

# Practical Course: Vision-based Navigation

## Winter Semester 2019

### Lecture 4. SfM

Vladyslav Usenko, Nikolaus Demmel,  
Prof. Dr. Daniel Cremers



# What We Will Cover Today

- Introduction to Visual SLAM
- Formulation of the SLAM Problem
- Full SLAM Posterior
- Bundle Adjustment (BA)
- Structure of the SLAM/BA Problem

# What is Visual SLAM?

- Visual simultaneous localization and mapping (VSLAM)...
  - Tracks the pose of the camera in a map, and simultaneously
  - Estimates the parameters of the environment map (f.e. reconstruct the 3D positions of interest points in a common coordinate frame)
- Loop-closure: Revisiting a place allows for drift compensation
  - How to detect a loop closure?

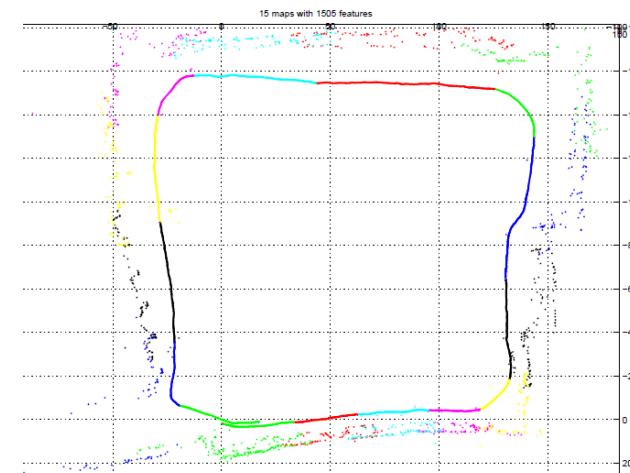
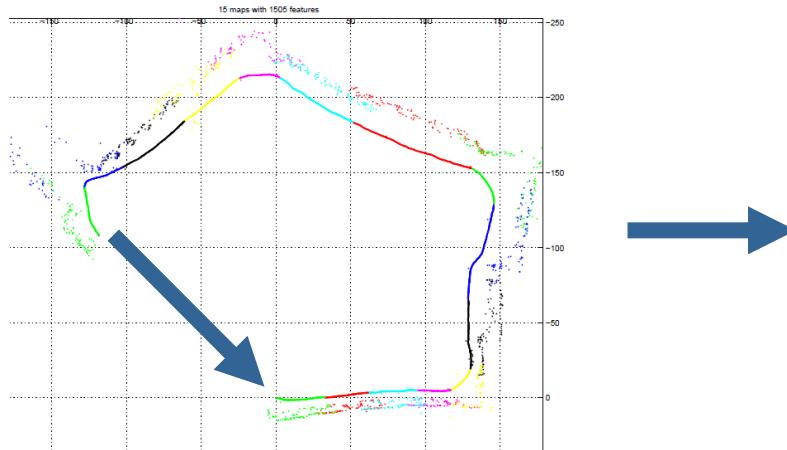
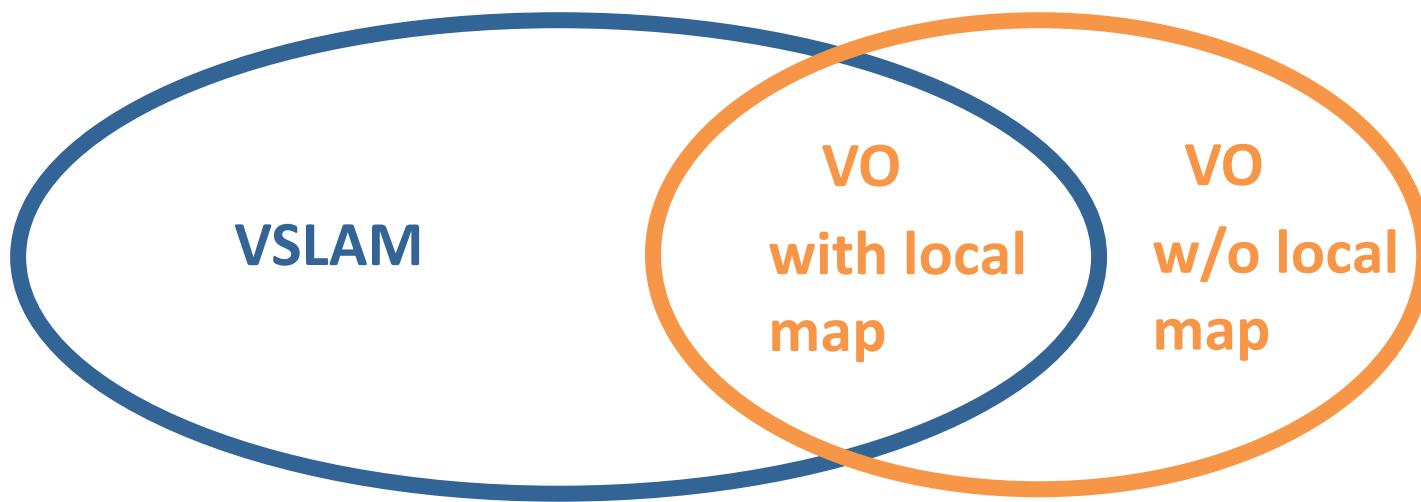


Image credit: Clemente et al., RSS 2007

# What is Visual SLAM?

- Visual simultaneous localization and mapping (VSLAM)...
  - Tracks the **pose** of the camera in a map, and simultaneously
  - Estimates the parameters of the **environment map** (f.e. reconstruct the 3D positions of interest points in a common coordinate frame)
- **Loop-closure:** Revisiting a place allows for drift compensation
  - How to detect a loop closure?
- **Global vs. local optimization methods**
  - Global: bundle adjustment, pose-graph optimization, etc.
  - Local: incremental tracking-and-mapping approaches, visual odometry with local maps. Often designed for real-time.
  - **Hybrids:** Real-time local SLAM + global optimization in a slower parallel process (f.e. LSD-SLAM)

