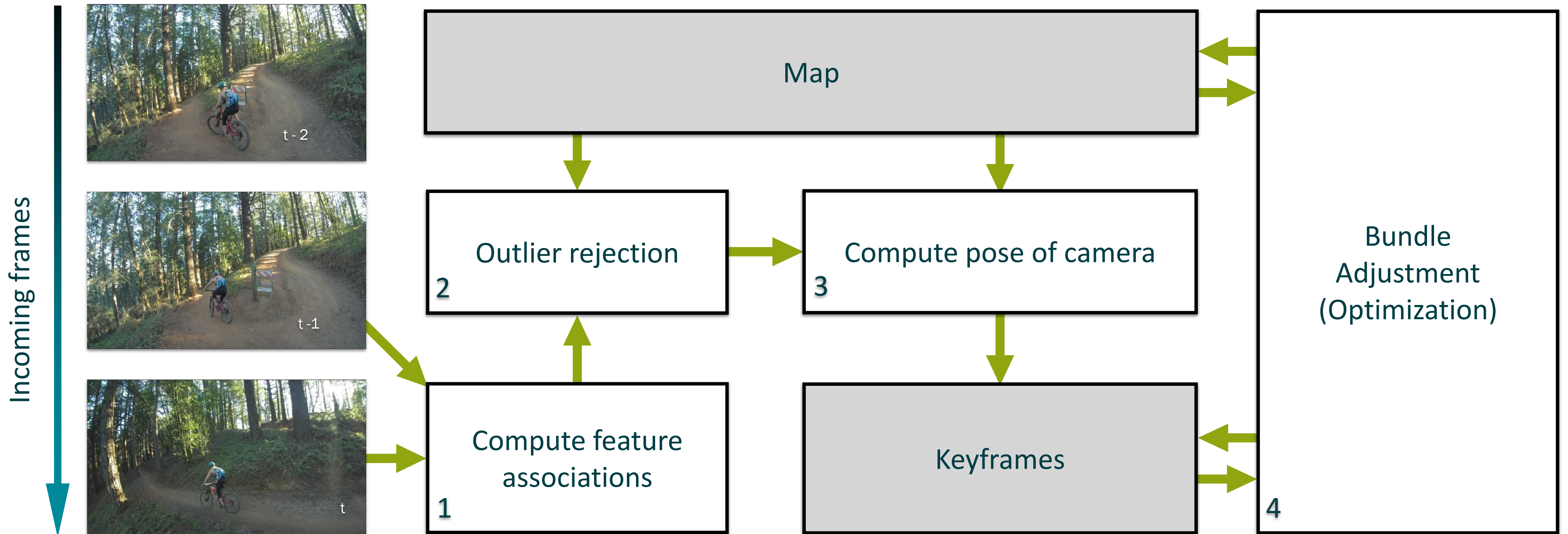
Recover state of a vehicle or  
sensor platform, usually over  
multiple time-steps.  
Recover camera pose with  
respect to map points.



Recover location of landmarks in  
some common reference frame.  
Triangulate map points  
using camera poses.

Bundle Adjustment is a form of optimization that does these steps *jointly*.

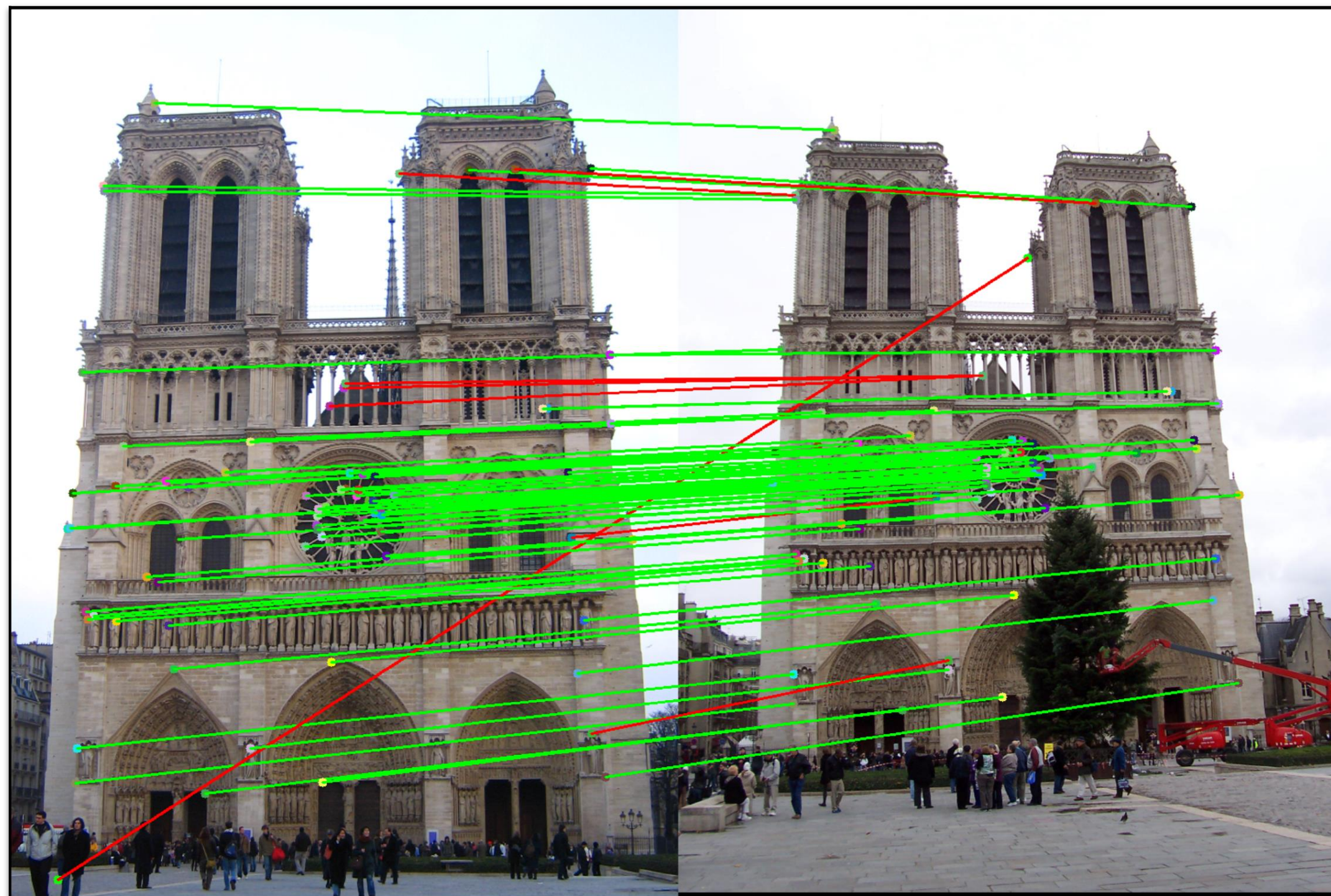
# Typical SLAM Pipeline w/ BA



Keyframes store our estimates of the ego-motion.

# How do we get feature associations?

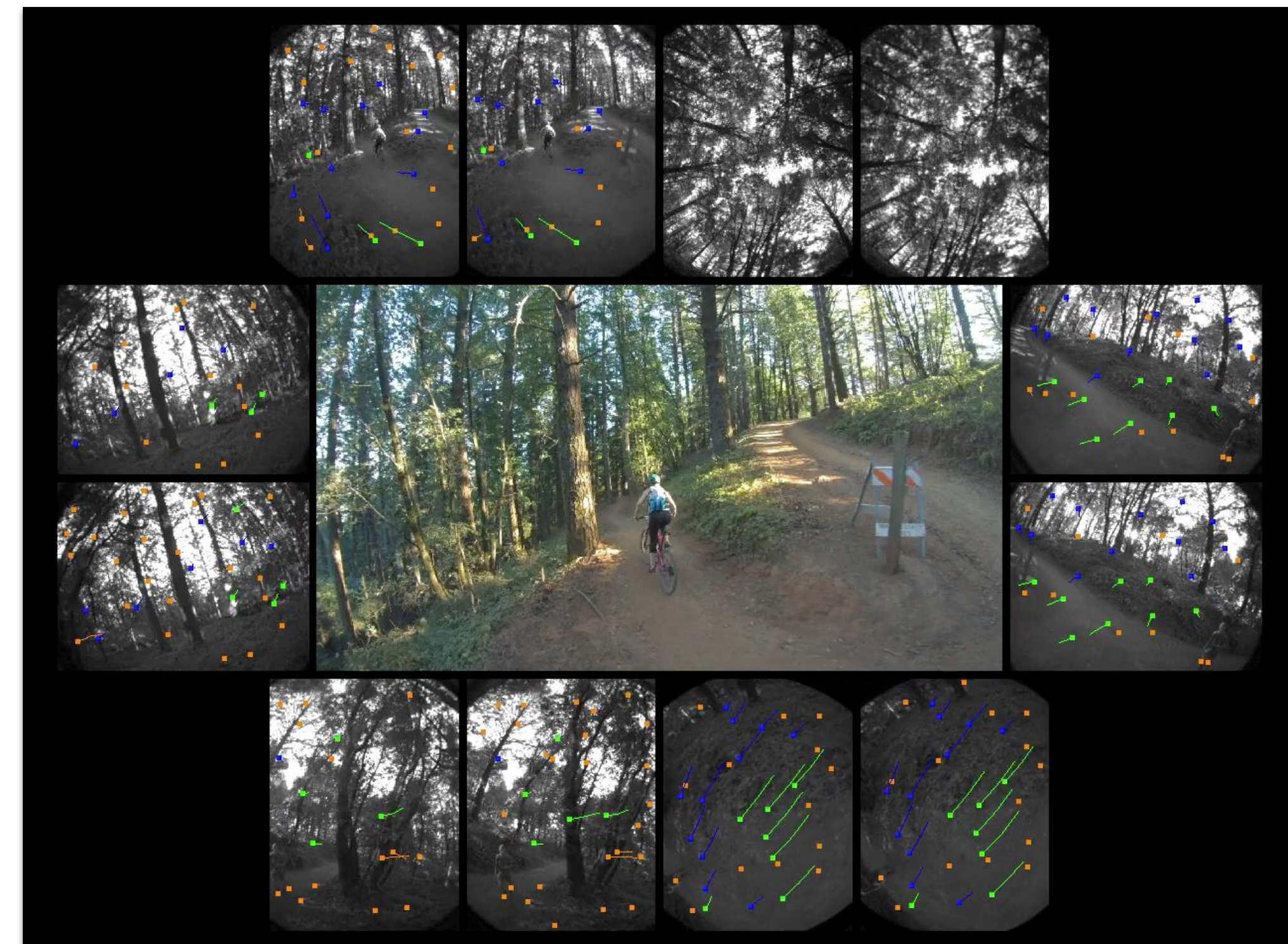
## Descriptor Matching



Examples: [SIFT](#), [KAZE](#), [ORB](#), [SuperPoint](#)

