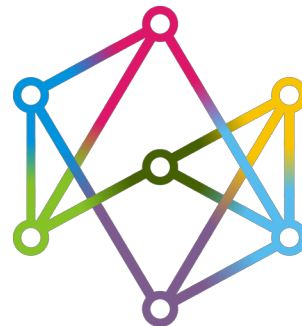


Vision Aided Navigation and SLAM: Overview II

Vadim Indelman



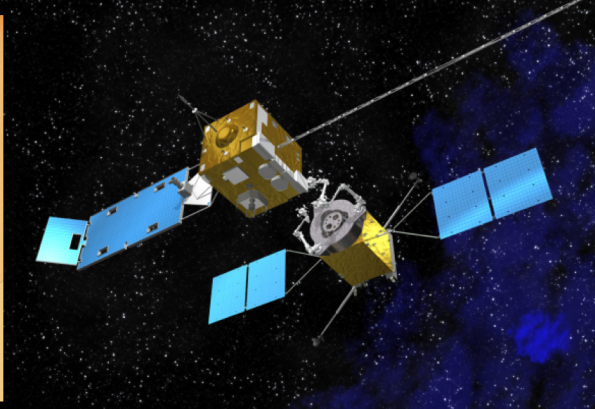
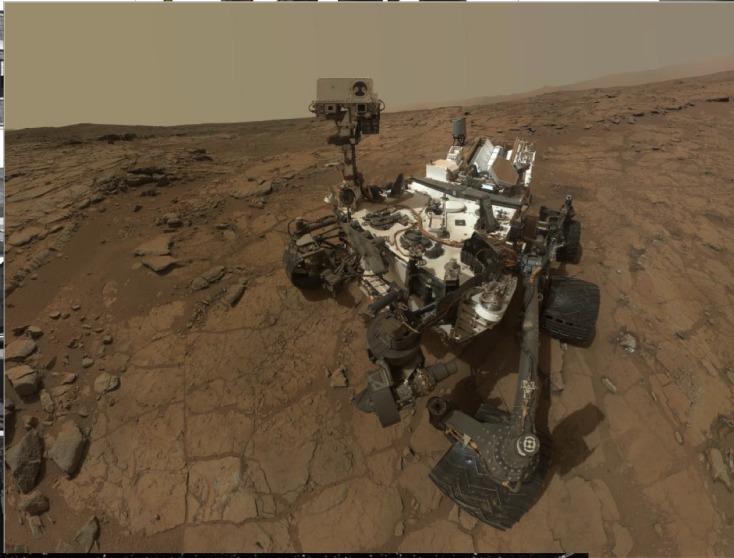
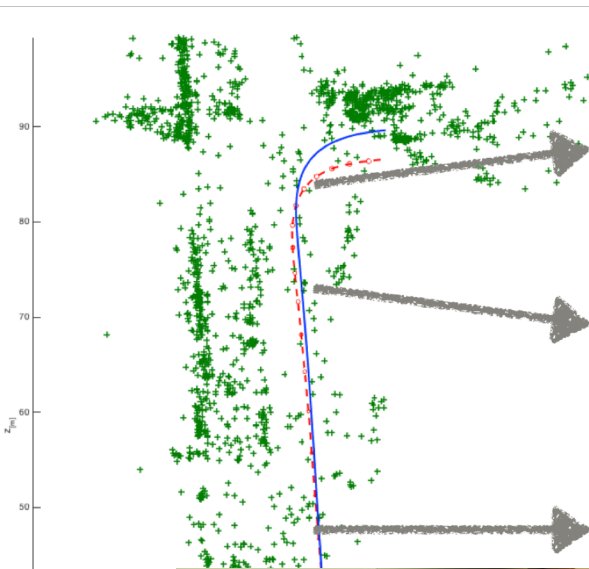
TECHNION
Israel Institute
of Technology



ANPL
Autonomous Navigation and
Perception Lab

Introduction

- Numerous applications



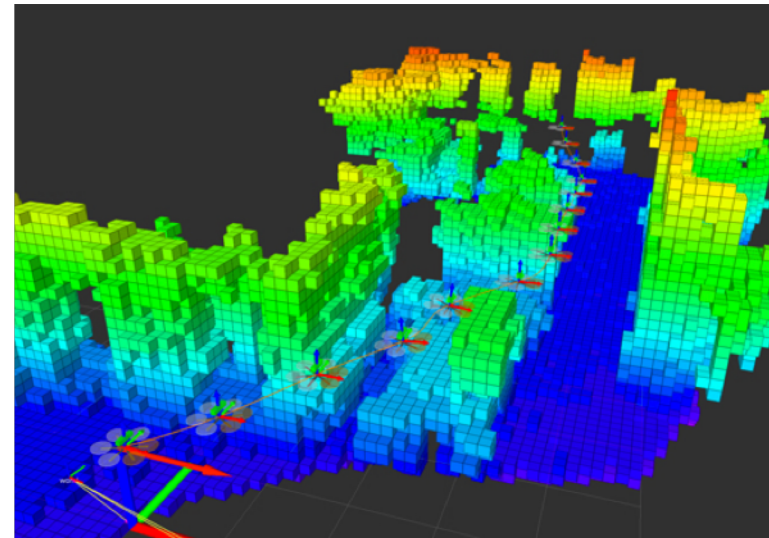
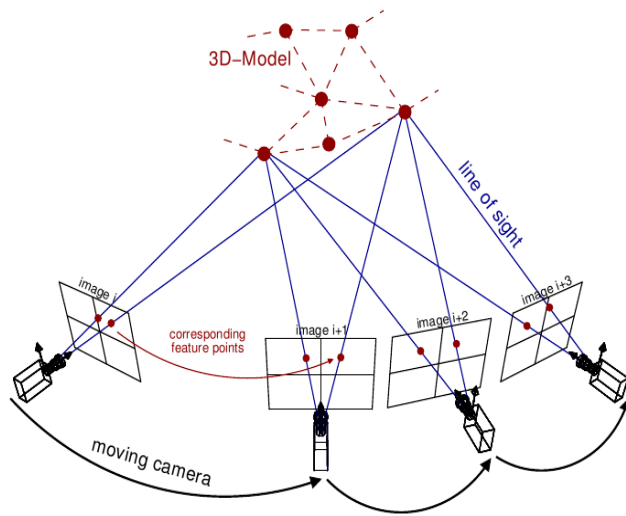
Copyright KMeL Robotics

kmelrobotics@gmail.com

Introduction

Autonomous navigation and perception in uncertain/unknown environments:

- **Perception and Inference:** Where am I? What is the surrounding environment? **This talk**
- **Planning Under Uncertainty & Active Perception:** Decide next action(s) given partial, noisy data



The “Big” Picture

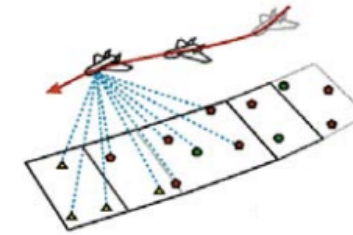
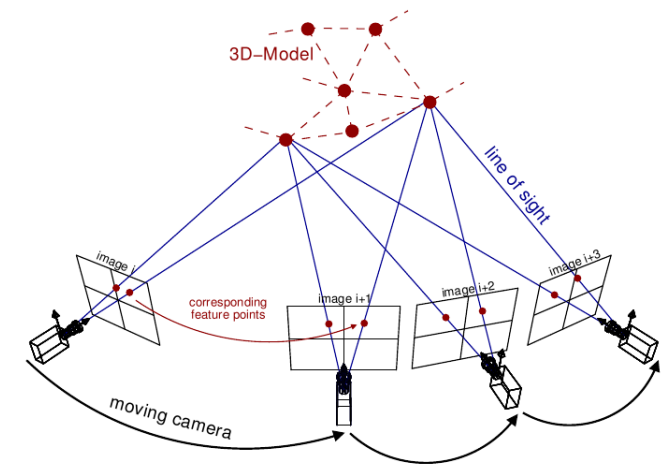


Image from [Kim et al. 2005]

- **Objective:**

- Estimate platform state and observed environment (e.g. 3D points)
- Environment is unknown, uncertain or dynamic



The “Big” Picture

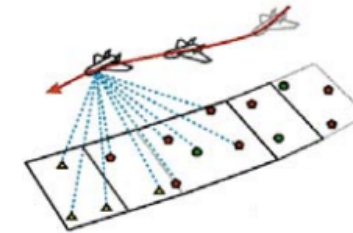
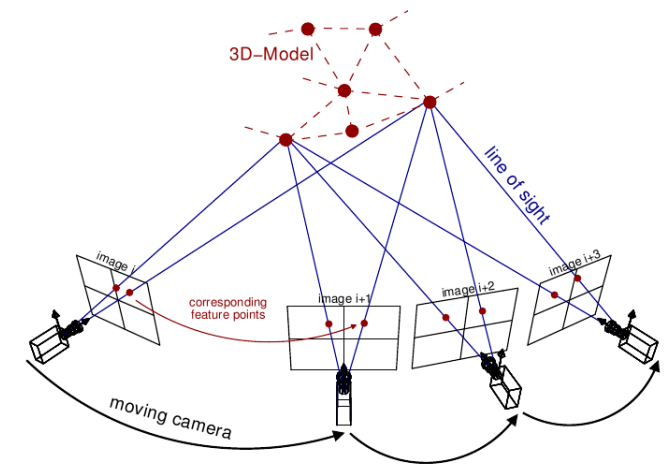


Image from [Kim et al. 2005]

- How?
 - Many images (sensor measurements)
 - Interest points (features) in each image
 - Track features from image to image, data association
 - Probabilistic inference over robot state and environment (e.g. 3D points)
- Additional sensors
- Online performance?