# VO vs. VSLAM



# Structure from Motion

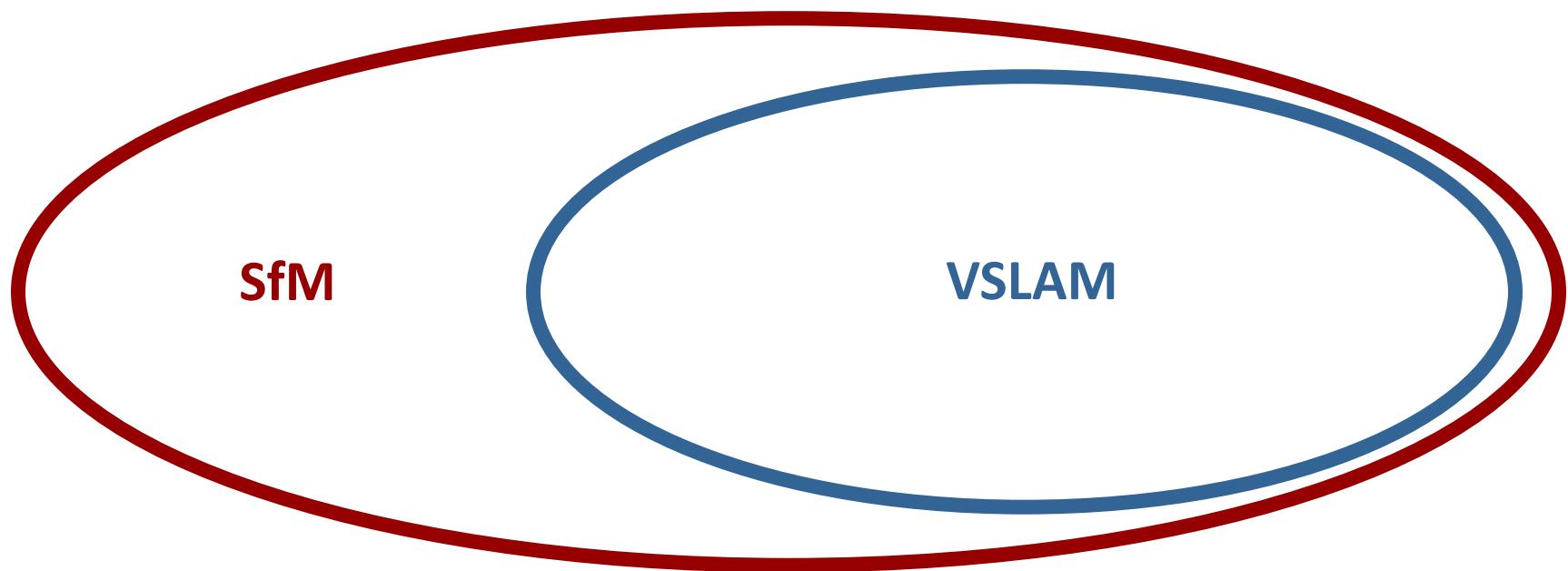
- Structure from Motion (SfM) denotes the joint estimation of
  - Structure, i.e. 3D reconstruction, and
  - Motion, i.e. 6-DoF camera poses,  
from a collection (i.e. unordered set) of images
- Typical approach: keypoint matching and bundle adjustment

# Structure from Motion



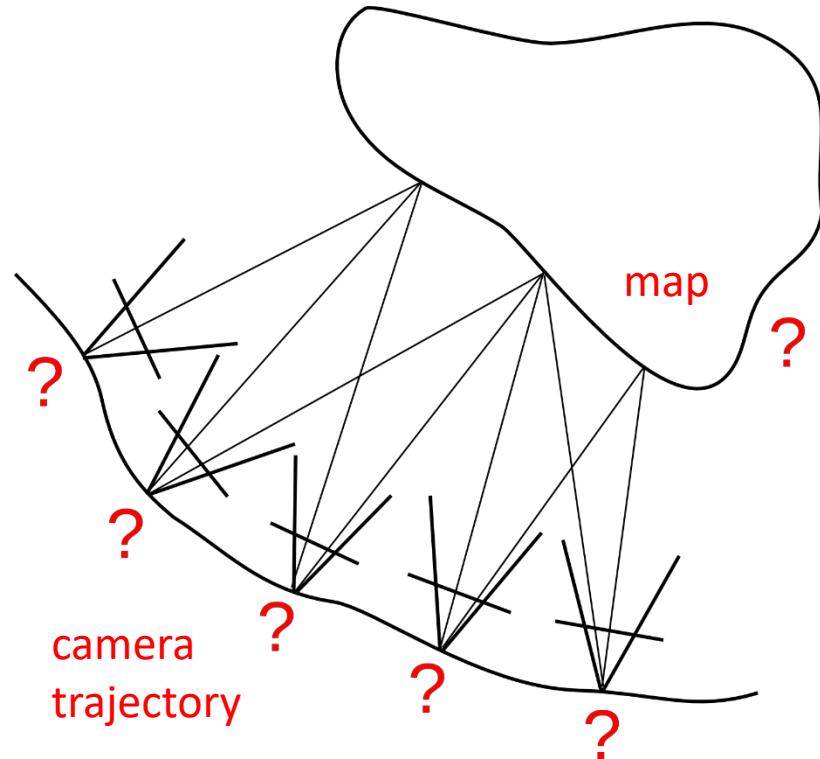
Agarwal et al., Building Rome in a Day, ICCV 2009, „Dubrovnik“ image set

# VSLAM vs. SfM



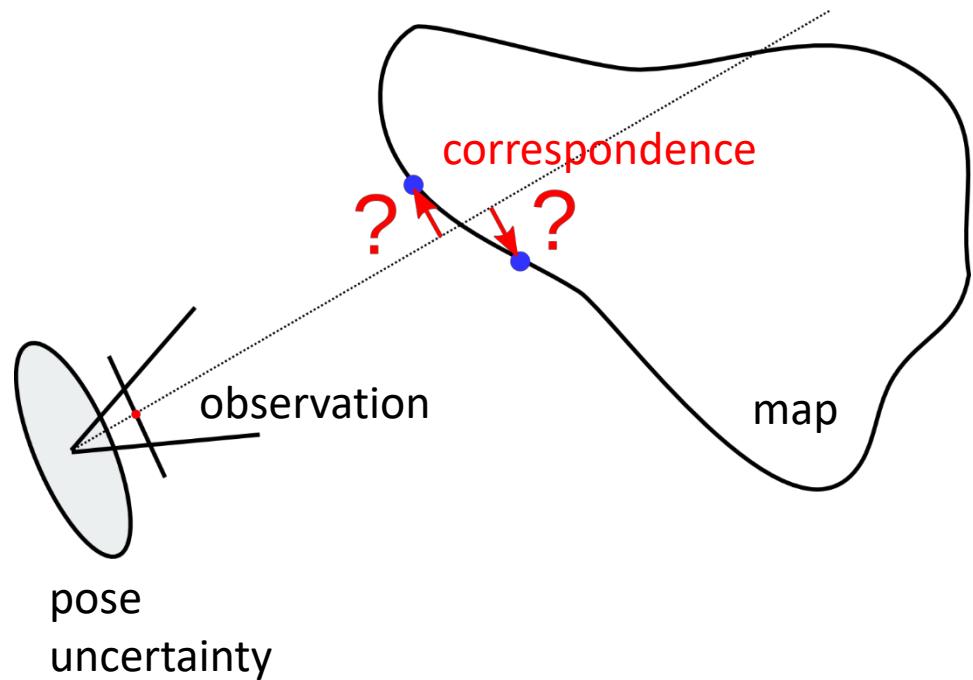
# Why is SLAM difficult?