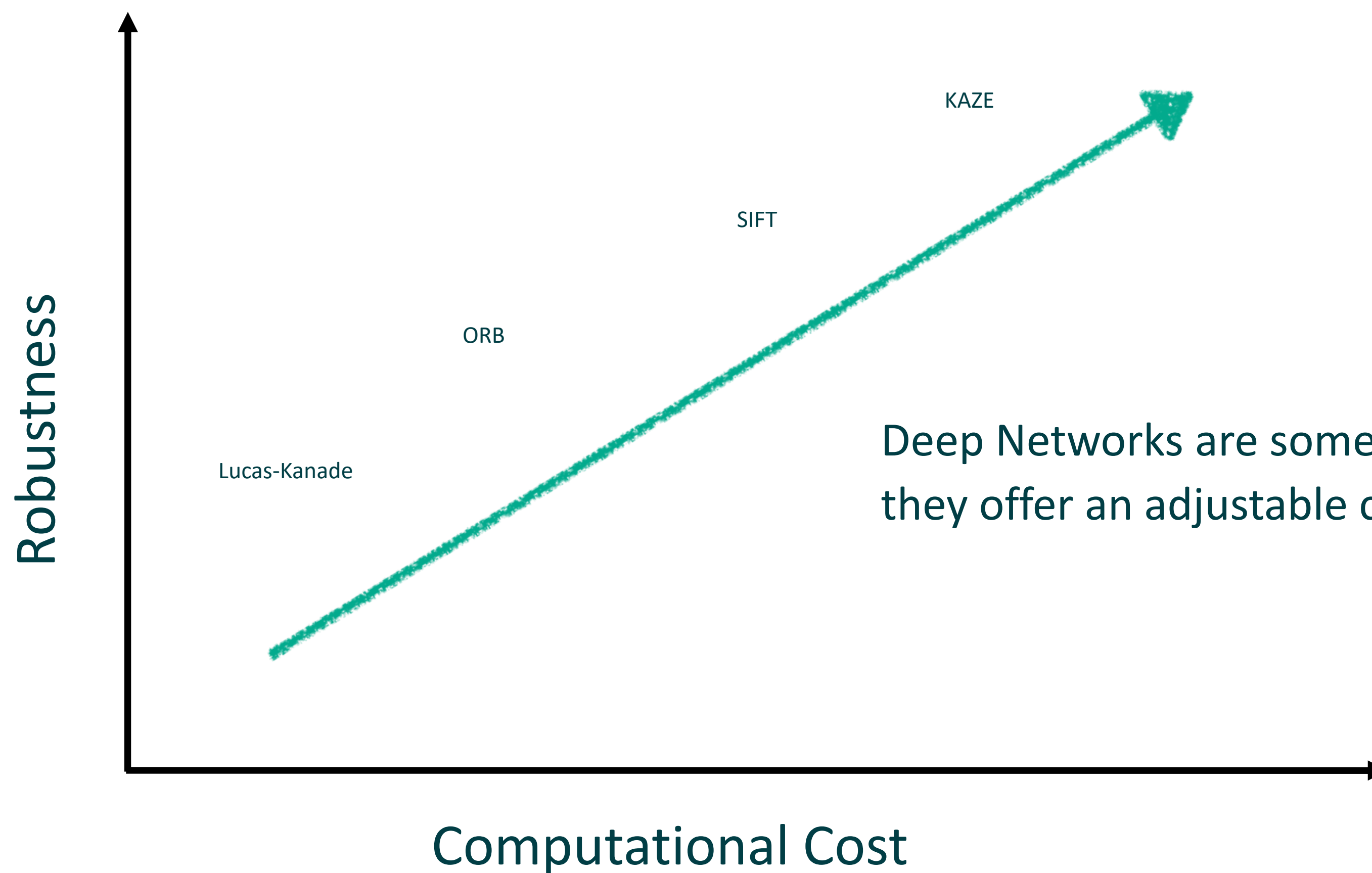
Image Source: [Georgia Tech](#)

## Feature Tracking / Flow



Examples: [Optical Flow](#), [Lucas Kanade Tracking](#), [FlowNet](#)

A “rule of thumb” principle to consider in selecting features (axes not to scale):



Deep Networks are somewhat difficult to place since they offer an adjustable cost-robustness trade-off.



# Outliers in Feature Association



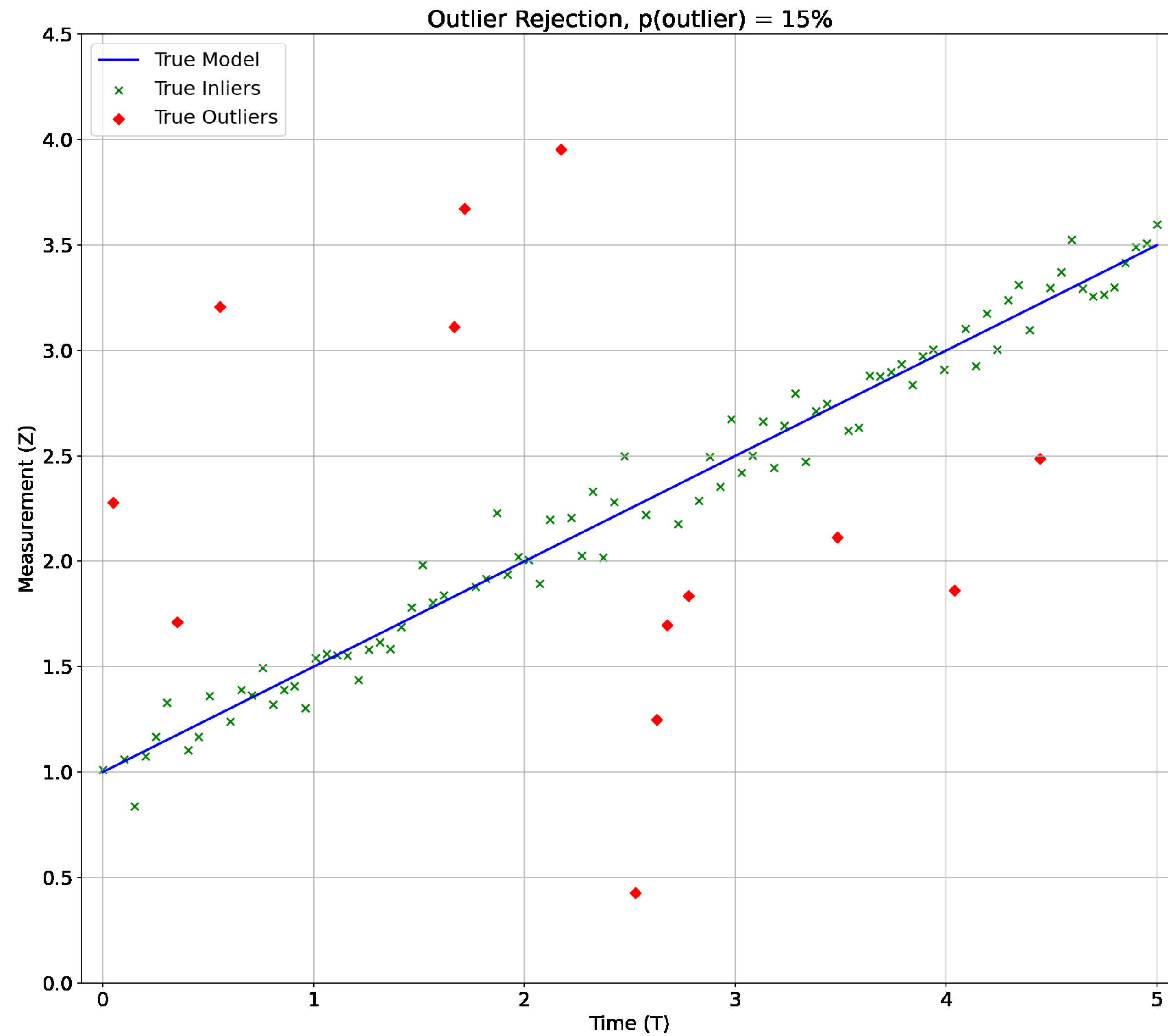
- *Outliers*: Data that does not agree with our sensor model.
- How do we deal with them?
- Let's review a (very simple) toy problem:

State: *alpha* and *beta*

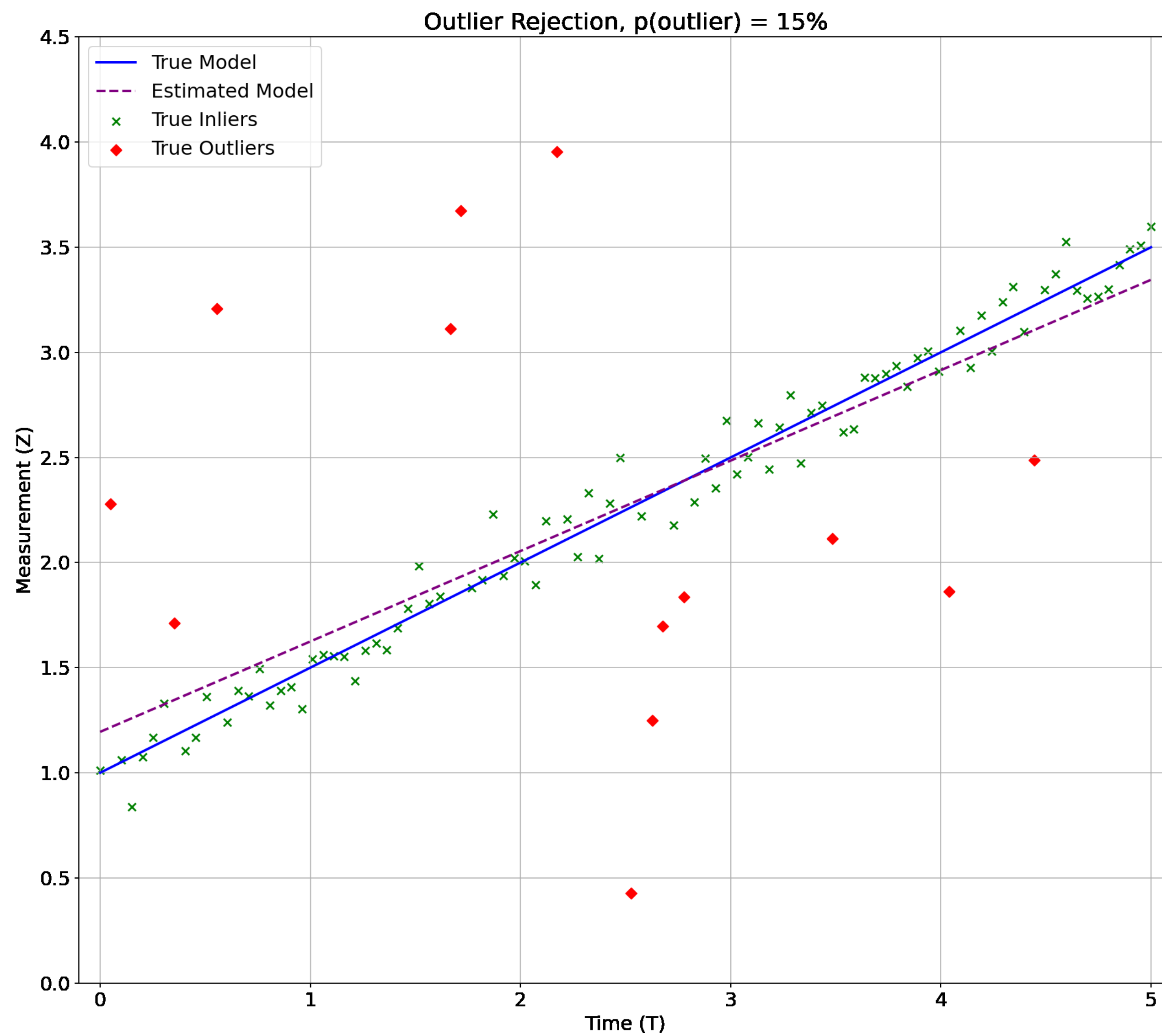
$$\mathbf{z} = h(\alpha, \beta) = \alpha t + \beta + \epsilon_z$$