Front-end

Back-end



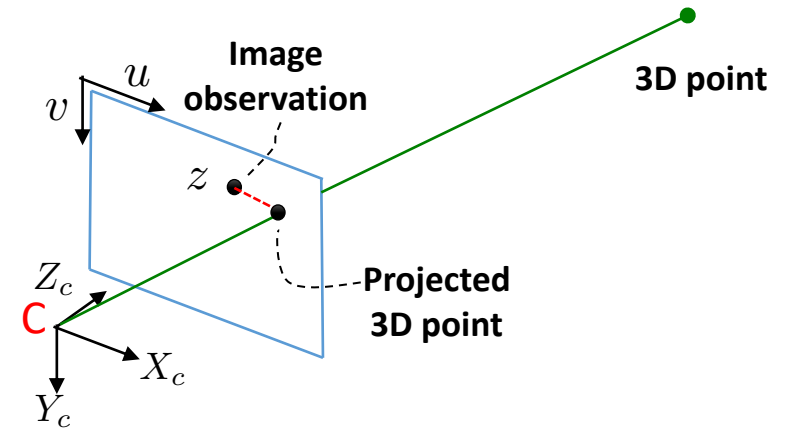
Outline

- Introduction
- Camera projective geometry
- Bundle Adjustment
- Incremental Smoothing and Mapping (iSAM) algorithms
- Visual-inertial SLAM and IMU pre-integration concept

Projection Matrix & Operator

- Projection of a 3D point $\mathbf{X}^w = (X_w, Y_w, Z_w)$:

$$\begin{pmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{pmatrix} = \underbrace{K}_{\text{Projection matrix}} \underbrace{[R \mid t]}_{\text{Camera pose}} \underbrace{\begin{pmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{pmatrix}}_{\text{3D point, landmark}}$$



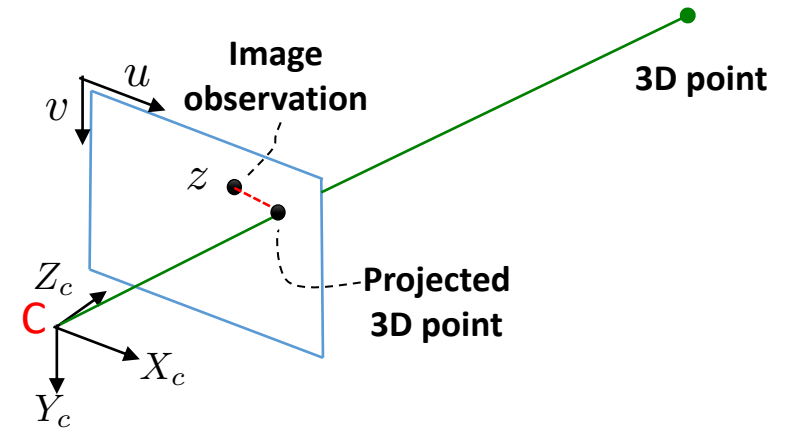
- To recover the pixel (u, v) : $\begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} \tilde{u} \\ \tilde{v} \end{pmatrix}^T$

- Projection of a 3D point (notation): $\pi(x, l) \doteq K [R \ t] l$

Known as
projection operator

Re-Projection Error

- Ideally:
$$\frac{(u, v)}{z} = \pi(x, l)$$
- In practice:
 - Image observations z are noisy
 - Incorrect/imprecise camera pose and 3D point



- Re-projection error:

$$\frac{z - \pi(x, l)}{\text{Image observation (pixel, feature)} \quad \text{Predicted measurement}}$$

- Assuming Gaussian image noise, measurement likelihood:

$$p(z|x, l) = \frac{1}{\sqrt{|2\pi\Sigma_v|}} \exp\left(-\frac{1}{2} \frac{\|z - \pi(x, l)\|_{\Sigma_v}^2}{\text{Re-projection error}}\right) \quad v \sim N(0, \Sigma_v)$$

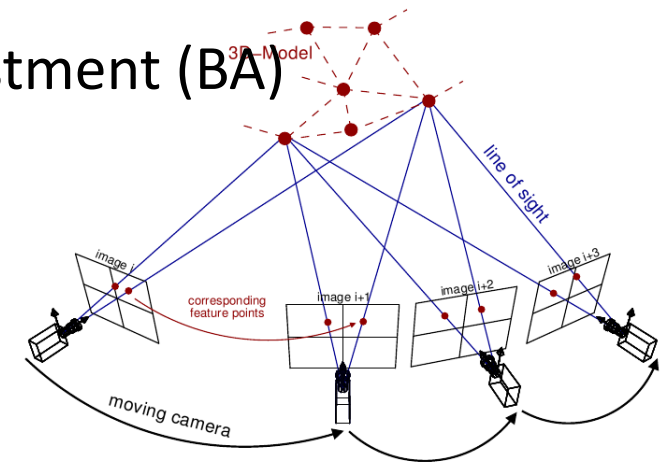
Re-projection error

Outline

- Introduction
- Camera projective geometry
- Bundle Adjustment
- Incremental Smoothing and Mapping (iSAM) algorithms
- Visual-inertial SLAM and IMU pre-integration concept

Visual SLAM and Bundle Adjustment

- Assume we are given a sequence of images
- **Objective:** Would like to infer camera poses and observed 3D points
 - Using **only** images as input (no additional sensors, **for now**)
 - Assume data association is given (very challenging by itself!)
- Problem known as
 - Computer vision: Structure from motion (SfM), Bundle adjustment (BA)
 - Robotics: Simultaneous localization and mapping (SLAM)



Bundle Adjustment

Consider the i -th image:

- A single image observation of a 3D point l_j

$$z_{i,j} = \pi(x_i, l_j) + v \quad \longrightarrow \quad p(z_{i,j} | x_i, l_j)$$

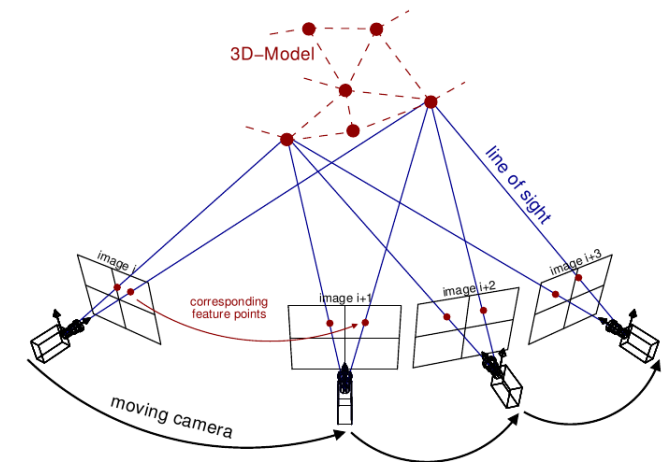
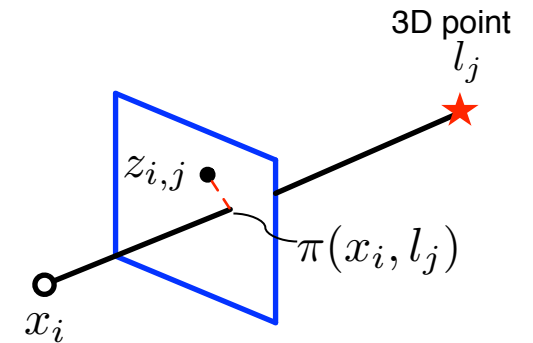
- Notations

- \mathcal{M}_i : indices of 3D points observed in image i
- Z_i : all image observations from image i

- The joint pdf over camera pose and observed 3D points:

$$p(x_i, \{l_j | j \in \mathcal{M}_i\} | Z_i) = \eta \prod_{j \in \mathcal{M}_i} p(z_{i,j} | x_i, l_j)$$

All 3D points
observed in image i

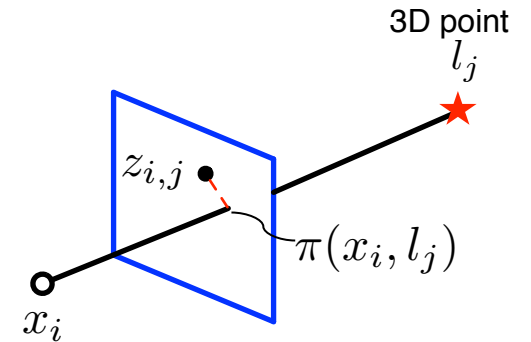


Bundle Adjustment

X : all camera poses (or platform states)
 L : observed 3D points (in any of the images)

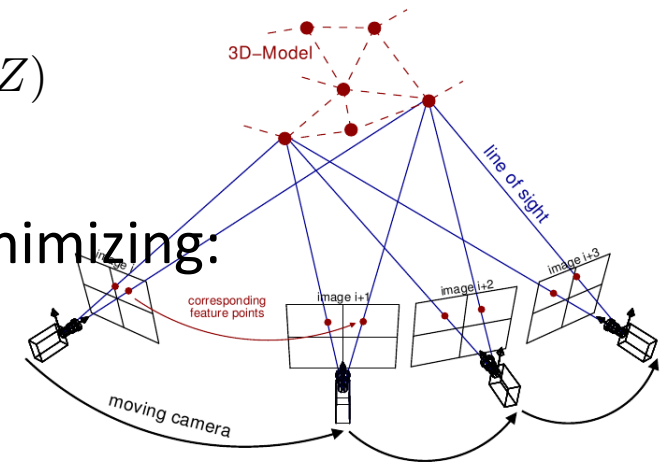
- Consider now N images
- Joint pdf for all image observations, in all camera frames:

$$p(X, L|Z) \propto \prod_i^N \prod_{j \in \mathcal{M}_i} p(z_{i,j} | x_i, l_j)$$



- Maximum a posteriori (MAP) estimate for $X^*, L^* = \arg \max_{X, L} p(X, L|Z)$
- Assuming Gaussian measurement likelihood – equivalent to minimizing:

$$J_{BA}(X, L) \doteq \sum_{i=1}^N \sum_{j \in \mathcal{M}_i} \|z_{i,j} - \pi(x_i, l_j)\|_{\Sigma}^2$$



- Approaches: Gauss-Newton, Levenberg-Marquardt, ...

Outline

- Introduction
- Camera projective geometry
- Bundle Adjustment
- Incremental Smoothing and Mapping (iSAM) algorithms
- Visual-inertial SLAM and IMU pre-integration concept

Back to visual SLAM & Vision Aided Navigation

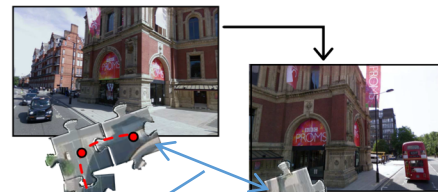
- SfM, BA
 - Sensors: camera (monocular/stereo)
 - Images can be unordered (e.g. downloaded from internet)
 - Cameras are sometimes uncalibrated
- SLAM & VAN
 - Variety of sensors: camera (monocular/stereo), laser scanner, odometry (IMU, wheel ..)
 - Imagery typically arrives in order (sequential)
 - **Online** operation is required

Loop Closure Observations

- Loop closure observations: Re-observation of a scene
- Essential for reducing drift – resets estimation errors to prior levels
- Challenging to identify



(a) Robust local motion estimation



(b) Mapping and loop-closure detection

The same scene!!