- Chicken-or-egg problem
  - Camera trajectory and map are unknown and need to be estimated from observations
  - Accurate localization requires an accurate map
  - Accurate mapping requires accurate localization



# Why is SLAM difficult?

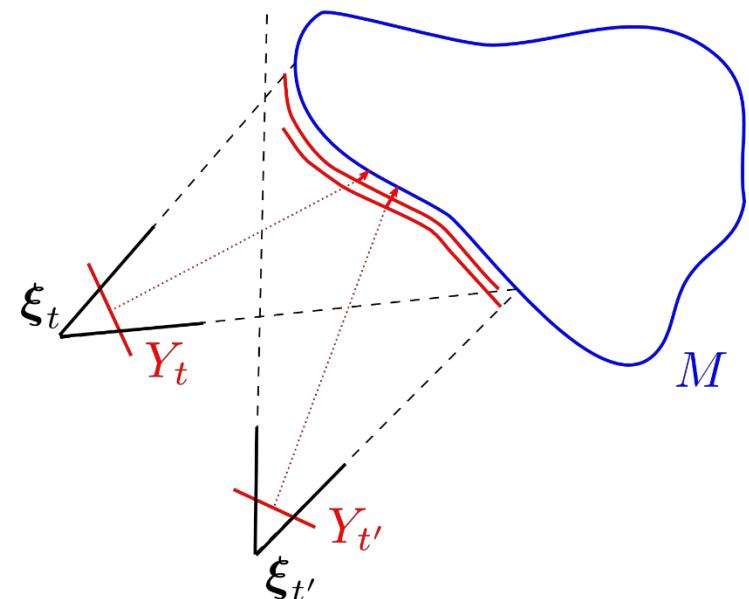
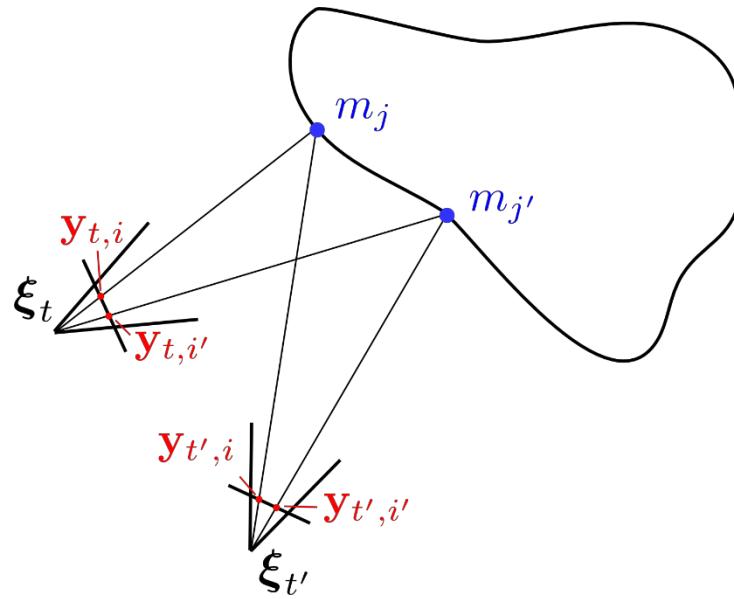
- Correspondences between observations and the map are unknown
- Wrong correspondences can lead to divergence of trajectory/map estimates
- Important to model uncertainties of observations and estimates in a probabilistic formulation of the SLAM problem



# Definition of Visual SLAM

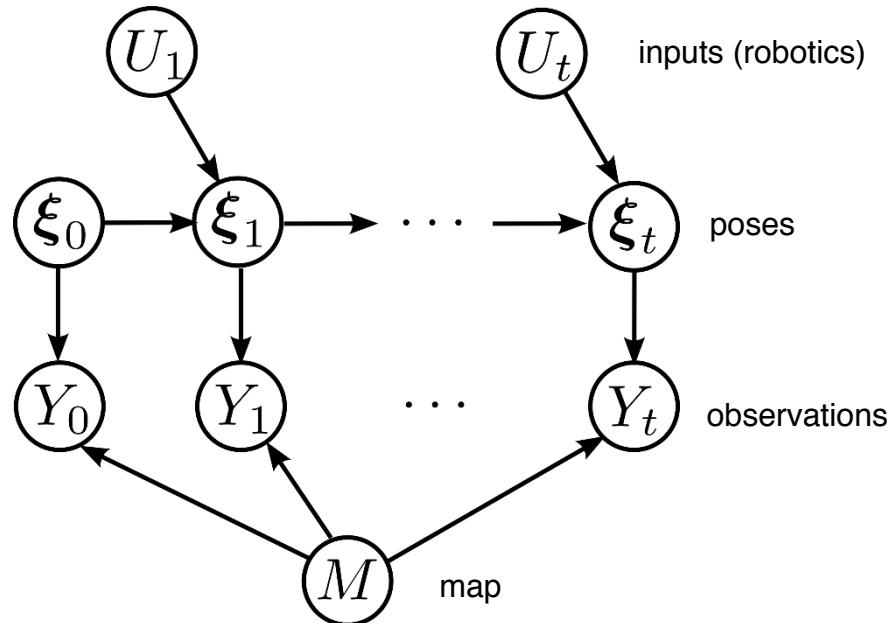
- Visual SLAM is the process of **simultaneously** estimating the **egomotion** of an object and the **environment map** using **only inputs from visual sensors** on the object and control inputs
- **Inputs:** images at discrete time steps  $t$ ,
  - **Monocular case:** Set of images  $I_{0:t} = \{I_0, \dots, I_t\}$
  - **Stereo case:** Left/right images  $I_{0:t}^l = \{I_0^l, \dots, I_t^l\}$      $I_{0:t}^r = \{I_0^r, \dots, I_t^r\}$
  - **RGB-D case:** Color/depth images  $I_{0:t} = \{I_0, \dots, I_t\}$      $Z_{0:t} = \{Z_0, \dots, Z_t\}$
- Robotics: **control inputs**  $U_{1:t}$
- **Output:**
  - **Camera pose** estimates  $\mathbf{T}_t \in \text{SE}(3)$  in world reference frame.  
For convenience, we also write  $\xi_t = \underline{\xi}(\mathbf{T}_t)$
  - **Environment map**  $M$