↑  
Measurement

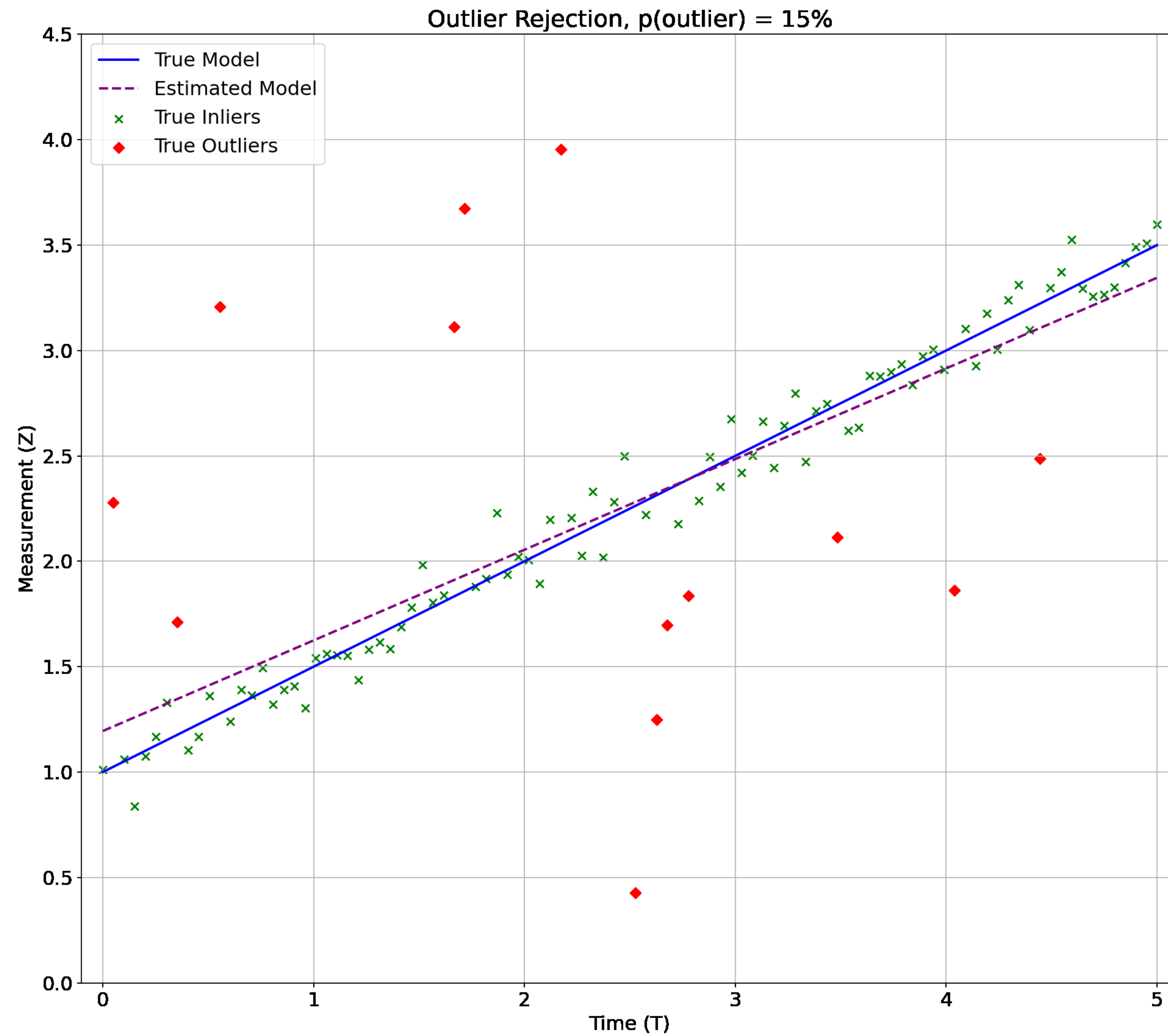
# Toy Problem



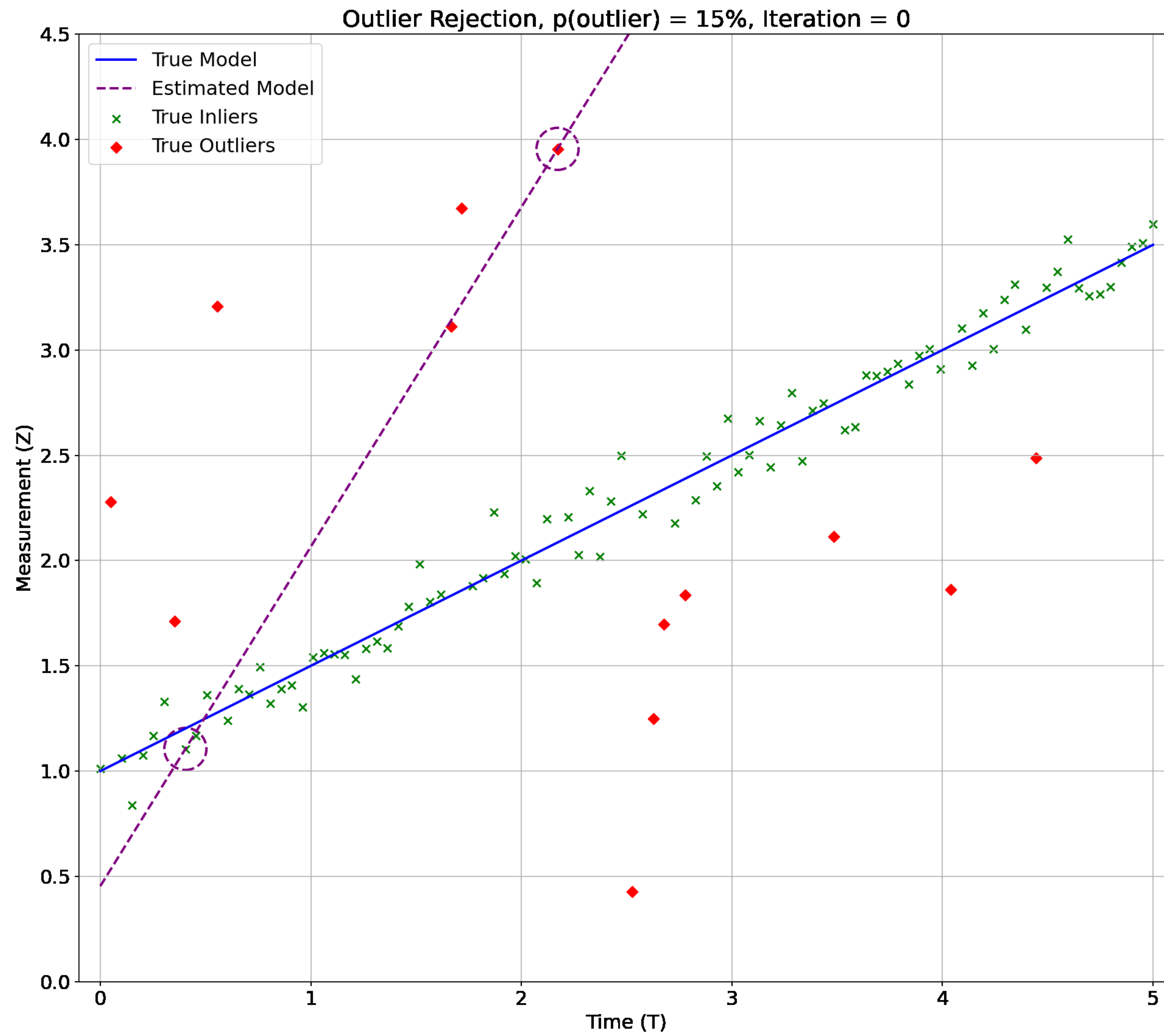
# Toy Problem



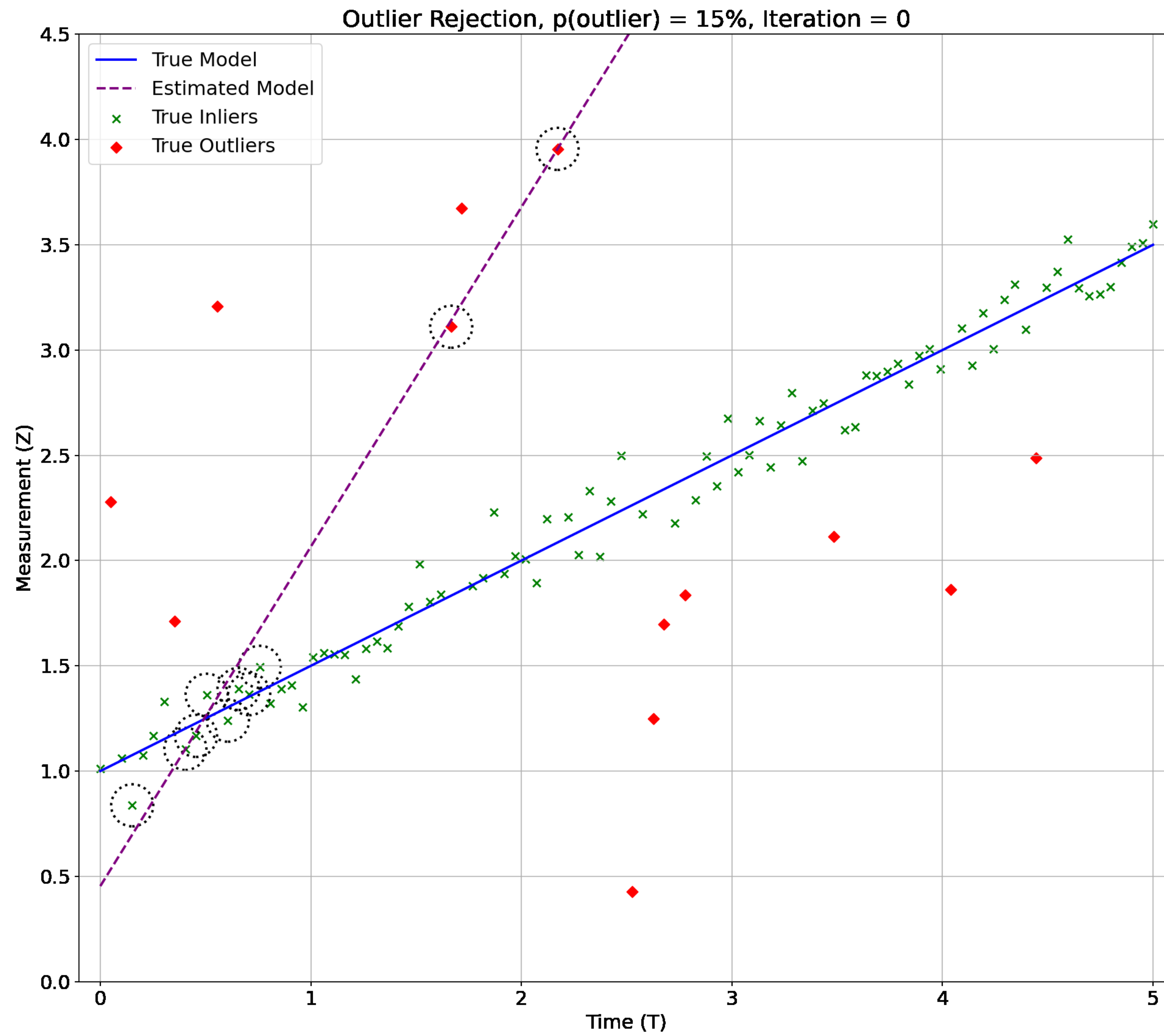
# RANSAC (Random Sample Consensus)