Estimation error

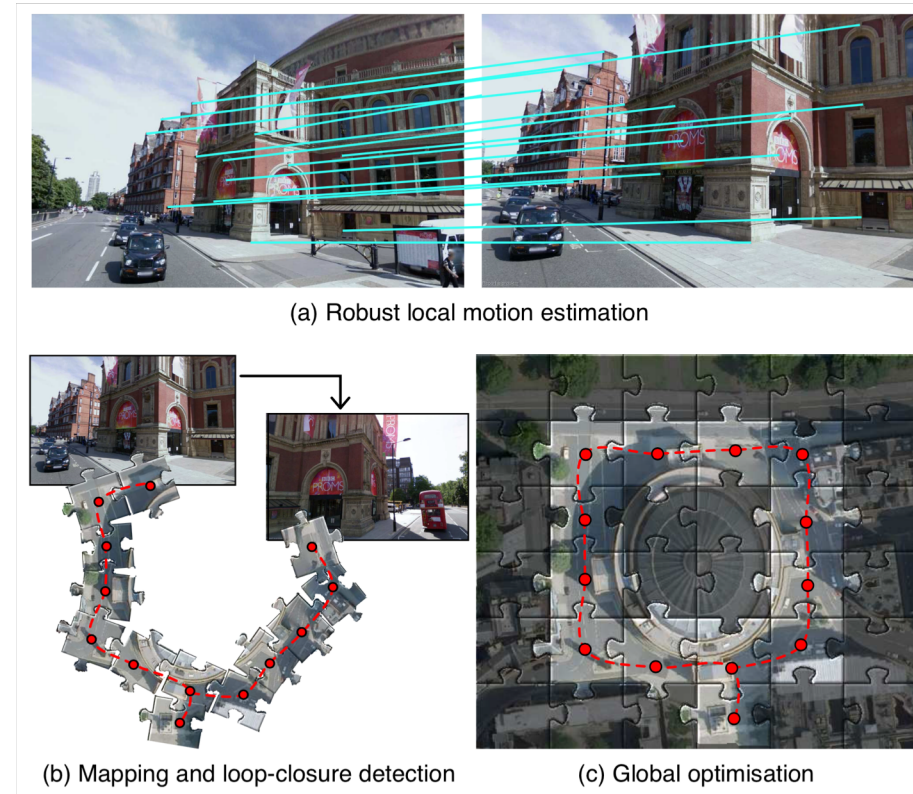


Loop Closure Observations

- Loop closure observations: Re-observation of a scene
- Essential for reducing drift – resets estimation errors to prior levels
- Challenging to identify
- Measurement equation:

$$z_{k,j} = h(x_k, l_j) + v$$

↑
3D point that has been observed
some time in the past



SLAM – Main Approaches

- Approaches differ in
 - Inference/Filtering techniques
 - Definition of what is estimated (state vector)

Common Inference Approaches

- EKF
- EIF (information form)
- Sparsity-aware optimization
- Particle filters

Latent Variables (State vector)

- Current state (e.g. pose) + landmarks
- Current state + past poses + landmarks
- Current state + past poses

Full SLAM

Pose SLAM

- Deep learning ...

Square Root Smoothing and Mapping (SAM)

Dellaert IJRR 2006

Smoothing and Mapping (SAM)

- Camera (or laser) observation model and likelihood

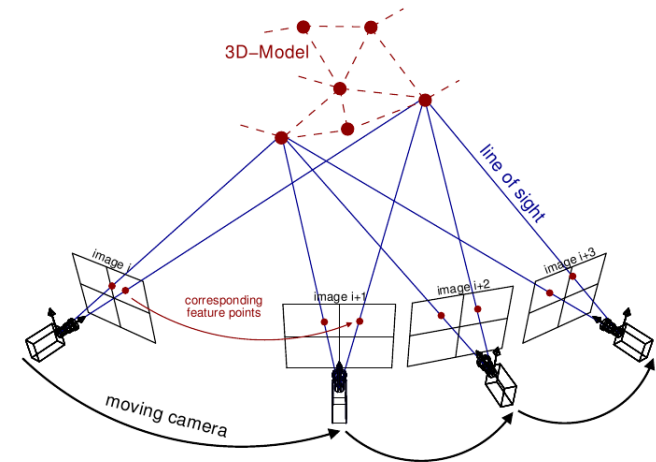
$$z_{k,j} = h(x_k, l_j) + v$$

$$p(z_{k,j} | x_k, l_j)$$

- 3D points l_j are random variables (unknown/uncertain environment)
- Need to be estimated, part of the inference process (similarly to SfM)

- Joint pdf:

$$p(x_{0:k}, L_k | u_{0:k-1}, z_{0:k})$$



Smoothing and Mapping (SAM)

- Full joint pdf:

$$p(x_{0:k}, L_k | u_{0:k-1}, z_{0:k}) = \eta p(x_0) \prod_i \left[p(x_i | x_{i-1}, u_{i-1}) \prod_{j \in \mathcal{M}_i} p(z_{i,j} | x_i, l_j) \right]$$

Similarity to BA?

- Maximum a posteriori (MAP) inference:

$$x_{0:k}^*, L_k^* = \arg \max_{x_{0:k}, L_k} p(x_{0:k}, L_k | u_{0:k-1}, z_{0:k})$$

- For Gaussian distributions, involves solving a nonlinear least-squares problem:

$$x_{0:k}^*, L_k^* = \arg \min_{x_{0:k}, L_k} \left\{ \|x_0 - \hat{x}_0\|_{\Sigma_0}^2 + \sum_i \left[\|x_i - f(x_{i-1}, u_{i-1})\|_{\Sigma_w}^2 + \sum_{j \in \mathcal{M}_i} \|z_{i,j} - h(x_i, l_j)\|_{\Sigma_v}^2 \right] \right\}$$

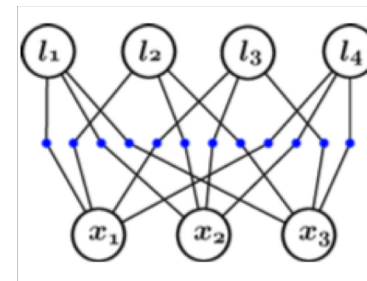
Smoothing and Mapping (SAM)

$$x_{0:k}^*, L_k^* = \arg \min_{x_{0:k}, L_k} \left\{ \|x_0 - \hat{x}_0\|_{\Sigma_0}^2 + \sum_i \left[\|x_i - f(x_{i-1}, u_{i-1})\|_{\Sigma_w}^2 + \sum_{j \in \mathcal{M}_i} \|z_{i,j} - h(x_i, l_j)\|_{\Sigma_v}^2 \right] \right\}$$

- Define $\Theta \doteq \{x_{0:k}, L_k\}$, linearize and collect terms

$$\Delta\Theta = \arg \min_{\Delta\Theta} \|\mathcal{A}\Delta\Theta - \check{b}\|^2$$

- Jacobian \mathcal{A} is a **big** & **sparse** matrix
- Example (camera-only, no motion model):
 - 3 cameras
 - 4 landmarks (3D points)



\mathcal{A}

	l_1	l_2	l_3	l_4	x_1	x_2	x_3
z_1	■				■		
z_2	■					■	
z_3	■						■
z_4		■				■	
z_5		■					■
z_6			■			■	
z_7			■				■
z_8				■		■	
z_9				■			■
z_{10}					■		
z_{11}						■	
z_{12}							■

Smoothing and Mapping (SAM)

$$x_{0:k}^*, L_k^* = \arg \min_{x_{0:k}, L_k} \left\{ \|x_0 - \hat{x}_0\|_{\Sigma_0}^2 + \sum_i \left[\|x_i - f(x_{i-1}, u_{i-1})\|_{\Sigma_w}^2 + \sum_{j \in \mathcal{M}_i} \|z_{i,j} - h(x_i, l_j)\|_{\Sigma_v}^2 \right] \right\}$$

- Define $\Theta \doteq \{x_{0:k}, L_k\}$, linearize and collect terms

$$\Delta\Theta = \arg \min_{\Delta\Theta} \|\mathcal{A}\Delta\Theta - \check{b}\|^2$$

- Jacobian \mathcal{A} is a **big** & **sparse** matrix
- How to recover MAP estimate efficiently, online?
- Sparsity-aware (incremental) optimization:
 - Solve via factorization (e.g. QR) and back-substitution
 - Update linearization point and repeat process until convergence

Solution via QR factorization

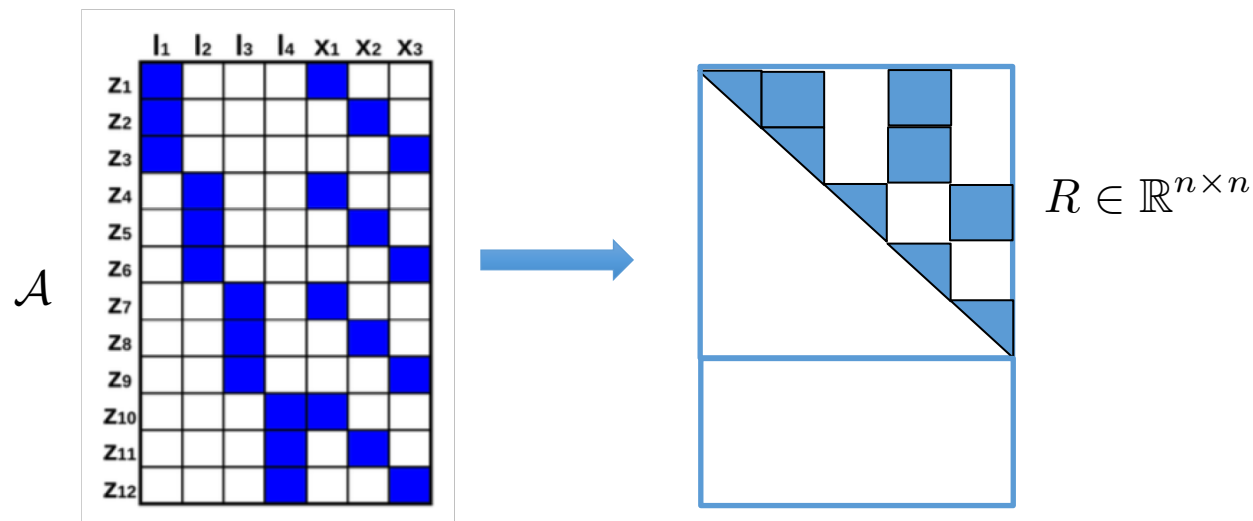
- Least squares problem:
$$\Delta\Theta = \arg \min_{\Delta\Theta} \left\| \mathcal{A}\Delta\Theta - \check{b} \right\|^2$$

- QR factorization:
$$Q^T \mathcal{A} = \begin{bmatrix} R \\ 0 \end{bmatrix} \quad Q^T \check{b} = \begin{bmatrix} d \\ e \end{bmatrix}$$

R - Upper triangular (sparse) matrix

Q - Orthogonal matrix

- It can be shown that:
$$\left\| \mathcal{A}\Delta\Theta - \check{b} \right\|_2^2 = \left\| R\Delta\Theta - d \right\|_2^2 + \underbrace{\left\| e \right\|_2^2}_{\text{Least-squares residual}}$$



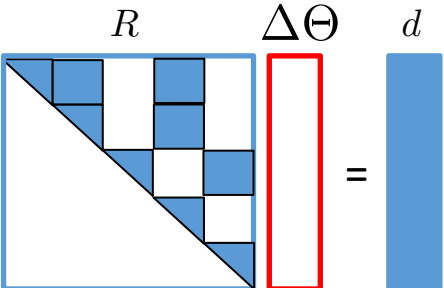
Solution via QR factorization

- Least squares problem:
$$\Delta\Theta = \arg \min_{\Delta\Theta} \|\mathcal{A}\Delta\Theta - \check{b}\|^2$$

- QR factorization:
$$Q^T \mathcal{A} = \begin{bmatrix} R \\ 0 \end{bmatrix} \quad Q^T \check{b} \doteq \begin{bmatrix} d \\ e \end{bmatrix}$$