# Map Observations in Visual SLAM



- With  $Y_t$  we denote observations of the environment map in image  $I_t$ , f.e.
  - Indirect point-based method:  $Y_t = \{\mathbf{y}_{t,1}, \dots, \mathbf{y}_{t,N}\}$  (2D or 3D image points)
  - Direct RGB-D method:  $Y_t = \{I_t, Z_t\}$  (all image pixels)
  - ...
- Involves data association to map elements  $M = \{m_1, \dots, m_S\}$ 
  - We denote correspondences by  $c_{t,i} = j, 1 \leq i \leq N, 1 \leq j \leq S$   
map  $i^{\text{th}}$  observation at time  $t$  to map element  $j$

# Probabilistic Formulation of Visual SLAM

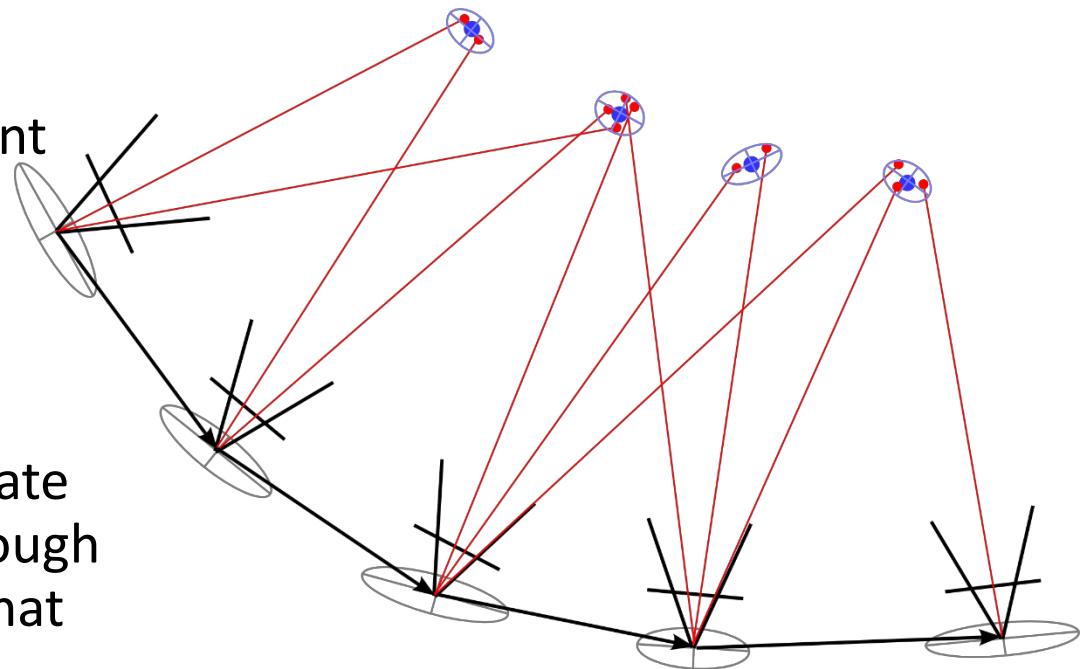


- SLAM **posterior** probability:  $p(\xi_{0:t}, M \mid Y_{0:t}, U_{1:t})$
- Observation **likelihood**:  $p(Y_t \mid \xi_t, M)$  with initial estimate of the map
- State-transition probability:  $p(\xi_t \mid \xi_{t-1}, U_t)$

# SLAM Graph Optimization

- Joint optimization for poses and map elements from image observations of map elements

- Common map element observations induce constraints between the poses



- Map elements correlate with each others through the common poses that observe them
- No temporal sequence: Bundle Adjustment

# Probabilistic Formulation

- SLAM posterior:  $p(\xi_{0:t}, M \mid Y_{0:t}, U_{1:t}, c_{0:t})$

- Observation likelihood:

$$p(Y_t \mid \xi_t, M, c_t) = p(Y_t \mid \xi_t, m_{ct})$$

$$p(Y_t \mid \xi_t, m_{ct}) = \prod_i p(y_{t,i} \mid \xi_t, m_{ct,i})$$

- State-transition probability:

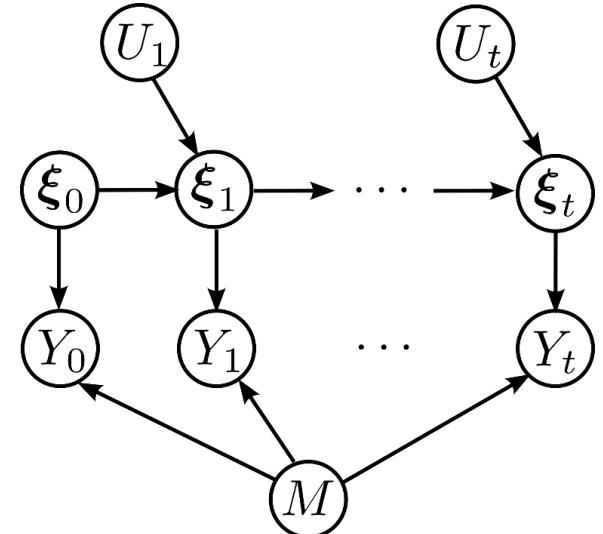
$$p(\xi_t \mid \xi_{t-1}, U_t)$$

- SLAM posterior can be **factorized**:

$$p(\xi_{0:t}, M \mid Y_{0:t}, U_{1:t}, c_{0:t}) = \eta p(Y_t \mid \xi_t, m_{ct}) p(\xi_{0:t}, M \mid Y_{0:t-1}, U_{1:t}, c_{0:t-1})$$

$$= \eta p(Y_t \mid \xi_t, m_{ct}) p(\xi_t \mid \xi_{t-1}, U_t) p(\xi_{0:t-1}, M \mid Y_{0:t-1}, U_{1:t-1})$$

$$= \eta' p(\xi_0) p(M) \prod_t p(Y_t \mid \xi_t, m_{ct}) p(\xi_t \mid \xi_{t-1}, U_t)$$