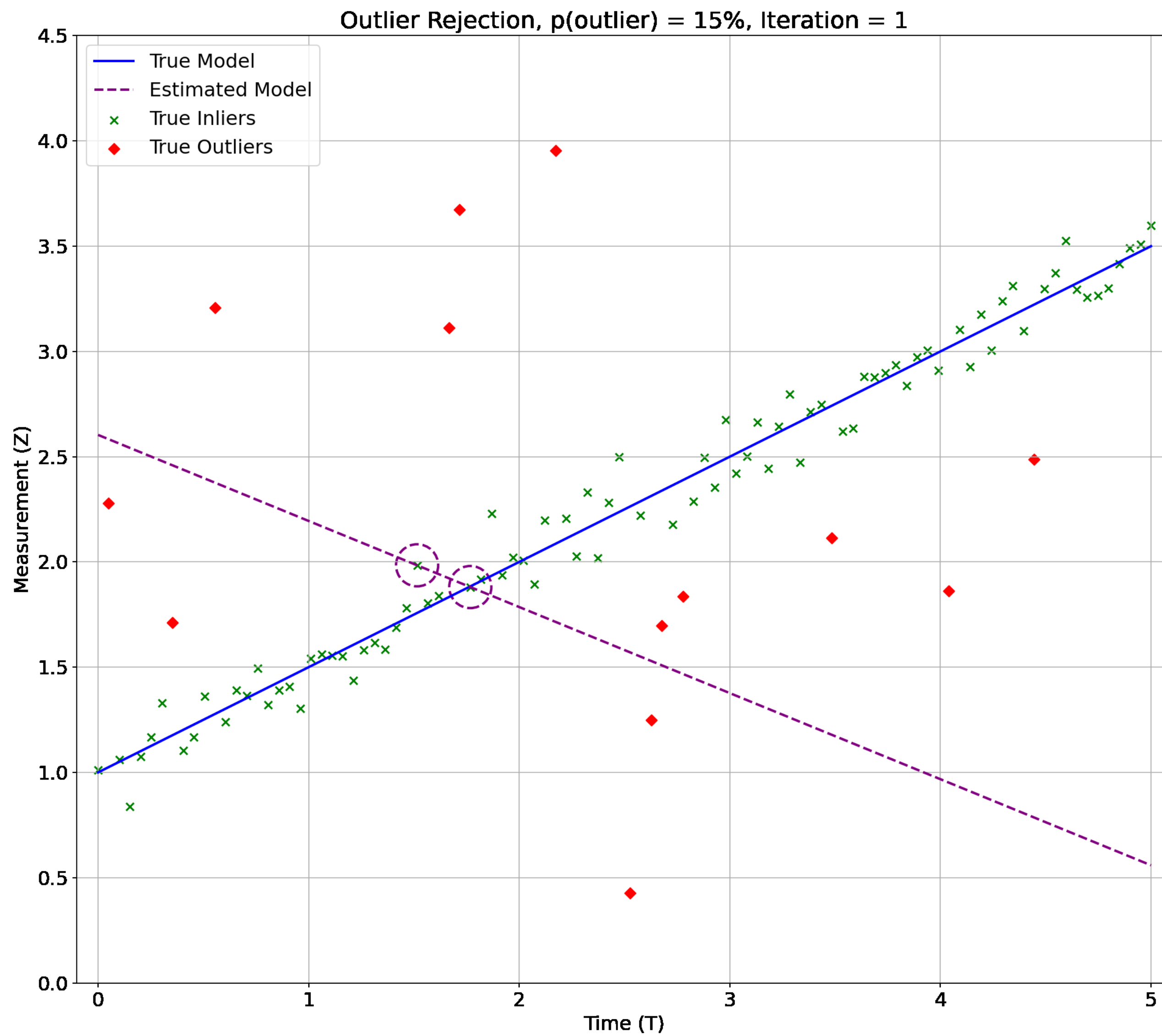
# RANSAC



# RANSAC

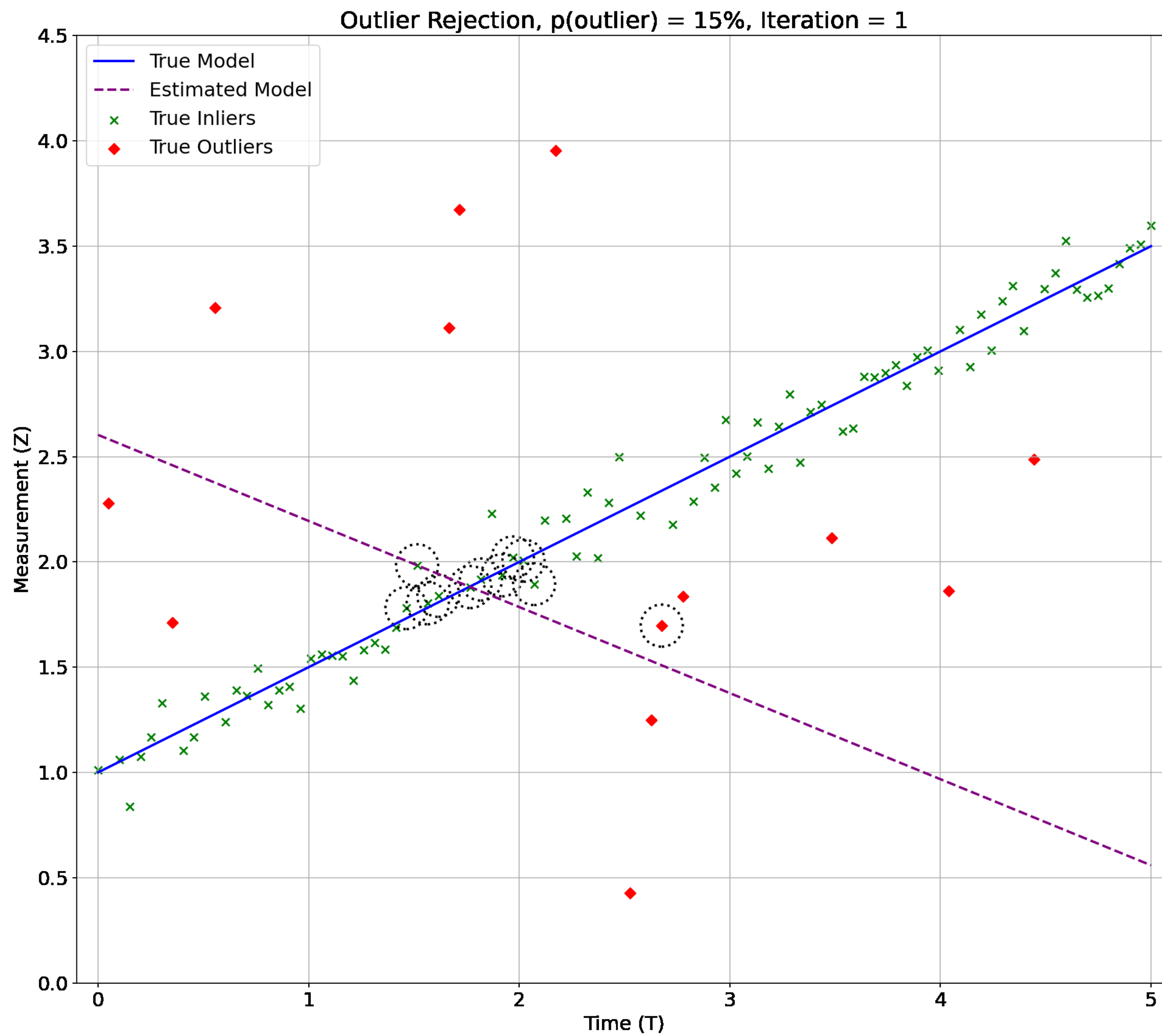


# RANSAC

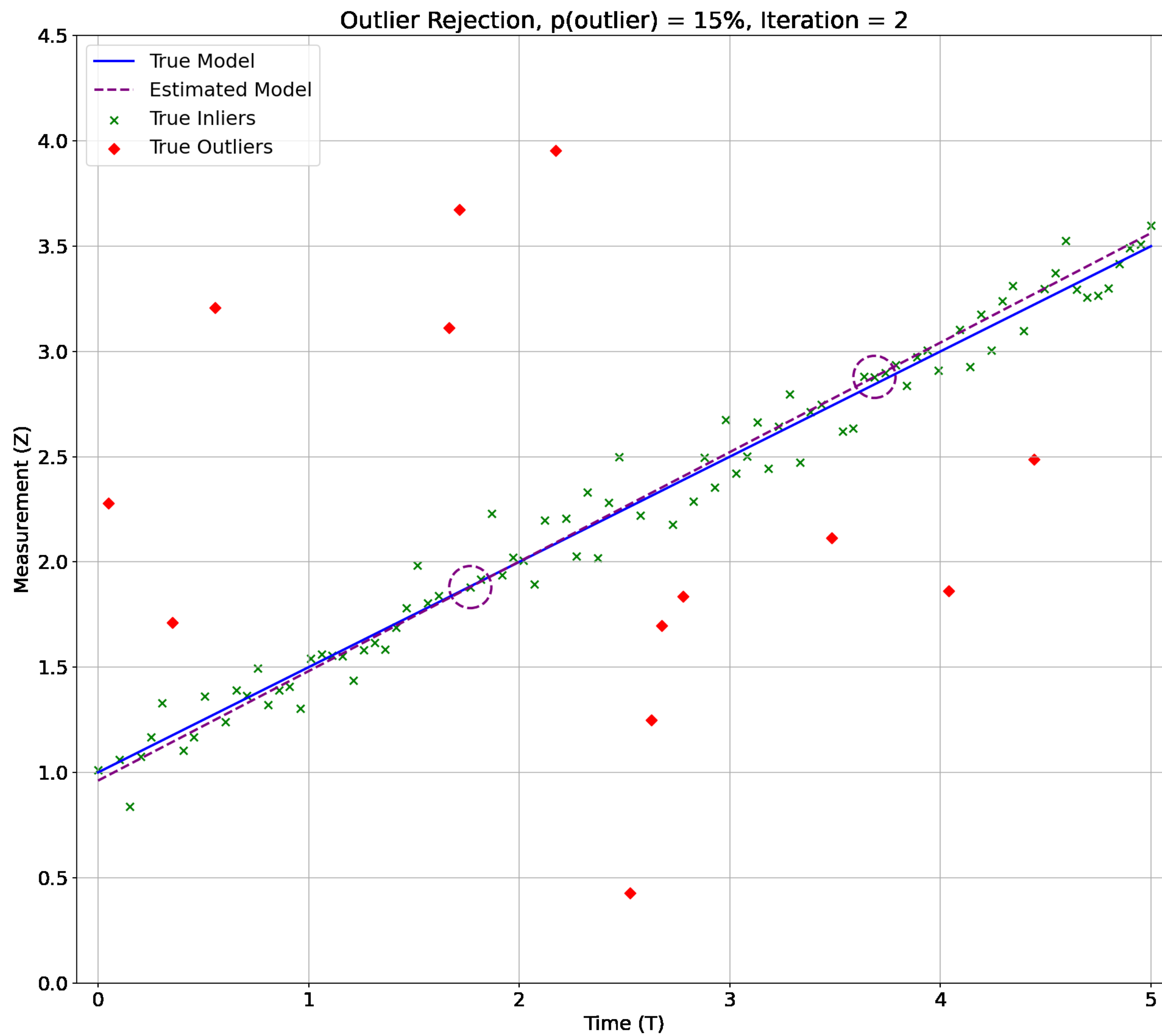


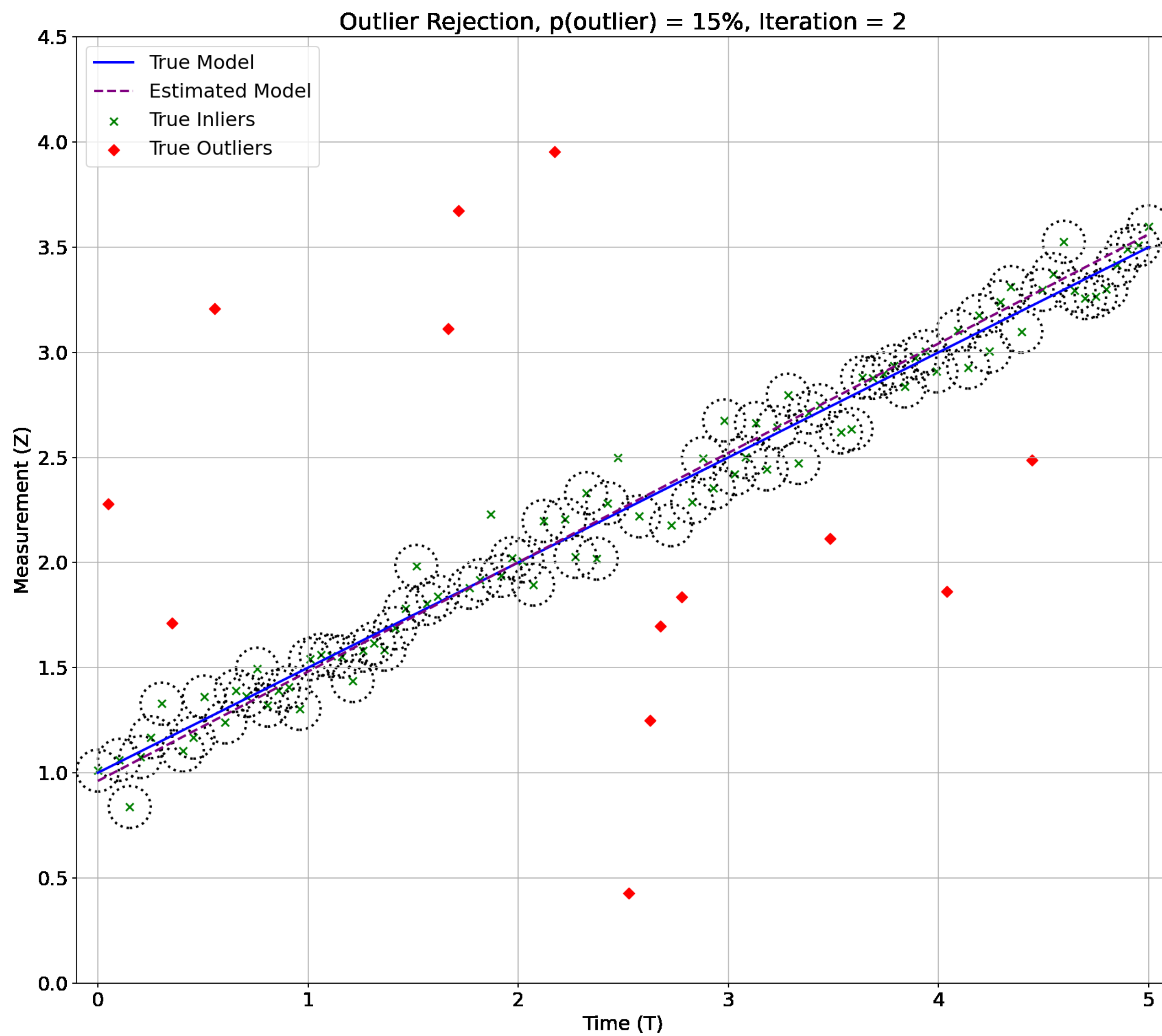


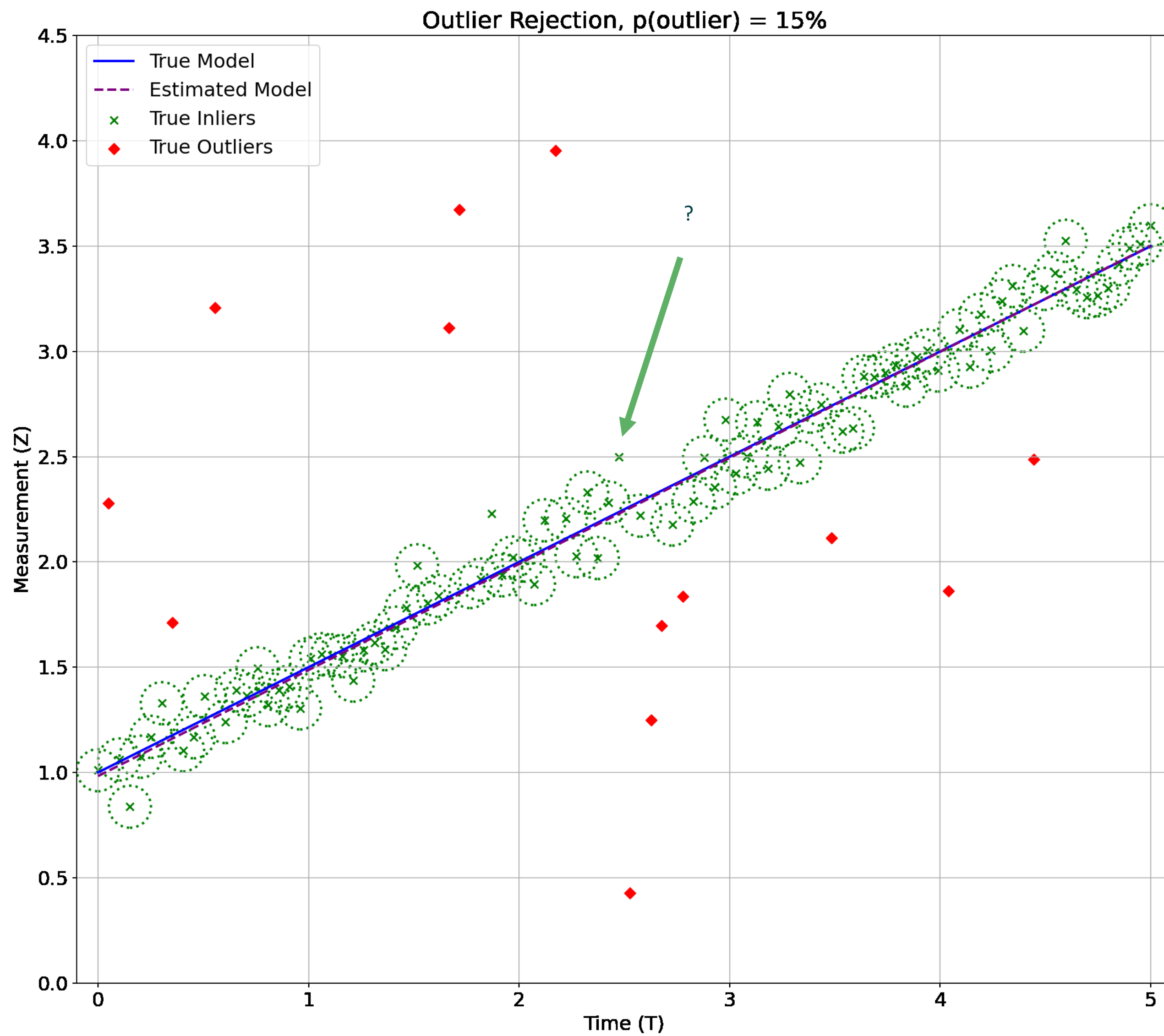
# RANSAC



# RANSAC







- Pros:
  - Dead simple to implement: Draw  $K$  examples, solve, count, repeat.
  - Easily wrap around an existing method.
  - Trivially parallelized. Have more CPU time? Sample more.
- Cons:
  - Relatively weak guarantees.
  - Can require a lot of iterations for high outlier fractions or models with a large  $K$ .
  - Hyper-parameters need tuning.