R - Upper triangular (sparse) matrix
 Q - Orthogonal matrix

- Least squares (LS) solution $\Delta\Theta^*$ is obtained via **back-substitution**:

$$R\Delta\Theta = d$$


Graphical Model Perspective

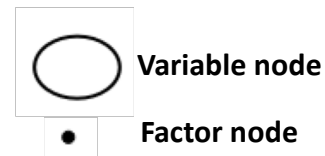
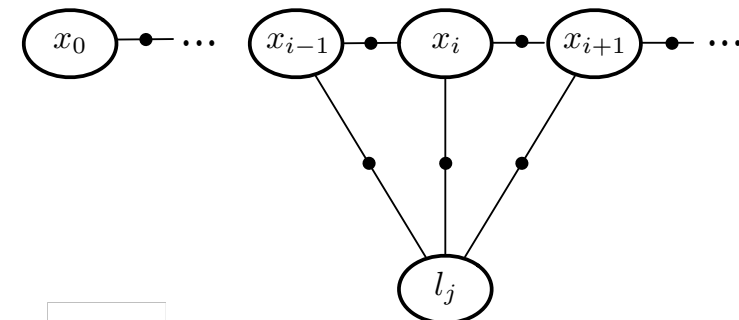
Factor Graph

- Bipartite undirected graph $G(\mathcal{F}, \Theta, \mathcal{E})$ with two node types
 - $\theta_j \in \Theta$: Variable nodes (correspond to states to be inferred)
 - $f_i \in \mathcal{F}$: Factor nodes (associated with process and measurement models)
 - $e_{ij} \in \mathcal{E}$: Edges always connect between variable and factor nodes
- Factor graph describes a factorization of the joint pdf in terms of process and measurement models

$$p(x_{0:k}, L | u_{0:k-1}, z_{0:k}) = \underbrace{\eta}_{\Theta} p(x_0) \prod_i \left[\underbrace{p(x_i | x_{i-1}, u_{i-1})}_{\text{factors } f_i(\Theta_i)} \prod_{j \in \mathcal{M}_i} \underbrace{p(z_{i,j} | x_i, l_j)}_{\text{factors } f_i(\Theta_i)} \right]$$



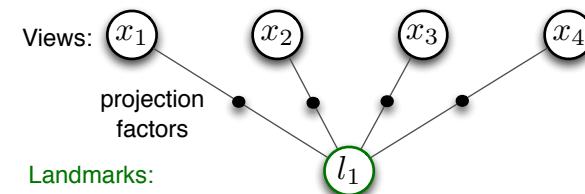
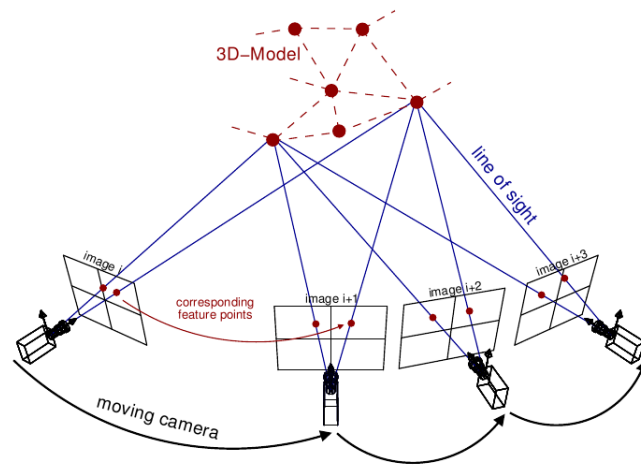
$$p(\Theta) \propto f(\Theta) \propto \prod f_i(\Theta_i)$$



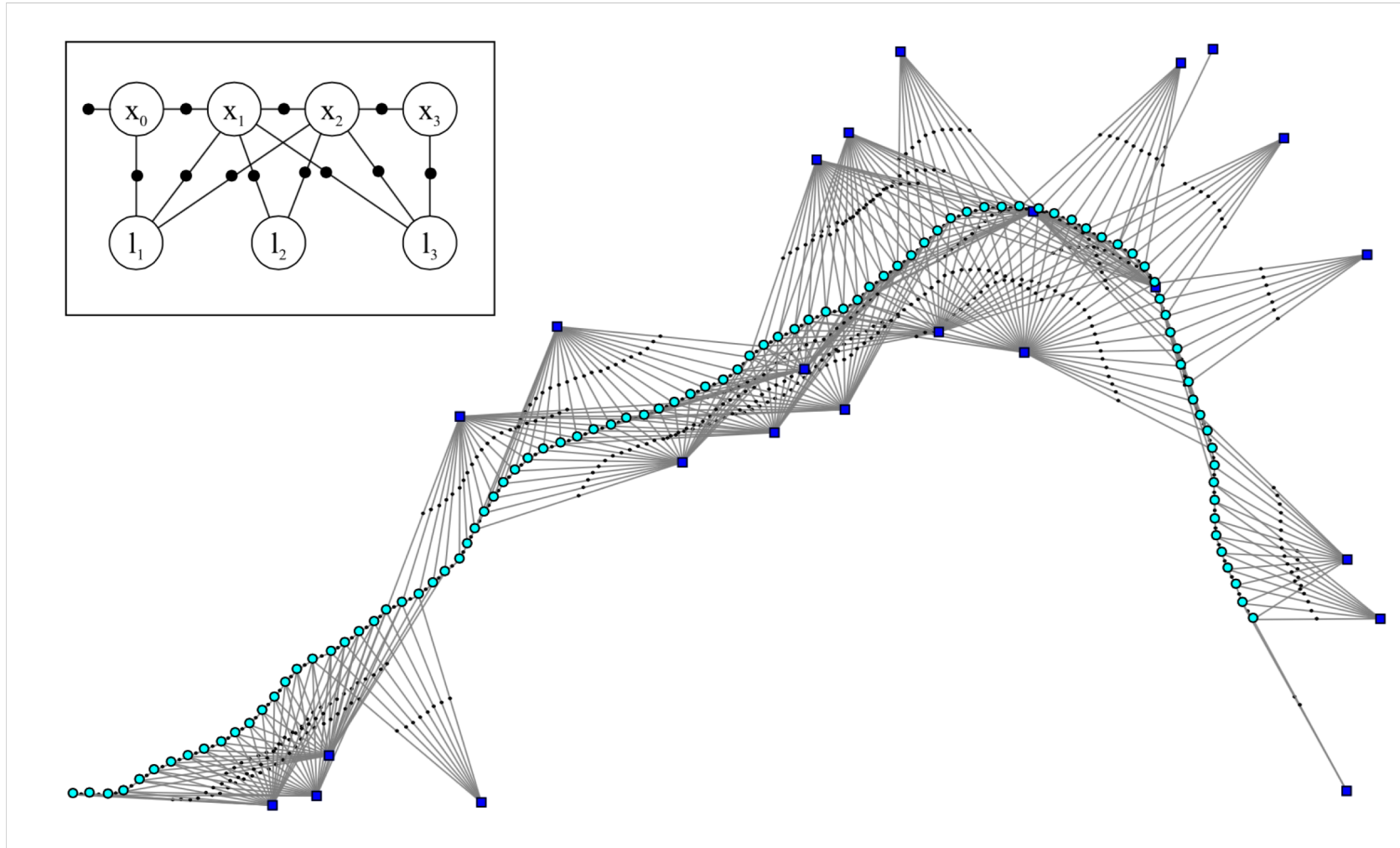
Factor Graph Representation for BA

$$p(X, L|Z) \propto \prod_i \prod_{j \in \mathcal{M}_i} p(z_{i,j} | x_i, l_j)$$

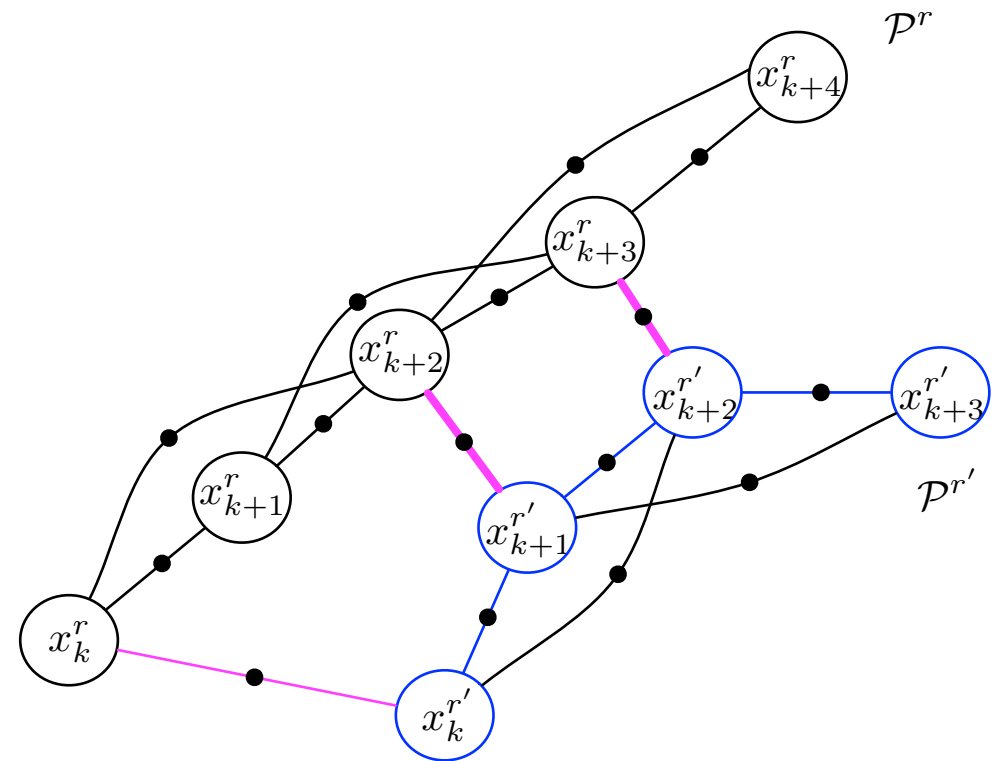
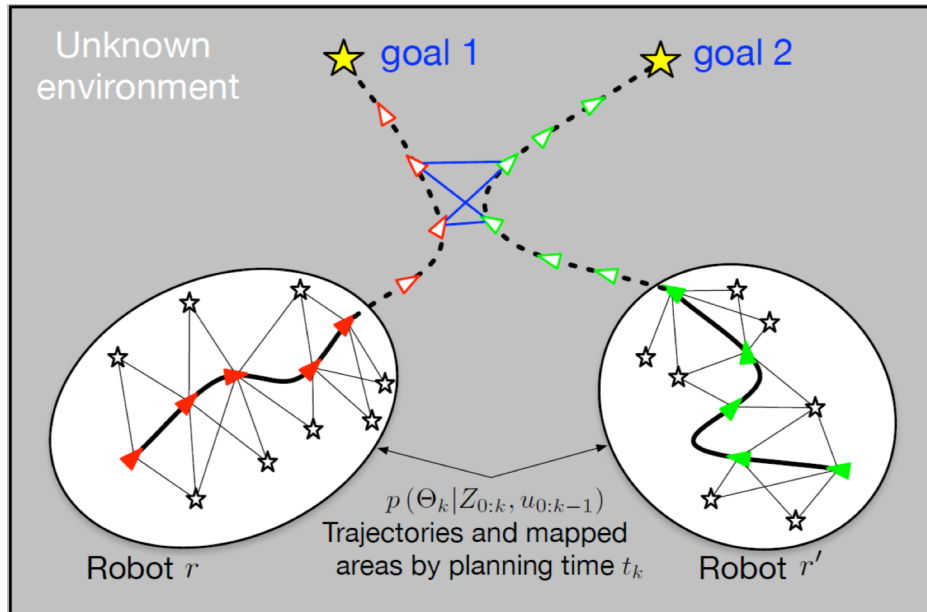
$$f_{proj}(x_i, l_j) \doteq \exp\left(-\frac{1}{2} \|z_{i,j} - \pi(x_i, l_j)\|_{\Sigma}^2\right)$$



Factor Graph – SLAM Problem



Factor Graph – Multi-Robot SLAM



Inference and Variable Elimination

- **Key insight:** Inference == Converting a factor graph to a Bayes net using the elimination alg.