recursively

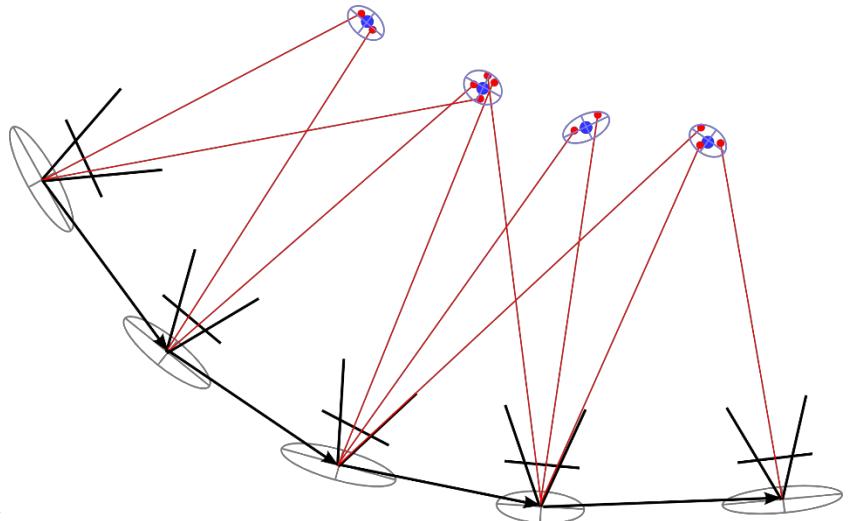
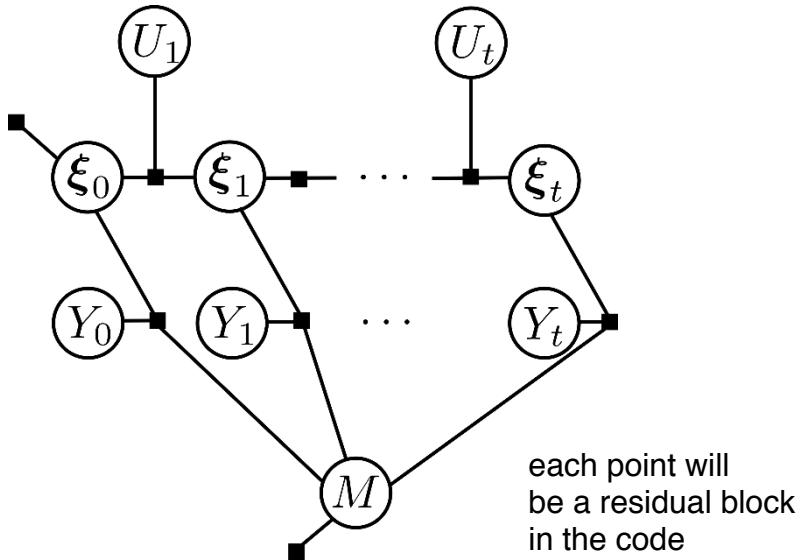
# Factor Graph

- Factor graph representation of the full SLAM posterior

$$p(\xi_{0:t}, M | Y_{0:t}, U_{1:t}, c_{0:t})$$

$$= \eta \ p(\xi_0) p(M) \prod_t p(Y_t | \xi_t, m_{ct}) p(\xi_t | \xi_{t-1}, U_t)$$

]

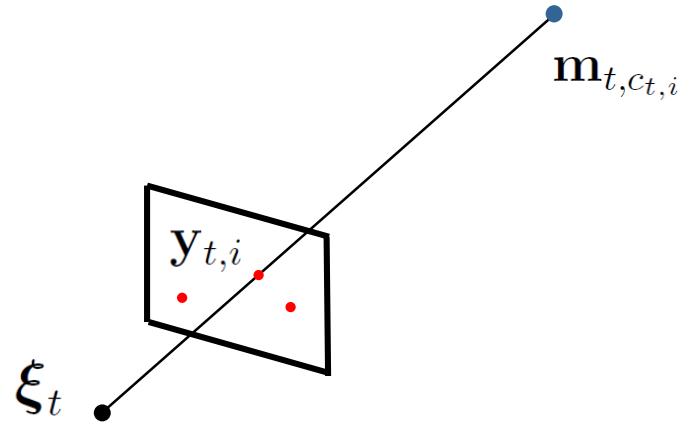


# Explicit Model

- $N_t$  noisy 2D point observation  
of 3D landmarks in each image,  
known data association

$$\mathbf{y}_{t,i} = \underline{h(\xi_t, \mathbf{m}_{t,c_{t,i}})} + \delta_t = \underline{\pi(\mathbf{T}(\xi_t)^{-1} \bar{\mathbf{m}}_{t,c_{t,i}})} + \delta_{t,i}$$

$$\delta_{t,i} \sim \mathcal{N}(0, \Sigma_{\mathbf{y}_{t,i}})$$



- No control inputs
- Gaussian prior on pose  $\underline{\xi_0} \sim \mathcal{N}(\xi^0, \Sigma_{0,\xi})$
- Uniform prior on landmarks

# Full SLAM Optimization as Energy Minimization

- Optimize negative log posterior probability (MAP estimation)

$$E(\xi_{0:t}, M) = \frac{1}{2} (\xi_0 \ominus \xi^0)^\top \Sigma_{0,\xi}^{-1} (\xi_0 \ominus \xi^0)$$

$$+ \frac{1}{2} \sum_{\tau=0}^t \sum_{i=1}^{N_\tau} (\mathbf{y}_{\tau,i} - h(\xi_\tau, \mathbf{m}_{c_{\tau,i}}))^\top \Sigma_{\mathbf{y}_{\tau,i}}^{-1} (\mathbf{y}_{\tau,i} - h(\xi_\tau, \mathbf{m}_{c_{\tau,i}}))$$

- Non-linear least squares!! We know how to optimize this..
- Remark: noisy state transitions based on control inputs add further residuals between subsequent poses

# Full SLAM Optimization as Energy Minimization

- Let's define the residuals on the full state vector

$$\mathbf{r}^0(\mathbf{x}) := \boldsymbol{\xi}_0 \ominus \boldsymbol{\xi}^0$$

$$\mathbf{r}_{t,i}^y(\mathbf{x}) := \mathbf{y}_{t,i} - h(\boldsymbol{\xi}_t, \mathbf{m}_{c_{t,i}})$$

$$\mathbf{x} := \begin{pmatrix} \boldsymbol{\xi}_0 \\ \vdots \\ \boldsymbol{\xi}_t \\ \mathbf{m}_1 \\ \vdots \\ \mathbf{m}_S \end{pmatrix}$$

t poses  
s landmarks

- Stack the residuals in a vector-valued function and collect the residual covariances on the diagonal blocks of a square matrix

$$\mathbf{r}(\mathbf{x}) := \begin{pmatrix} \mathbf{r}^0(\mathbf{x}) \\ \mathbf{r}_{0,1}^y(\mathbf{x}) \\ \vdots \\ \mathbf{r}_{t,N_t}^y(\mathbf{x}) \end{pmatrix} \quad \mathbf{W} := \begin{pmatrix} \Sigma_{0,\boldsymbol{\xi}}^{-1} & 0 & \cdots & 0 \\ 0 & \Sigma_{\mathbf{y}_{0,1}}^{-1} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \Sigma_{\mathbf{y}_{t,N_t}}^{-1} \end{pmatrix}$$

- Rewrite error function as  $\underline{E(\mathbf{x})} = \frac{1}{2}\mathbf{r}(\mathbf{x})^\top \mathbf{W} \mathbf{r}(\mathbf{x})$