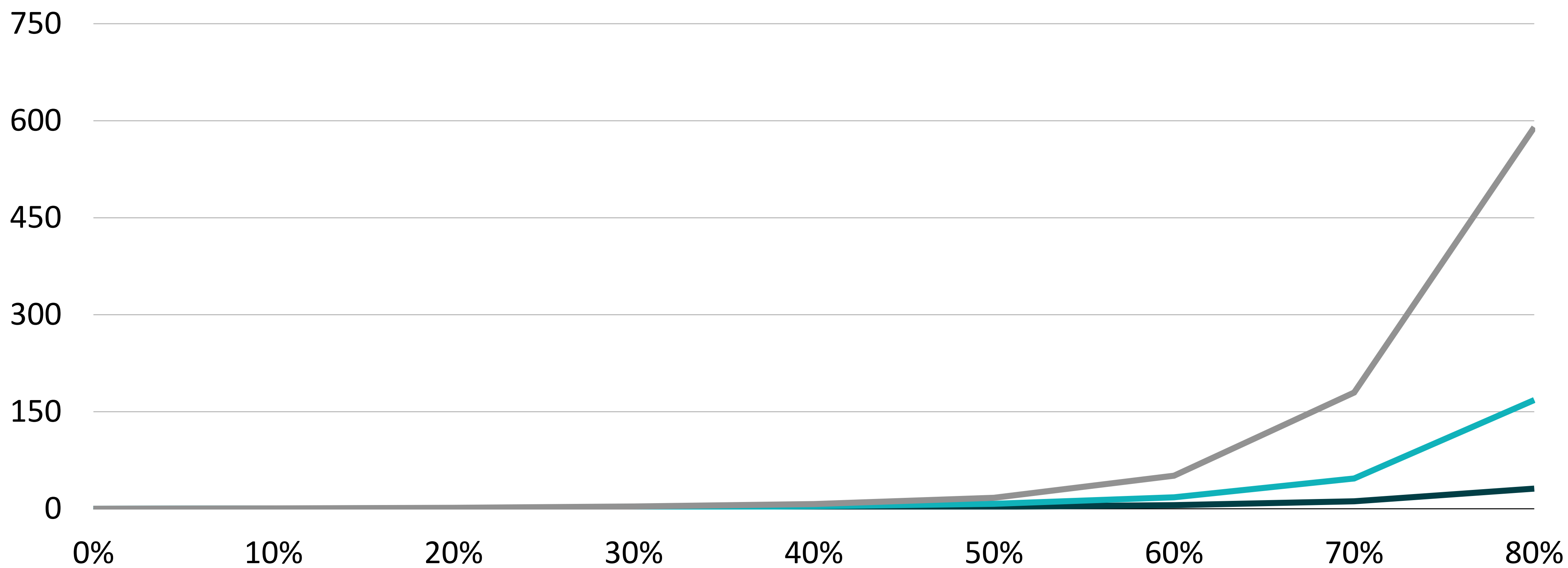
... *but*, still quite useful in practice.