- Factor graph:

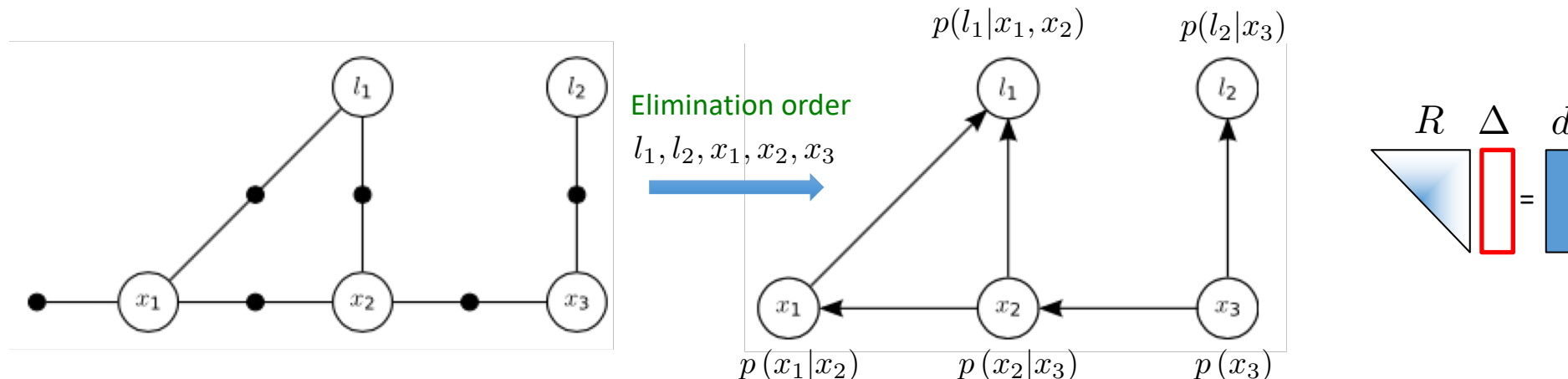
$$f(\Theta) = \prod_j f_j(\theta_j) = f(x_1) f(x_1, x_2) f(x_2, x_3) f(l_1, x_1) f(l_1, x_2) f(l_2, x_3)$$

- Represents the joint pdf (e.g.)

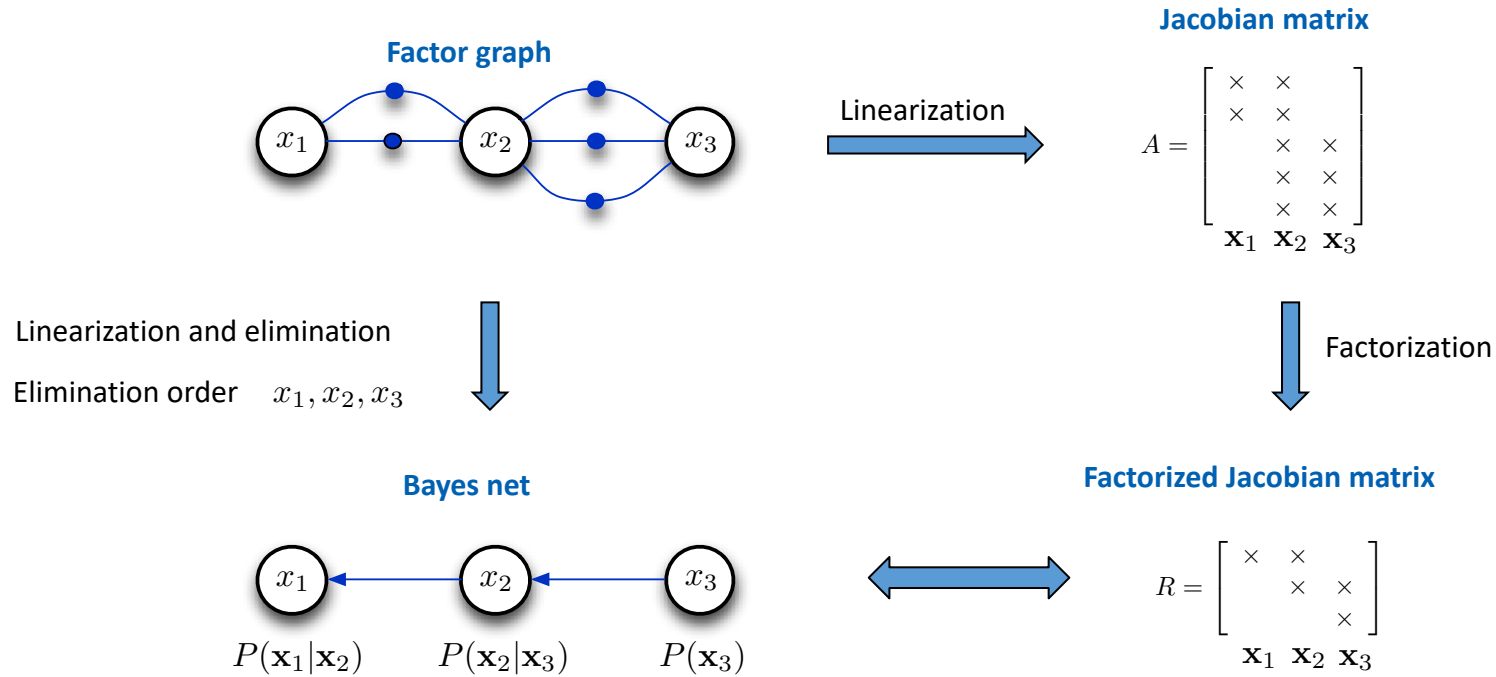
$$p(x_{1:3}, l_{1:2} | u_{1:2}, z_{1:3}) = \eta p(x_1) p(x_2 | x_1, u_1) p(x_3 | x_2, u_2) p(z_1 | x_1, l_1) p(z_2 | x_2, l_1) p(z_3 | x_3, l_2)$$

- Final result – Bayes net, corresponds to the factorization:

$$p(l_1 | x_1, x_2) p(l_2 | x_3) p(x_1 | x_2) p(x_2 | x_3) p(x_3)$$



Example



$$\Delta^* = \arg \min_{\Delta} (A\Delta - b)$$

$$A = QR \quad d \doteq Q^T b$$

$$\Delta^* = \arg \min_{\Delta} (R\Delta - d)$$

$$R \Delta = d$$

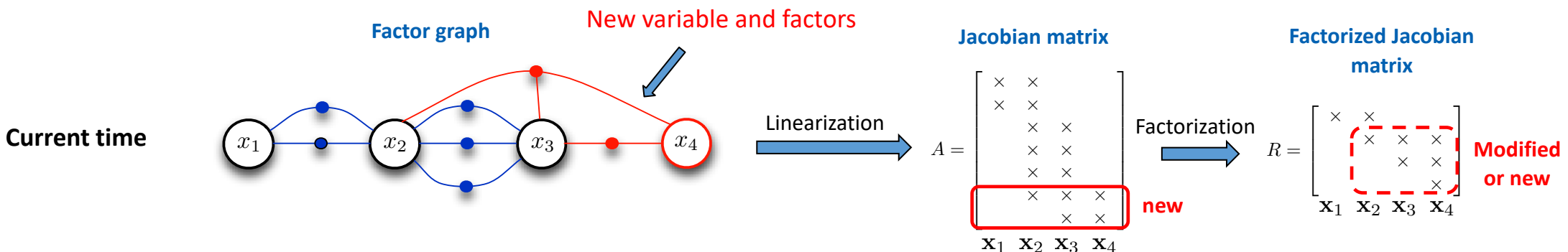
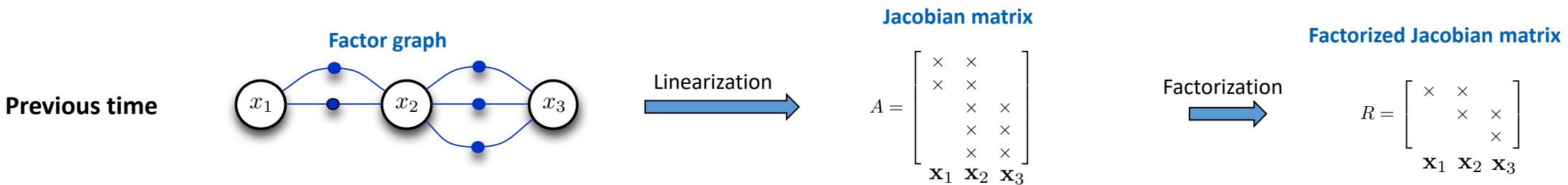
Thus Far: Smoothing and Mapping (SAM)

- Efficient, sparsity-aware nonlinear optimization
- Inference in graphical models
 - Joint pdf can be represented by a **factor graph**
 - Factorization (calculating the R matrix) is equivalent to variable elimination, represented by a **Bayes net**
- Still, **batch** algorithm:
 - Each time, solves the entire NLS problem
 - Online performance?