# Recap: Gauss-Newton Method

- Idea: Approximate Newton's method to minimize  $E(x)$ 
  - Approximate  $E(x)$  through linearization of residuals

$$\begin{aligned}\tilde{E}(x) &= \frac{1}{2} \tilde{r}(x)^\top W \tilde{r}(x) \\ &= \frac{1}{2} (\mathbf{r}(x_k) + \mathbf{J}_k(x - x_k))^\top W (\mathbf{r}(x_k) + \mathbf{J}_k(x - x_k)) \quad \mathbf{J}_k := \nabla_x \mathbf{r}(x)|_{x=x_k} \\ &= \frac{1}{2} \mathbf{r}(x_k)^\top W \mathbf{r}(x_k) + \underbrace{\mathbf{r}(x_k)^\top W \mathbf{J}_k}_{=: b_k^\top} (x - x_k) + \frac{1}{2} (x - x_k)^\top \underbrace{\mathbf{J}_k^\top W \mathbf{J}_k}_{=: H_k} (x - x_k)\end{aligned}$$

- Find root of  $\nabla_x \tilde{E}(x) = b_k^\top + (x - x_k)^\top H_k$  using Newton's method, i.e.

$$\nabla_x \tilde{E}(x) = 0 \text{ iff } x = x_k - \underbrace{H_k^{-1} b_k}_{=:}$$

- Pros:
  - Faster convergence (approx. quadratic convergence rate)
- Cons:
  - Divergence if too far from local optimum ( $H$  not positive definite)
  - Solution quality depends on initial guess

# Structure of the Bundle Adjustment Problem

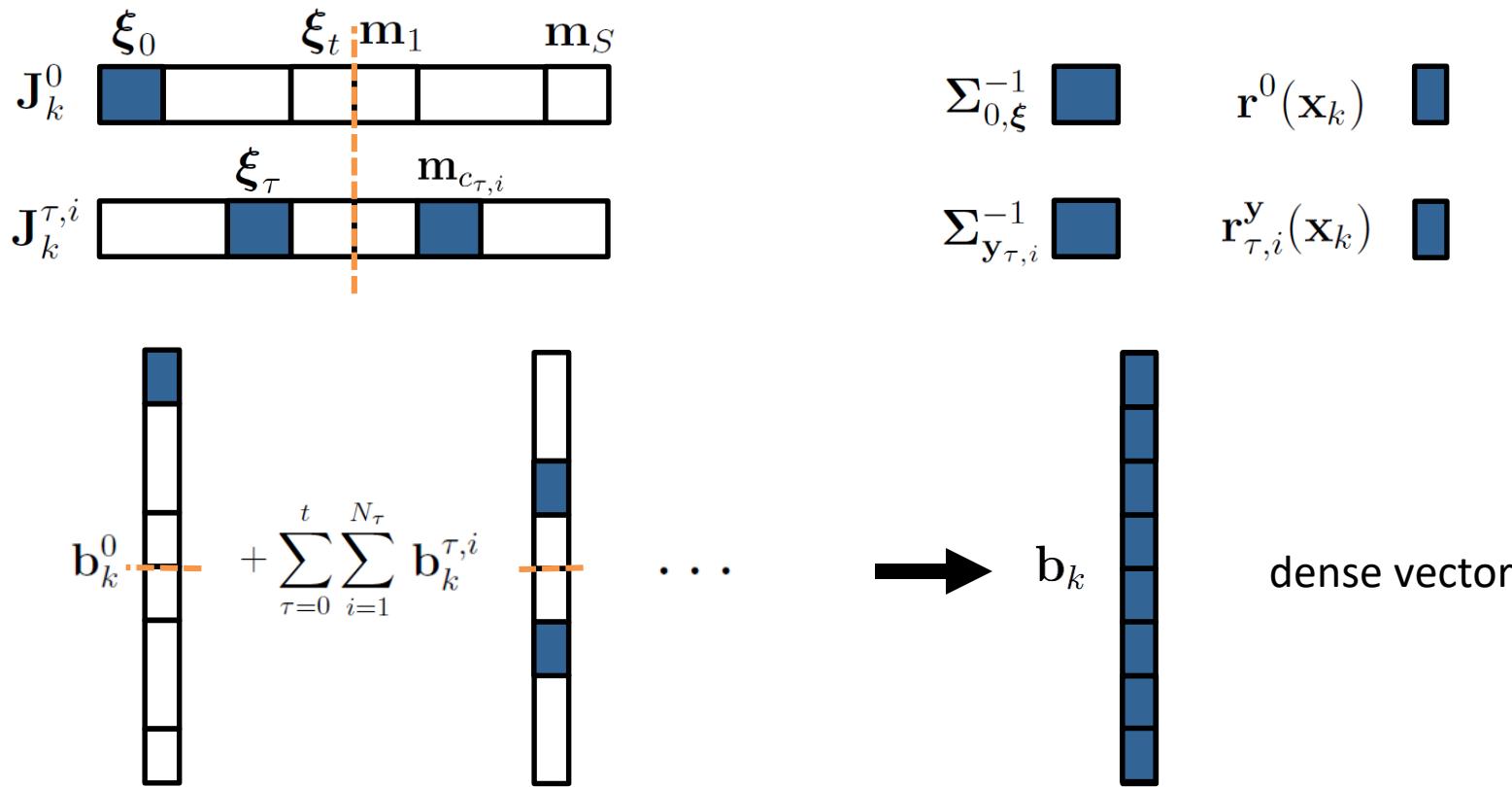
- $\mathbf{b}_k$  and  $\mathbf{H}_k$  sum terms from individual residuals:

$$\mathbf{b}_k = \mathbf{b}_k^0 + \sum_{\tau=0}^t \sum_{i=1}^{N_\tau} \mathbf{b}_k^{\tau,i} = (\mathbf{J}_k^0)^\top \Sigma_{0,\xi}^{-1} \mathbf{r}^0(\mathbf{x}_k) + \sum_{\tau=0}^t \sum_{i=1}^{N_\tau} (\mathbf{J}_k^{\tau,i})^\top \Sigma_{\mathbf{y}_{\tau,i}}^{-1} \mathbf{r}_{\tau,i}^{\mathbf{y}}(\mathbf{x}_k)$$

$$\mathbf{H}_k = \mathbf{H}_k^0 + \sum_{\tau=0}^t \sum_{i=1}^{N_\tau} \mathbf{H}_k^{\tau,i} = (\mathbf{J}_k^0)^\top \Sigma_{0,\xi}^{-1} (\mathbf{J}_k^0) + \sum_{\tau=0}^t \sum_{i=1}^{N_\tau} (\mathbf{J}_k^{\tau,i})^\top \Sigma_{\mathbf{y}_{\tau,i}}^{-1} (\mathbf{J}_k^{\tau,i})$$

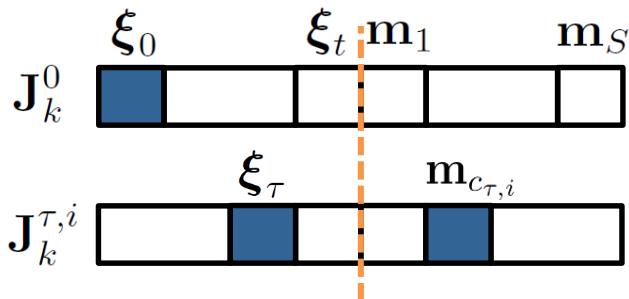
- What is the structure of these terms?