— 2 Pts

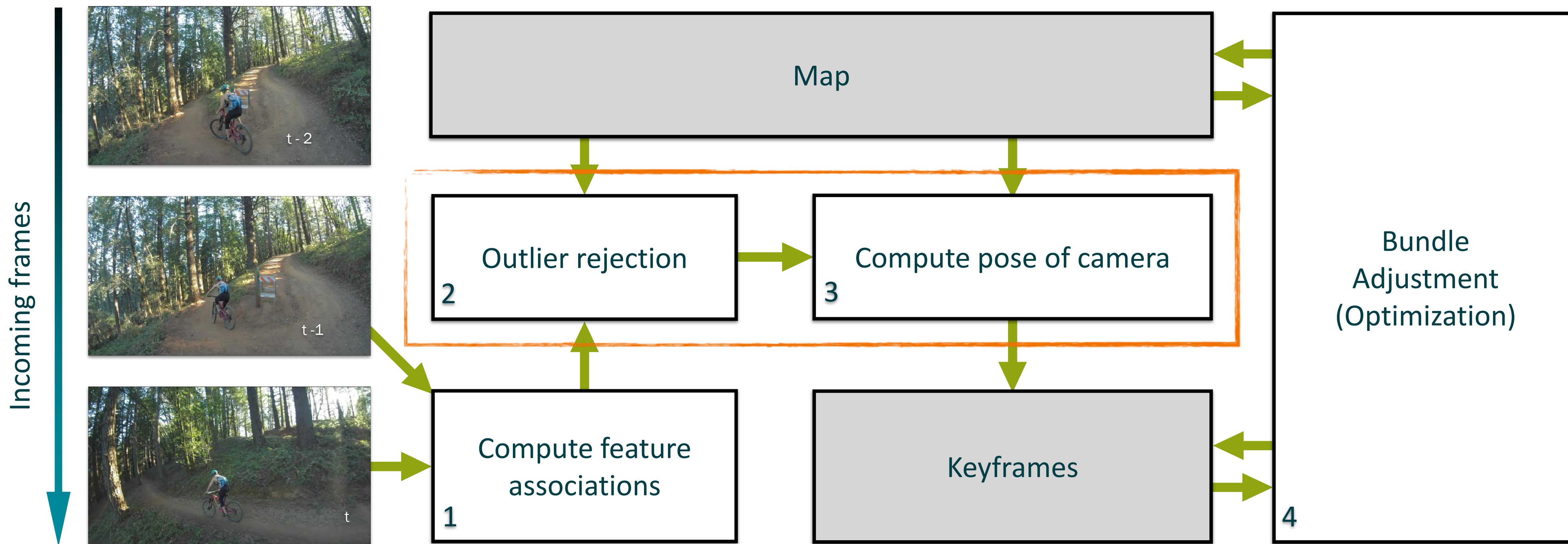
— 3 Pts

— 4 Pts

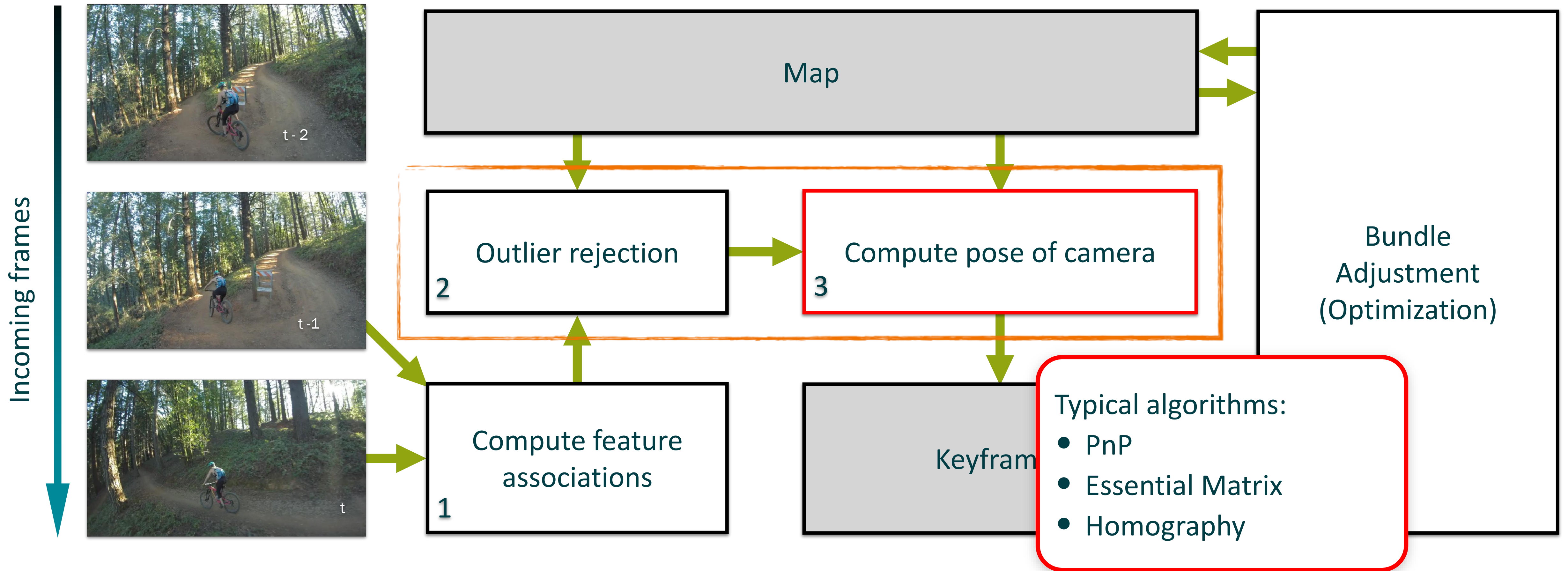
## # Iterations vs. Outlier Fraction



# Typical SLAM Pipeline w/ BA

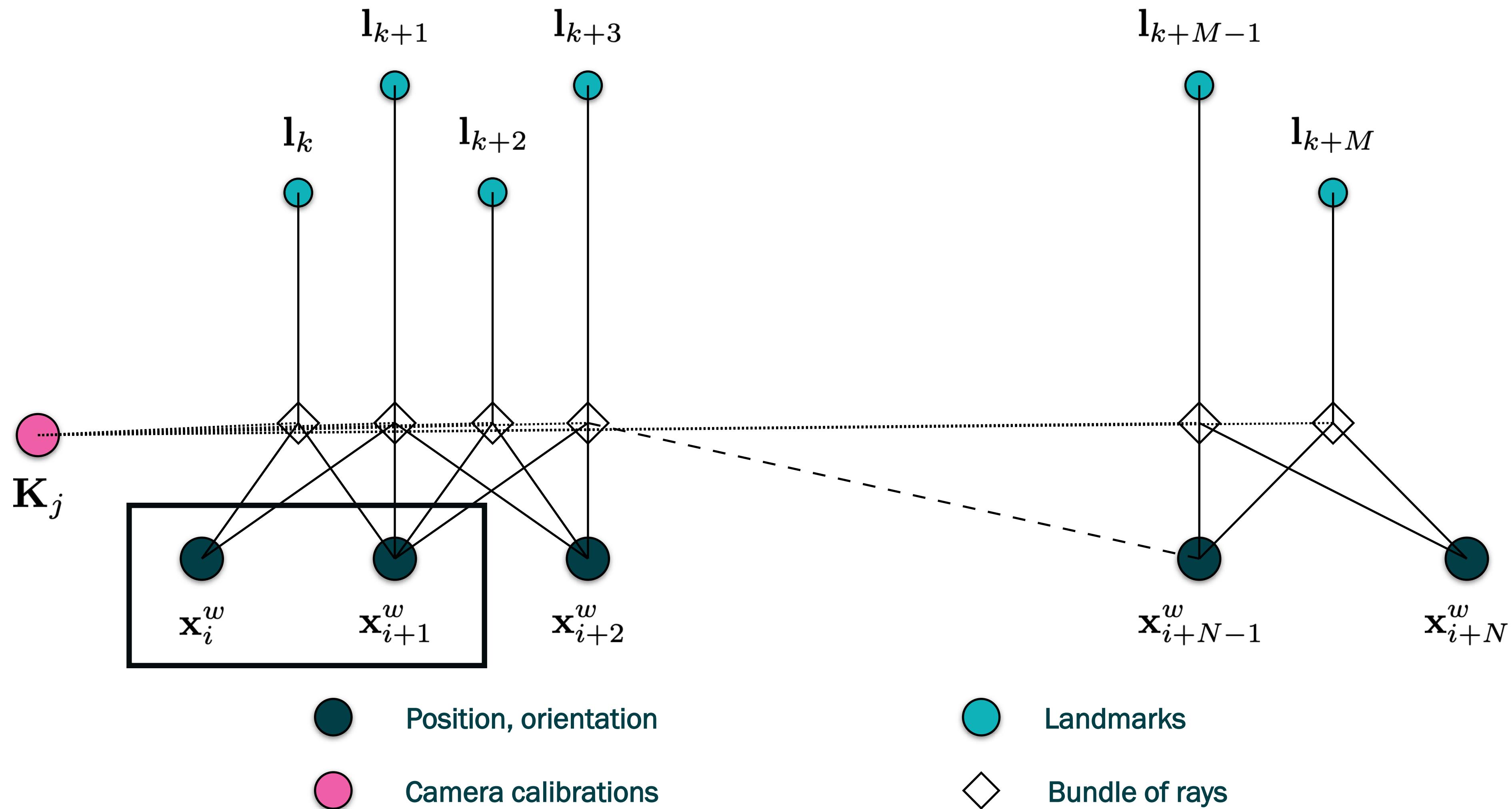


# Typical SLAM Pipeline w/ BA

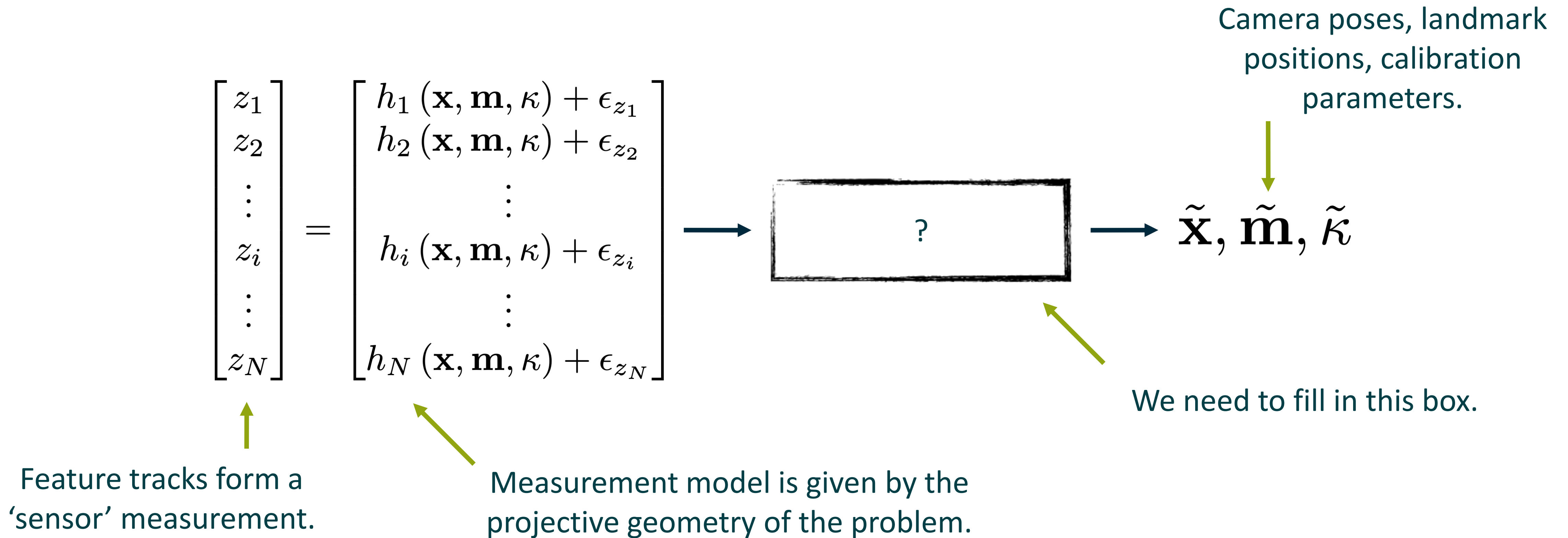




# BA as a Factor Graph



- How do we actually recover the states, given the measurements and our model?



- We can use a technique called *Nonlinear Least Squares* to do this.
- There are many ways to formulate SLAM problems generally, and we cannot review them all in the time allotted.
- *However*, this method is widely applicable, typically fast, and is straightforward to implement.
- For a much more comprehensive review, I highly recommend: [State Estimation for Robotics](#), Tim Barfoot, 2015 (Free online)

- We will convert our measurement models into a system of equations.
- Prior to that, we will make an additional assumption - that the measurement noise is drawn from a zero-mean gaussian.

$$\epsilon_z \propto N(\mu = \mathbf{0}, \Sigma_z)$$

- We will also assume we have an *initial guess* for our states. In a time recursive system, this could come from the previous frame.

We re-write our measurements as a residual functions:

$$\mathbf{z}_{ij} = h(\mathbf{x}_i^w, \mathbf{l}_j, \mathbf{K}) + \epsilon_{z_{ij}} \longrightarrow \mathbf{f}_{ij} = h(\mathbf{x}_i^w, \mathbf{l}_j, \mathbf{K}) - \mathbf{z}_{ij}$$

↖ The  $i$ 'th camera pose, observing the  $j$ 'th landmark.

And concatenate these into a large vector:

$$\mathbf{f}(\mathbf{x}^w, \mathbf{l}, \mathbf{K}) = \begin{bmatrix} \mathbf{f}_{ij} \\ \vdots \\ \mathbf{f}_{i+m, j+n} \\ \vdots \\ \mathbf{f}_{i+M, j+N} \end{bmatrix}$$

$$\|\mathbf{f}(\mathbf{x}^w, \mathbf{l}, \mathbf{K})\|_{\Sigma}^2 = \sum_i^M \sum_j^N \|\mathbf{f}_{ij}\|_{\Sigma_{z_{ij}}}^2 \delta_{ij}$$

We take the squared Mahalanobis norm, weighting by our assumed measurement uncertainty.

Our ‘best estimate’ will occur when the objective function is minimized:

$$\mathbf{y} := (\mathbf{x}^w, \mathbf{l}, \mathbf{K})$$

$$\tilde{\mathbf{y}} = \arg \min_{\mathbf{y}} \|\mathbf{f}(\mathbf{y})\|_{\Sigma}^2$$