Incremental Smoothing & Mapping (iSAM)

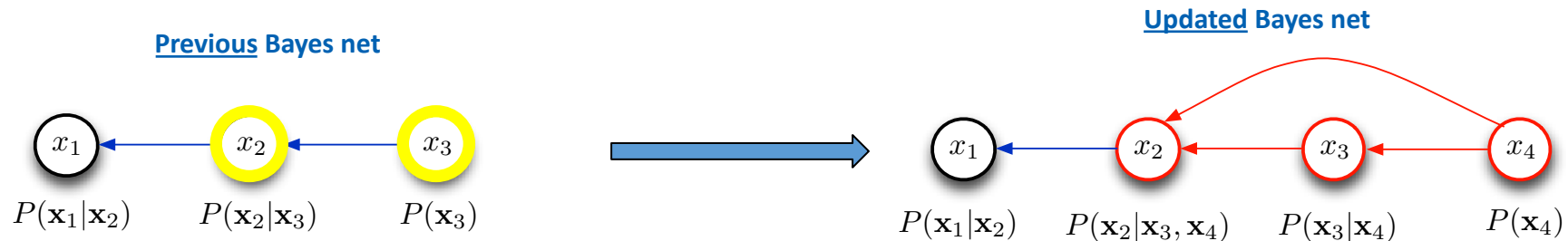
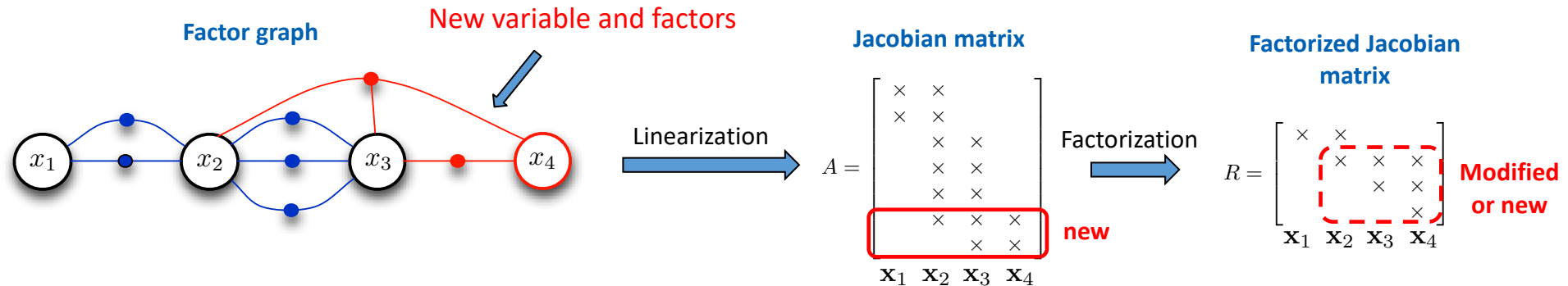
Each time new data is obtained (e.g. measurement):

- Previous approach (SAM): performs factorization from scratch
- **Is this really required?**
 - Typically, only a **very small** subset of entries is updated
 - Key idea – identify and update only these entries, i.e. **update** factorization (e.g. QR update)



Incremental Smoothing & Mapping (iSAM)

- Adding new variable (camera pose or navigation state) and factors



- What should be re-calculated?
- Nodes in all paths that lead from the last-eliminated node to nodes involved in new factors

$$\Delta^* = \arg \min_{\Delta} (A\Delta - b)$$

$$A = QR \quad d \doteq Q^T b$$

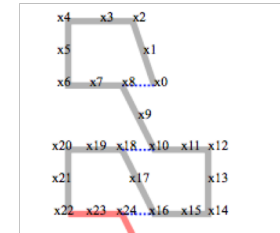
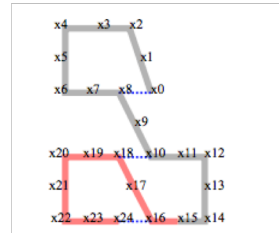
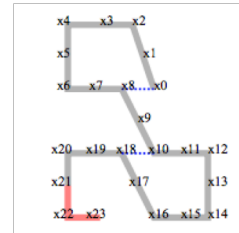
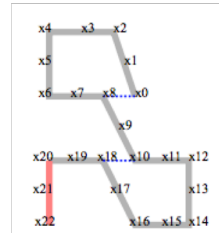
$$\Delta^* = \arg \min_{\Delta} (R\Delta - d)$$

$$R \quad \Delta \quad d$$

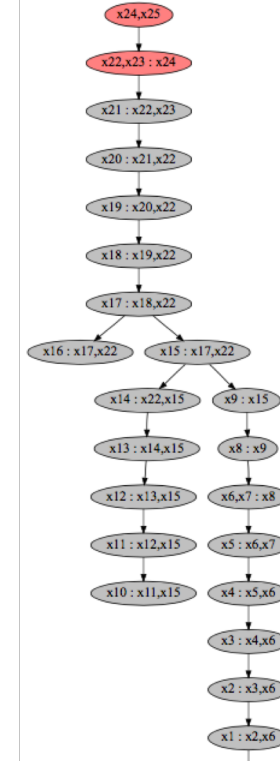
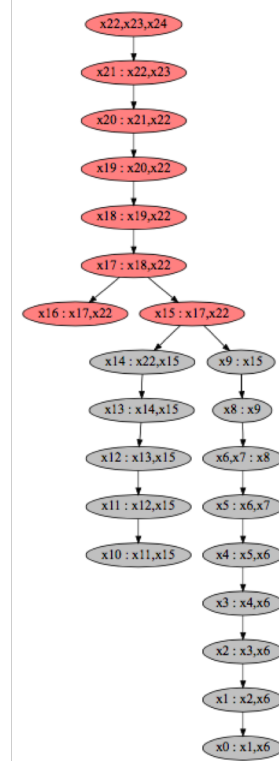
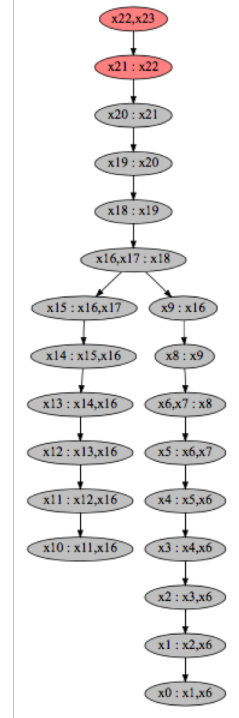
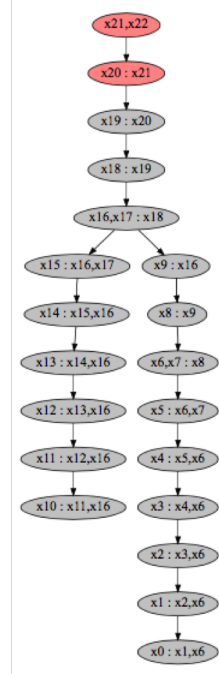
Incremental Smoothing & Mapping (iSAM)

time →

Trajectory:



Bayes tree:
(calculated from Bayes net)



Only red parts
are re-calculated

Example: Incremental Light Bundle Adjustment

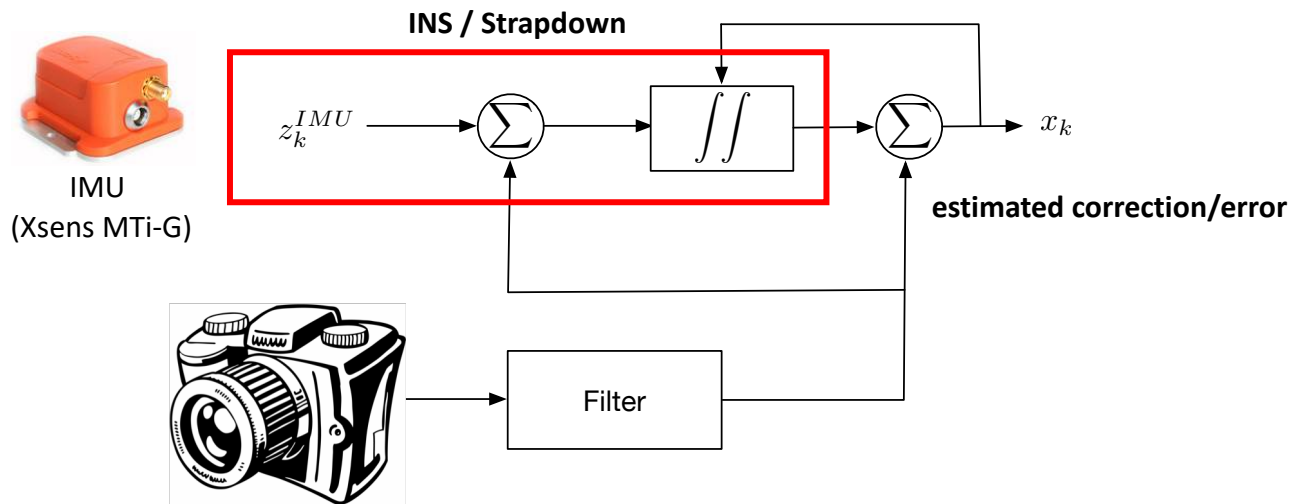


Outline

- Introduction
- Camera projective geometry
- Bundle Adjustment
- Incremental Smoothing and Mapping (iSAM) algorithms
- Visual-inertial SLAM and IMU pre-integration concept

Vision-Aided Navigation (VAN)

- Common approach:
 - Integrate IMU measurements **outside** the filter (in real time)
 - Use external sensors to correct solution
 - Does not support re-linearization (of past IMU measurements)



Vision-Aided Navigation (VAN)

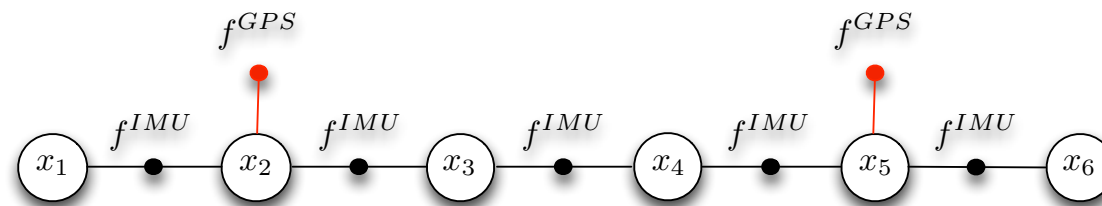
- Common approach:
 - Integrate IMU measurements **outside** the filter (in real time)
 - Use external sensors to correct solution
 - Does not support re-linearization (of past IMU measurements)
- An alternative – **visual-inertial bundle adjustment**
 - Concept: one big optimization
 - Better accuracy, improved estimation consistency
 - Real time performance?
 - Incremental factorization
 - IMU Pre-Integration

Information Fusion via Incremental Smoothing (iSAM algorithms)

Indelman et al. RAS 2013

Information Fusion via iSAM2.0

- The same concept applies also to multiple sensors, possibly operating at different rates
- The joint pdf and the factor graph should be accordingly
- Information fusion from different sensors:
 - Sensors introduce appropriate factors to the graph
 - Calculate MAP estimate via incremental smoothing (iSAM)
- Conceptually, factor graph representing IMU and GPS measurements:



IMU Only – Inertial Navigation

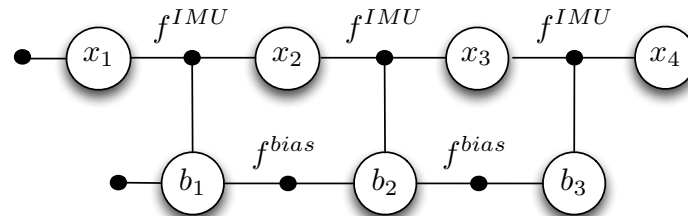
- Joint pdf
 - Formulated in terms of discrete inertial navigation equations (numerical integration)
 - Considers a basic model for IMU calibration parameters (for simplicity)



$$p(X_k, B_k | Z_k) \propto \prod_i^k p(x_{k+1} | x_k, b_k, z_k^{IMU}) p(b_{k+1} | b_k)$$

$\doteq f^{IMU}$ $\doteq f^{bias}$

Factor graph



IMU Only – Inertial Navigation

- Joint pdf
 - Formulated in terms of discrete inertial navigation equations (numerical integration)
 - Considers a basic model for IMU calibration parameters (for simplicity)

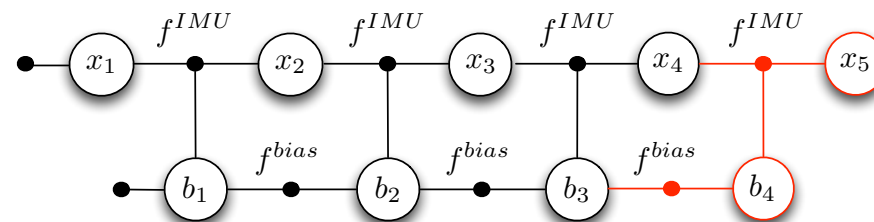


$$p(X_k, B_k | Z_k) \propto \prod_i^k p(x_{k+1} | x_k, b_k, z_k^{IMU}) p(b_{k+1} | b_k)$$

$\doteq f^{IMU}$ $\doteq f^{bias}$

- Assume MAP estimate at time t_k has been calculated
- What involves recovering MAP estimate at time t_{k+1} ?