# Structure of the Bundle Adjustment Problem



$$\mathbf{b}_k = \mathbf{b}_k^0 + \sum_{\tau=0}^t \sum_{i=1}^{N_\tau} \mathbf{b}_k^{\tau,i} = (\mathbf{J}_k^0)^\top \Sigma_{0,\xi}^{-1} \mathbf{r}^0(\mathbf{x}_k) + \sum_{\tau=0}^t \sum_{i=1}^{N_\tau} (\mathbf{J}_k^{\tau,i})^\top \Sigma_{\mathbf{y}_{\tau,i}}^{-1} \mathbf{r}_{\tau,i}^{\mathbf{y}}(\mathbf{x}_k)$$

# Structure of the Bundle Adjustment Problem



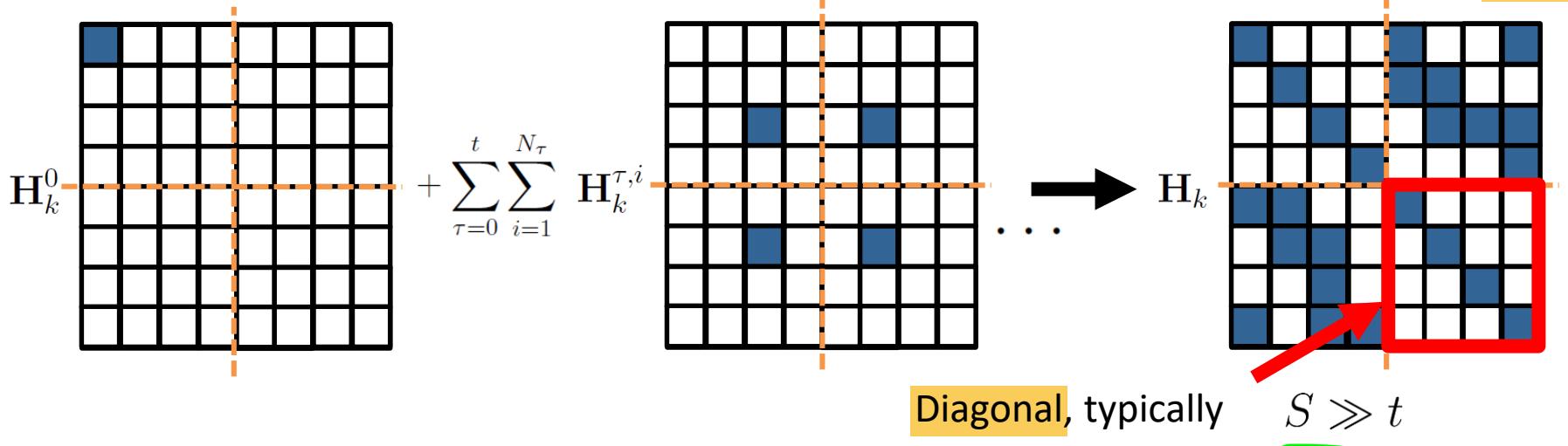
$$\Sigma_{0,\xi}^{-1} \quad \boxed{\text{blue square}}$$

$$\mathbf{r}^0(\mathbf{x}_k) \quad \boxed{\text{blue rectangle}}$$

$$\Sigma_{\mathbf{y}_{\tau,i}}^{-1} \quad \boxed{\text{blue square}}$$

$$\mathbf{r}_{\tau,i}^{\mathbf{y}}(\mathbf{x}_k) \quad \boxed{\text{blue rectangle}}$$

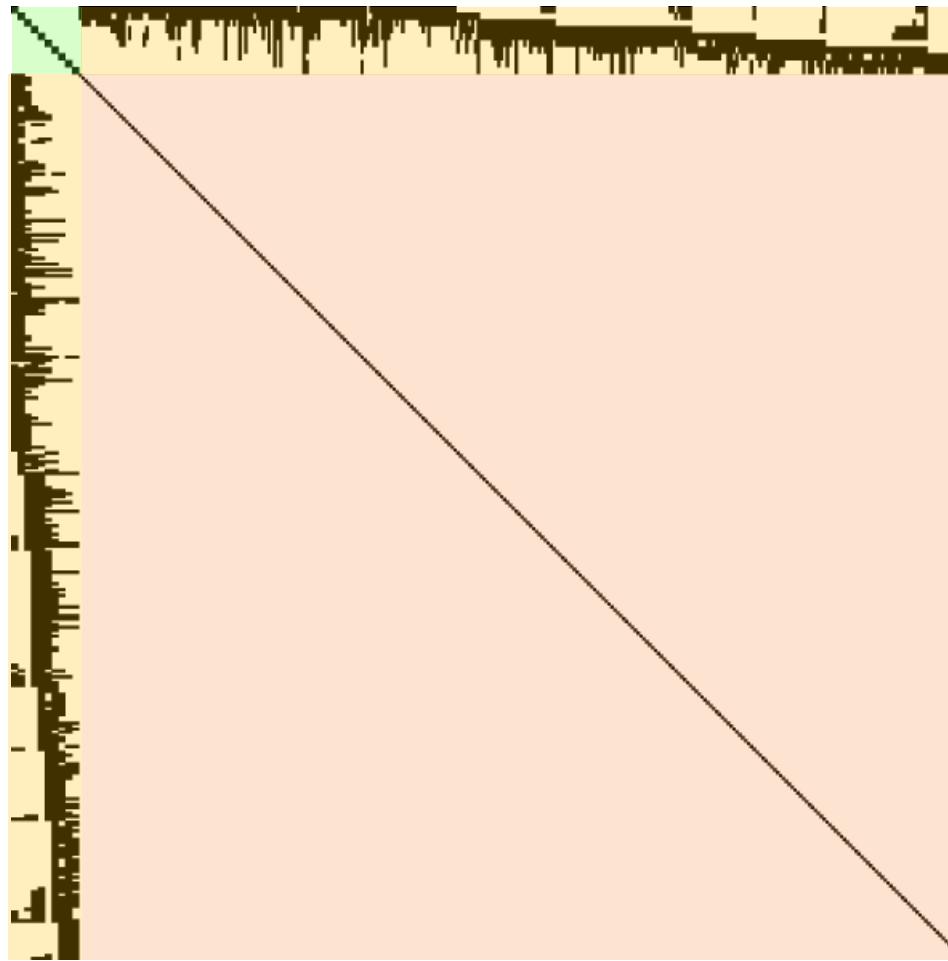
Sparse!