Because  $f$  is usually going to be non-linear for most SLAM problems, we end up *linearizing* the problem and taking a series of steps.

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \delta \mathbf{y}_k$$

$$\delta \mathbf{y}_k = \arg \min_{\delta \mathbf{y}} \|\mathbf{f}(\mathbf{y}_k) + \mathbf{J} \delta \mathbf{y}\|_{\Sigma}^2$$

Jacobian  $\mathbf{J}$  is the linearization of  $f$  about our initial guess.

The solution at each iteration:

$$\delta \mathbf{y}_k = \left( \mathbf{J}^T \boldsymbol{\Sigma}_z^{-1} \mathbf{J} \right)^{-1} \mathbf{J}^T \boldsymbol{\Sigma}_z^{-1} \mathbf{f}(\mathbf{y}_k)$$

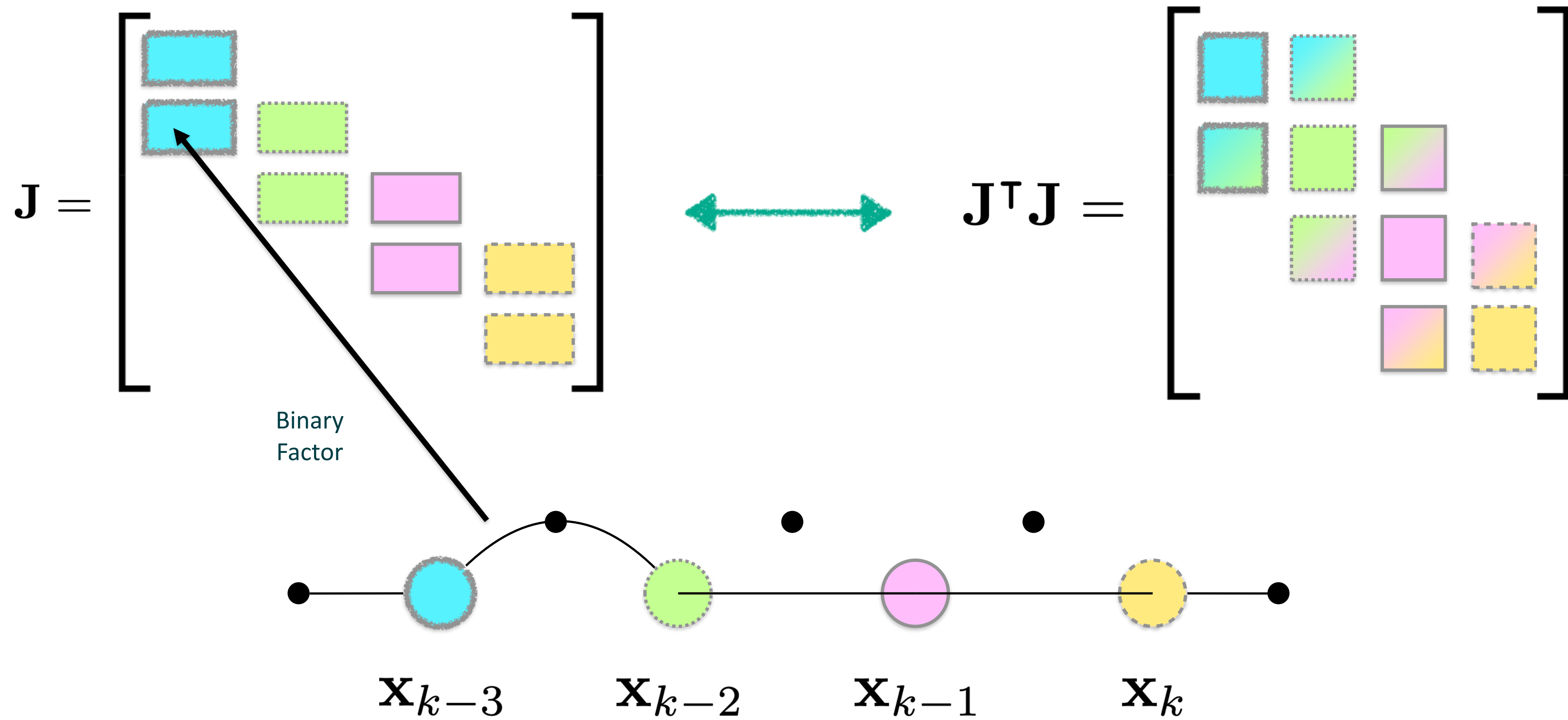
First order approximation of the Hessian. Inversion has complexity  $O(|y|^3)$

Each residual is weighted by its inverse uncertainty.

When linearized about the converged solution, the inverted Hessian doubles as a first order approximation of the *marginal covariance* of our estimate: \*

$$\boldsymbol{\Sigma}_{\tilde{\mathbf{y}}} \approx \left( \mathbf{J}^T \boldsymbol{\Sigma}_z^{-1} \mathbf{J} \right)^{-1}$$

\* See [Barfoot](#), Chapters 3 and 4.





- In the *linearized* form, the problem is ‘easy’ to solve.
  - Reduces to iterated application of *weighted least squares*.
  - Generally, cost of solving for updates is *cubic* in the number of states:
    - However, in some problems (like BA) there is sparsity we can leverage to improve this.
- Huge number of problems can be cast this way (given an initial guess).
- Can run in a fixed memory footprint → suitable for embedded use case.
- With the appropriate  $\Sigma$  weights we can show the NLS produces an *approximate* estimate of the uncertainty in our solution.

- Remember our assumptions:
  - We needed an initial guess to linearize the system. If the guess is poor, the gradient used in the optimizer will steer our solution in the wrong direction.
  - Additionally, the covariance estimate we get out is only as good as the linearization point.
- We also assumed Gaussian noise on the measurements.
  - Outliers must be removed, or they will dominate the optimization.