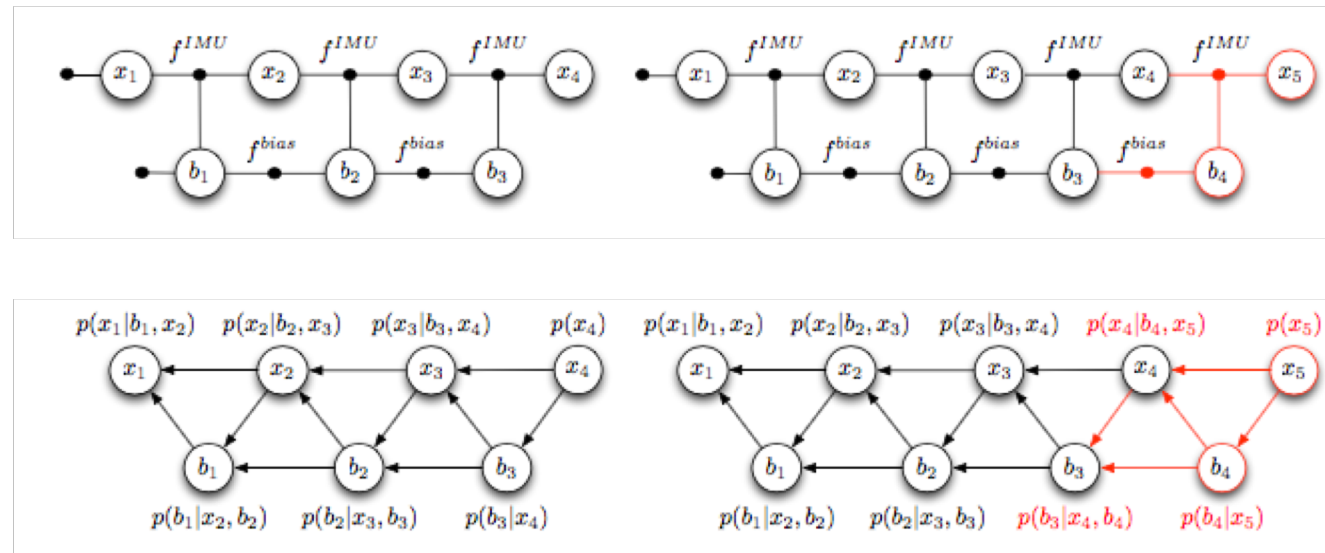
Factor graph



IMU Only – Inertial Navigation

- Corresponding factor graph and Bayes net:

Only 2 last variables should be re-eliminated



... but, what happens in presence of additional sensors?

Visual-Inertial Bundle Adjustment (SLAM, VAN)

- IMU + single camera:

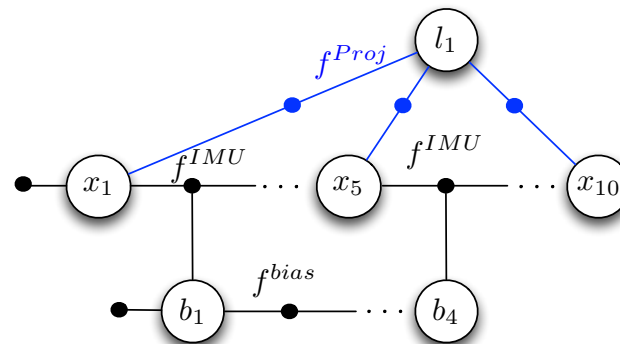
$$p(X_k, B_k, L_k | Z_k) \propto \prod_i^k p(x_{k+1} | x_k, b_k, z_k^{IMU}) \underset{\doteq f^{IMU}}{=} p(b_{k+1} | b_k) \underset{\doteq f^{bias}}{=} \prod_{j \in \mathcal{M}_i} p(z_{i,j} | x_i, l_j) \underset{\doteq f^{proj}}{=}$$

$$\Delta^* = \arg \min_{\Delta} (A\Delta - b)$$

$$A = QR \quad d \doteq Q^T b$$

$$\Delta^* = \arg \min_{\Delta} (R\Delta - d)$$

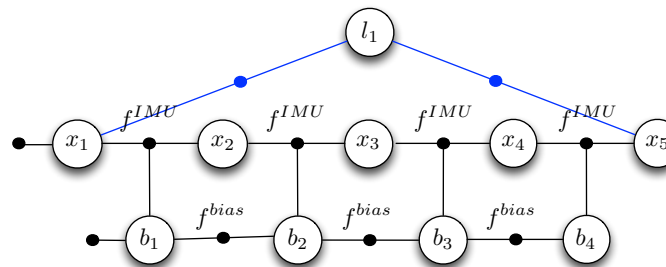
Factor graph



Incorporating High Rate Sensors

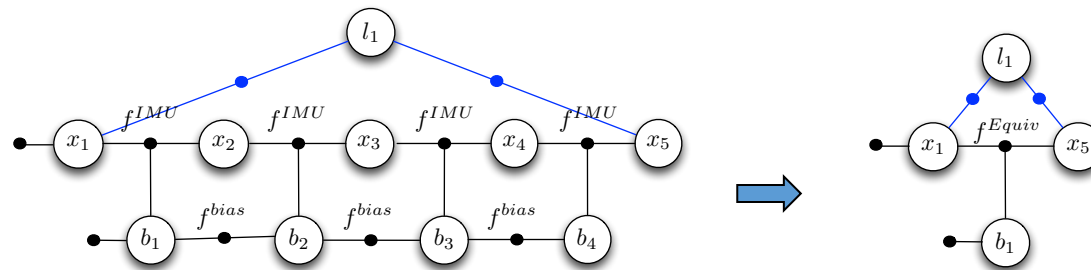
- **Challenges**

- Many variables to re-eliminate each time new measurements from other sensors (e.g. camera) come in
- Recover MAP estimate at IMU rate?
- How to avoid adding state variables to the optimization at IMU rate?



IMU Pre-Integration

- How to avoid adding state variables to the optimization at IMU rate?
 - **Pre-integrate** IMU observations in body frame of last keyframe [Lupton et al. 2012]
 - Navigation states can be added at camera rate



IMU Pre-Integration

- In navigation (global) frame:

$$v_{t_2}^n = v_{t_1}^n + \int_{t_1}^{t_2} (C_{bt}^n (f_t^b - bias_f^{obs}) + g^n) dt$$

Similar concept for position & orientation

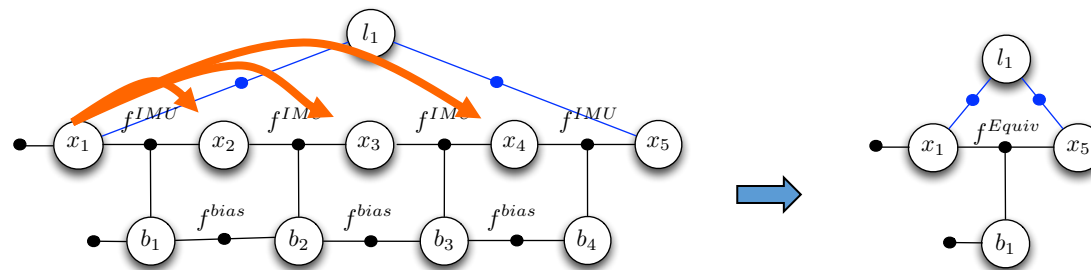
- Instead - **in body frame of the last pose** within optimization:

$$v_{t_2}^n = v_{t_1}^n + \int_{t_1}^{t_2} g^n dt + C_{bt_1}^n \int_{t_1}^{t_2} (C_{bt}^{bt_1} (f_t^b - bias_f^{obs})) dt$$



IMU Pre-Integration

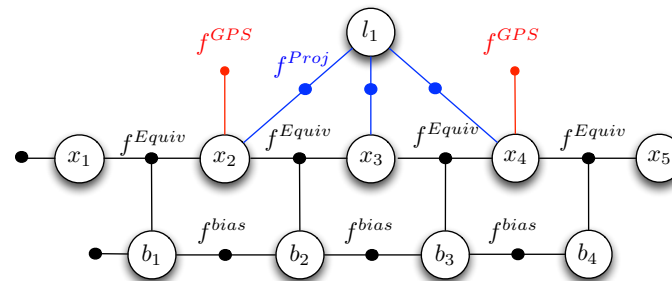
- How to avoid adding state variables to the optimization at IMU rate?
 - **Pre-integrate** IMU observations in body frame of last keyframe [Lupton et al. 2012]
 - Navigation states can be added at camera rate
- Real time performance – predict solution using pre-integrated IMU information and current MAP estimate
 - Without adding new variables to optimization
- Expressing in relative frame (i.e. body frame) allows to **avoid** re-playing all observations when re-linearizing!



Available as open source!
(part of gtsam)