- It is worth considering the effect of linearization on our uncertainty estimate.
- For a Gaussian variable  $u$  and non-linear vector function  $g$ , we can approximate:

$$\mathbf{u} \propto N(\mu_{\mathbf{u}}, \Sigma_{\mathbf{u}})$$

$$\mathbf{v} = g(\mathbf{u})$$

Because we linearized, the fidelity of our first-order  $\Sigma$  relies on this approximation.

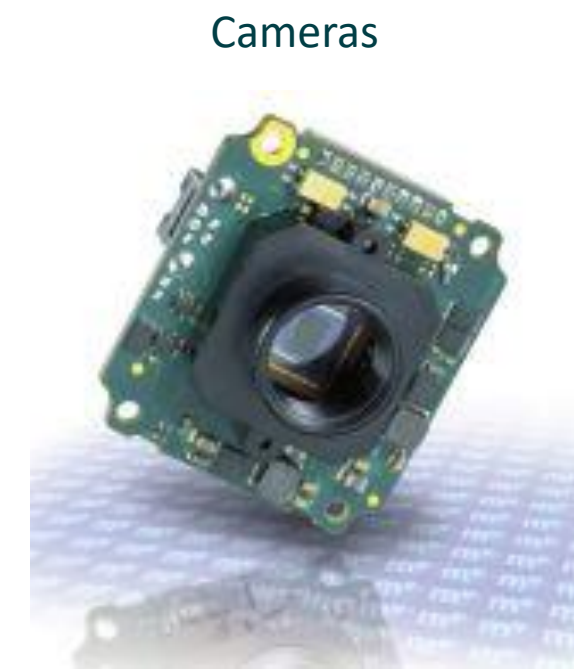
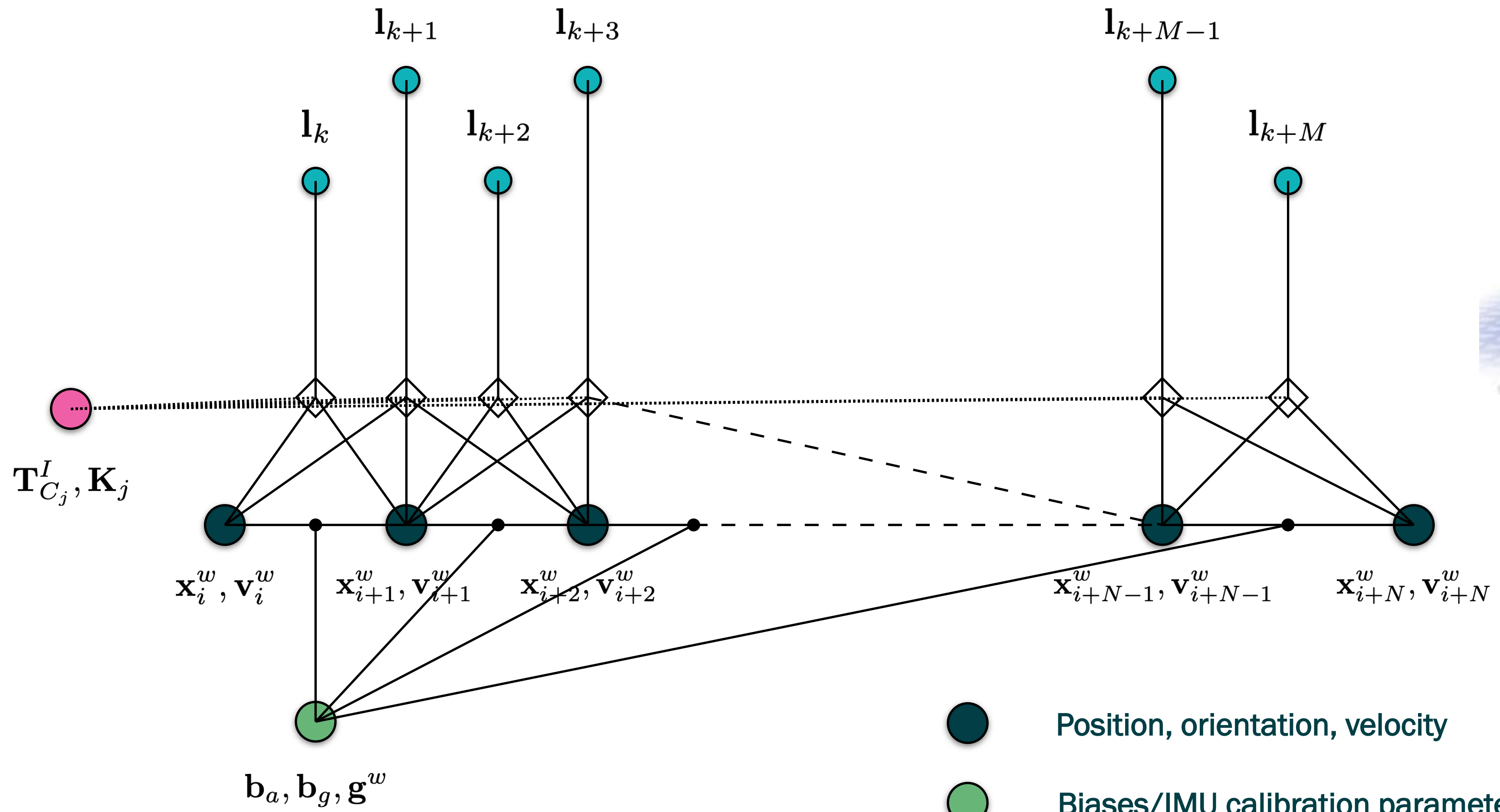
$$\mathbf{v} \approx N\left(g(\mathbf{u}), \frac{\delta g}{\delta \mathbf{u}} \Sigma_{\mathbf{u}} \frac{\delta g}{\delta \mathbf{u}}^T\right)$$

See [Barfoot](#), Chapter 2.

- Some relevant tools:
  - [GTSAM](#), open source package created by [Frank Dellaert](#) et al.
    - Allows specification of problem in factor graph format, built for SLAM.
  - [G2O](#)
    - Includes solutions for SLAM and BA.
  - [Ceres Solver](#), produced by Google
    - General non-linear least-squares optimizer.
  - Python
    - [scipy.optimize.least\\_squares](#)

- BA can operate at small and large scale.
  - Small: A few image frames on a mobile phone.
  - Large: Tens of thousand of images at city-scale.
- Fairly straightforward to implement.
- But:
  - Robust association *may* require expensive descriptors.
  - After feature association, we must devote nontrivial compute to outlier rejection.
  - Update rate limited to camera frame rate (slow).

# VIO as a Factor Graph: BA + IMU



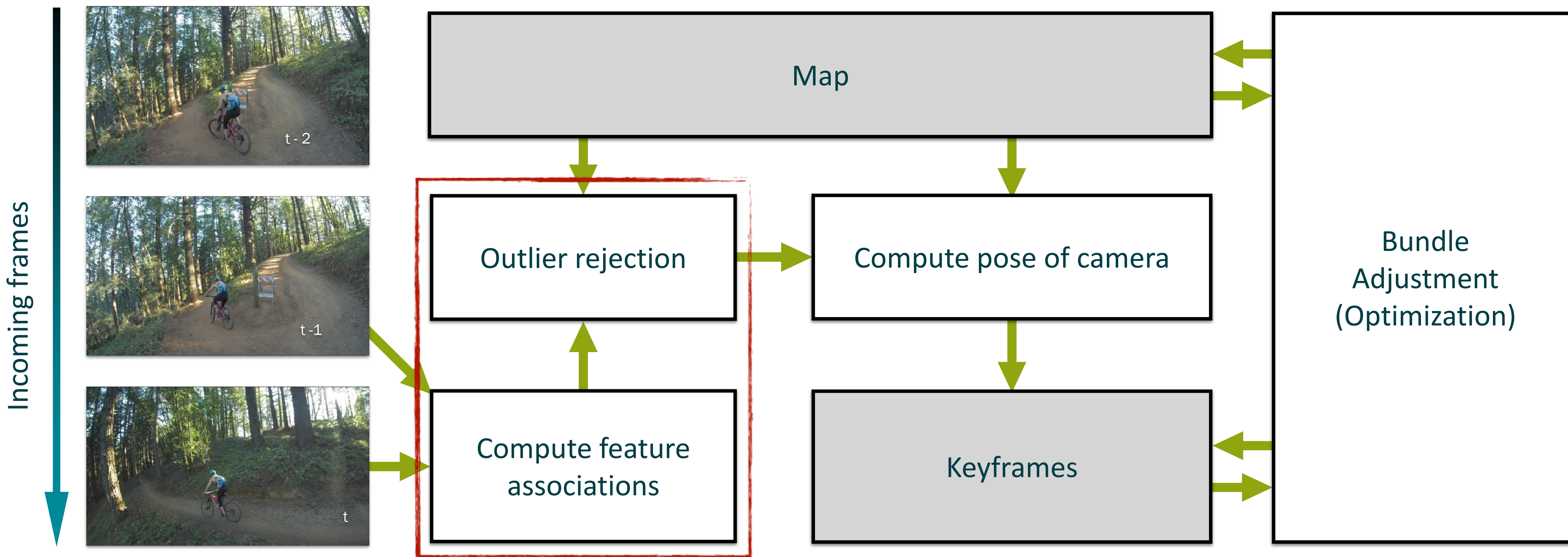
Source: [MatrixVision](#)



Source: [Lord MicroStrain](#)

- Position, orientation, velocity
- Landmarks
- Biases/IMU calibration parameters
- Camera calibrations
- ◇ Bundle of rays
- IMU, motion model

- One of the most successful adaptations of vision research to the market.
  - Present in smart phones, AR/VR headsets, drones, autonomous vehicles.
- Camera and IMU are highly complementary:
  - Camera:
    - Low update rate, high compute cost, subject to outlier data.
    - Able to relocalize accurately at large distances.
  - IMU:
    - High update rate, low compute cost, few outliers (maybe saturation).
    - Accurate over short intervals, but drifts over time.
    - Able to recover attitude with respect to global reference frame (gravity).



IMU can deliver substantial value here.



- Existing open-source implementations (not exhaustive):
  - [OpenMVG](#)
  - [COLMAP](#) Offline SFM and Multi-view Stereo (MVS)
  - [CMVS](#) Multi-view Stereo
  - [ORB-SLAM2](#) Real-time SLAM featuring BA optimization
  - [PTAM](#) One of the earliest functional visual-SLAM demos
  - [VINS-Mono](#) VIO, runs on a mobile device
  - [Basalt](#) VIO
  - [ROVIO](#) VIO, example of a *direct method*

- Additional Reading:
  - [State Estimation for Robotics](#) (Barfoot, 2015)
  - [Factor Graphs for Robot Perception](#) (Dellaert and Kaess, 2017)
  - [Visual Odometry](#), (Scaramuzza and Fraundorfer, 2011)
  - [Probabilistic Robotics](#), (Thrun, Burgard, and Fox, 2005)
  - [GTSAM](#) Software Library
  
- Questions? Feel free to reach out: [gareth@skydio.com](mailto:gareth@skydio.com)