Information Fusion via iSAM2.0

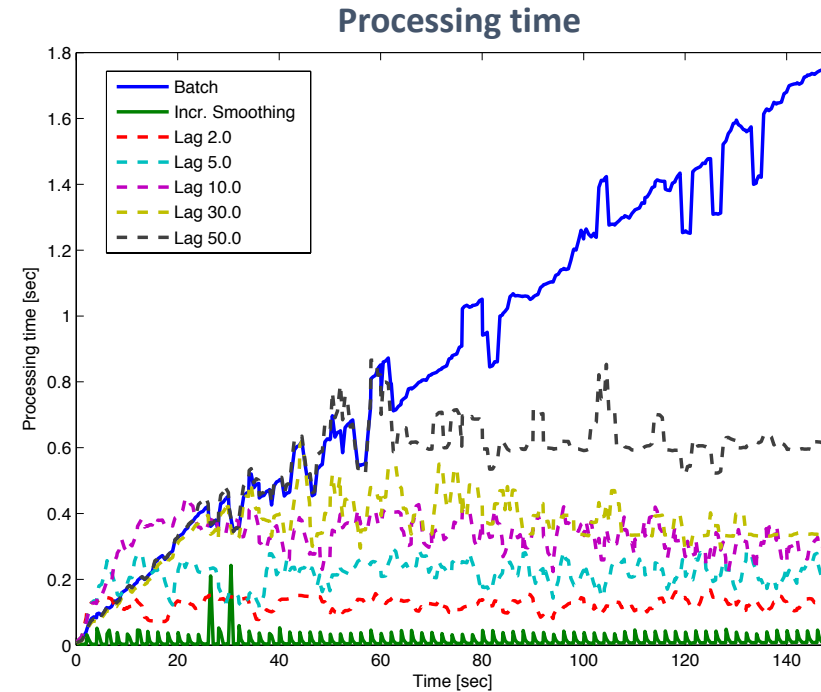
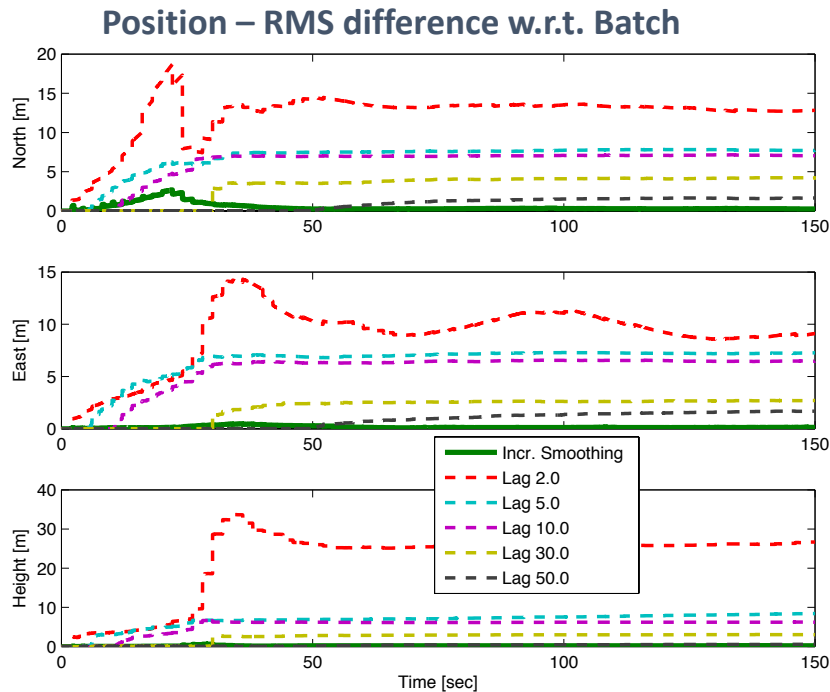
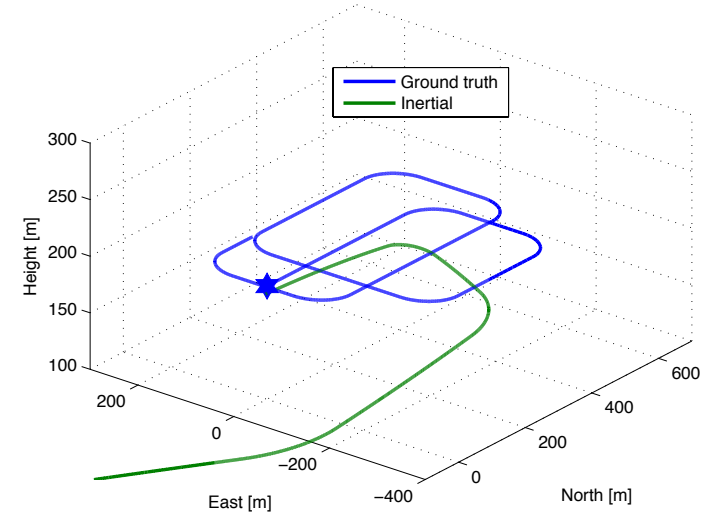
- Any other sensors can be similarly incorporated
For example:
 - IMU
 - GPS
 - Camera + explicit estimation of 3D points



Information Fusion via iSAM2.0

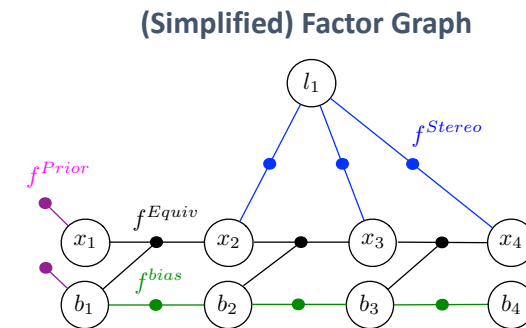
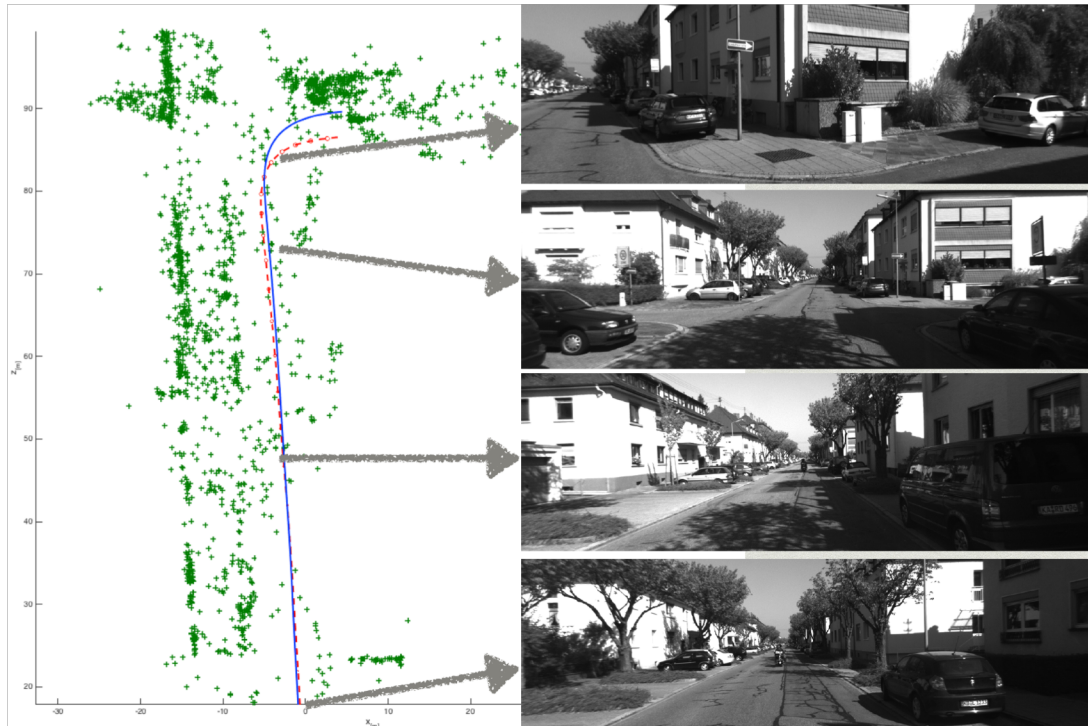
- Aerial scenario (simulation, Monte-Carlo study)
 - IMU, Monocular camera, Magnetometer
 - Short-track features only
 - Initial navigation errors

$$X_k^*, B_k^*, L_k^* = \arg \min_{X_k, B_k, L_k} J(X_k, B_k, L_k)$$



Information Fusion via iSAM2.0

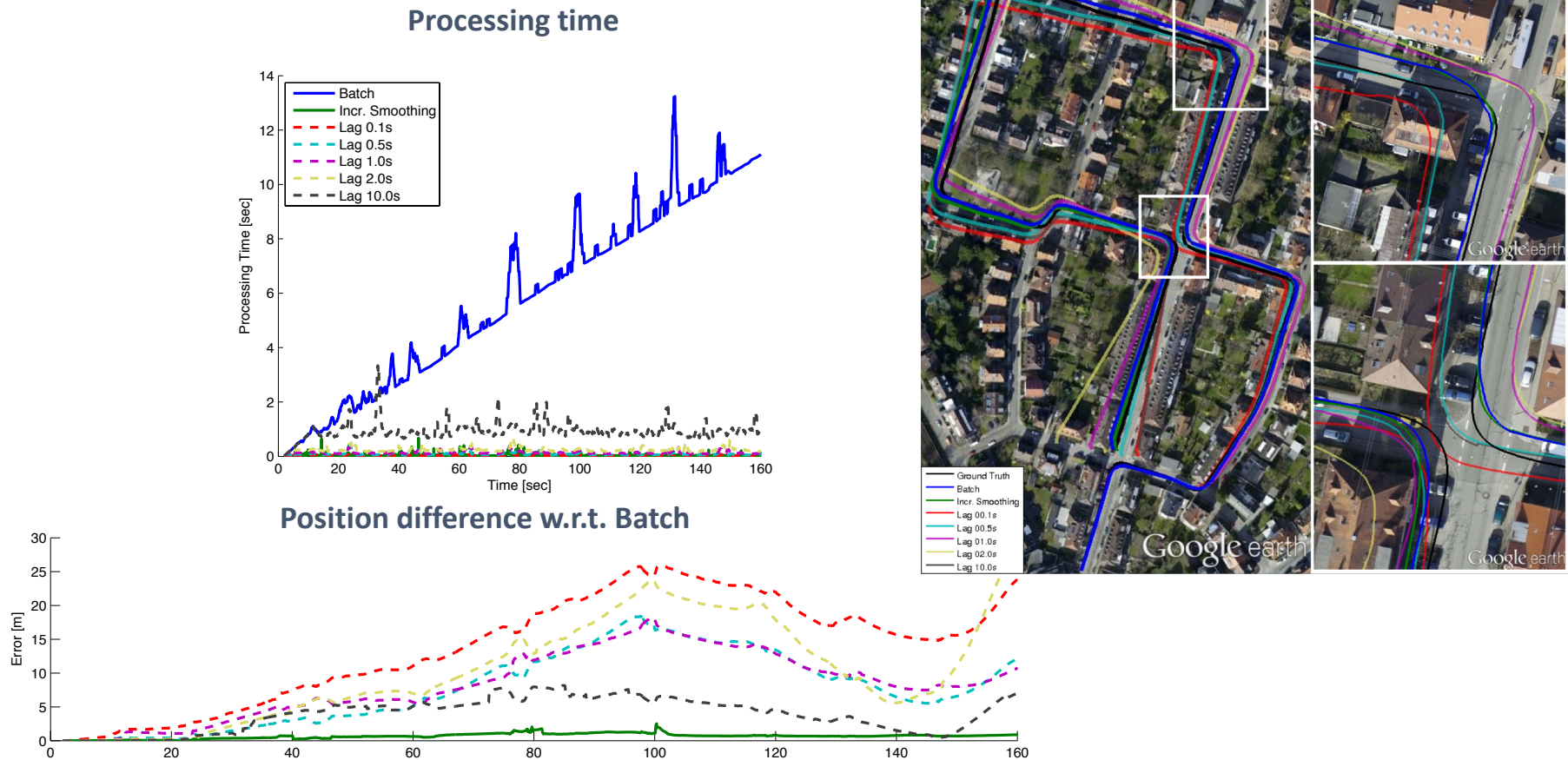
- KITTI Vision Benchmark
 - IMU
 - Stereo Camera (no loop closures)



Information Fusion via iSAM2.0

- KITTI Vision Benchmark
 - IMU
 - Stereo Camera (no loop closures)

$$X_k^*, B_k^*, L_k^* = \arg \max_{X_k, B_k, L_k} p(X_k, B_k, L_k | Z_k)$$



Images from Indelman13ras: "Information Fusion in Navigation Systems via Factor Graph Based Incremental Smoothing"

Summary

SLAM and VAN Overview:

- Front-end & Back-end
- Projection operator, re-projection error
- Bundle adjustment
- Smoothing and Mapping (SAM)
- Incremental SAM (iSAM)
- Information fusion with iSAM
- IMU Pre-Integration