$$\mathbf{H}_k = \mathbf{H}_k^0 + \sum_{\tau=0}^t \sum_{i=1}^{N_\tau} \mathbf{H}_k^{\tau,i} = (\mathbf{J}_k^0)^\top \Sigma_{0,\xi}^{-1} (\mathbf{J}_k^0) + \sum_{\tau=0}^t \sum_{i=1}^{N_\tau} (\mathbf{J}_k^{\tau,i})^\top \Sigma_{\mathbf{y}_{\tau,i}}^{-1} (\mathbf{J}_k^{\tau,i})$$

# Example Hessian of a BA Problem

Pose dimensions  
(10 poses)



Landmark  
dimensions  
(982 landmarks)

# Exploiting the Sparse Structure

- Idea:

Apply the Schur complement to solve the system in a partitioned way

$$\underbrace{H_k \Delta x}_{} = -b_k \quad \longrightarrow \quad \begin{pmatrix} H_{\xi\xi} & H_{\xi m} \\ H_{m\xi} & H_{mm} \end{pmatrix} \begin{pmatrix} \Delta x_\xi \\ \Delta x_m \end{pmatrix} = - \begin{pmatrix} b_\xi \\ b_m \end{pmatrix}$$

$$\longrightarrow \underbrace{\Delta x_\xi}_{=} = - (H_{\xi\xi} - H_{\xi m} \underbrace{H_{mm}^{-1}}_{=} H_{m\xi})^{-1} (b_\xi - H_{\xi m} \underbrace{H_{mm}^{-1}}_{=} b_m)$$

$$\longrightarrow \Delta x_m = - \underbrace{H_{mm}^{-1}}_{=} (b_m + H_{m\xi} \underbrace{\Delta x_\xi}_{=})$$

more easier to invert diagonal matrices

- Is this any better?

# Exploiting the Sparse Structure

- What is the structure of the two sub-problems ?

$$\Delta \mathbf{x}_\xi = - \underbrace{(\mathbf{H}_{\xi\xi} - \mathbf{H}_{\xi m} \mathbf{H}_{mm}^{-1} \mathbf{H}_{m\xi})^{-1}}_{\text{Reduced pose Hessian}} \underbrace{(\mathbf{b}_\xi - \mathbf{H}_{\xi m} \mathbf{H}_{mm}^{-1} \mathbf{b}_m)}$$

- Poses:

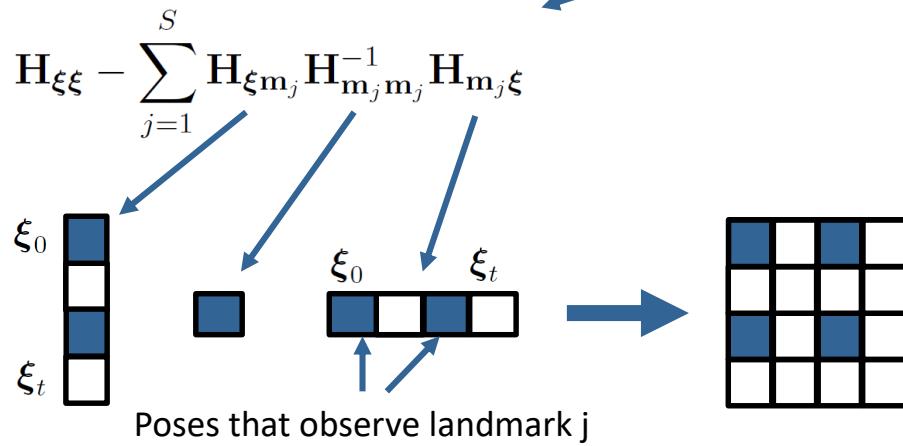
$$\mathbf{H}_{\xi\xi} - \mathbf{H}_{\xi m} \mathbf{H}_{mm}^{-1} \mathbf{H}_{m\xi} = \mathbf{H}_{\xi\xi} - \sum_{j=1}^S \mathbf{H}_{\xi m_j} \mathbf{H}_{m_j m_j}^{-1} \mathbf{H}_{m_j \xi}$$
$$\mathbf{b}_\xi - \mathbf{H}_{\xi m} \mathbf{H}_{mm}^{-1} \mathbf{b}_m = \mathbf{b}_\xi - \sum_{j=1}^S \mathbf{H}_{\xi m_j} \mathbf{H}_{m_j m_j}^{-1} \mathbf{b}_{m_j}$$

Reduced pose Hessian

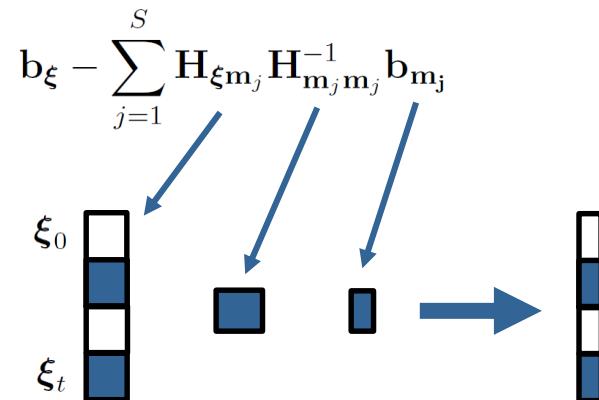
# Exploiting the Sparse Structure

- What is the structure of the two sub-problems ?

- Poses:  $\Delta \mathbf{x}_\xi = - \underbrace{(\mathbf{H}_{\xi\xi} - \mathbf{H}_{\xi m} \mathbf{H}_{mm}^{-1} \mathbf{H}_{m\xi})^{-1}}_{\text{Sparse Matrix}} (\mathbf{b}_\xi - \mathbf{H}_{\xi m} \mathbf{H}_{mm}^{-1} \mathbf{b}_m)$



$$\mathbf{H}_{\xi\xi} - \sum_{j=1}^S \mathbf{H}_{\xi m_j} \mathbf{H}_{m_j m_j}^{-1} \mathbf{H}_{m_j \xi} = \text{Sparse Matrix}$$



$$\mathbf{b}_\xi - \sum_{j=1}^S \mathbf{H}_{\xi m_j} \mathbf{H}_{m_j m_j}^{-1} \mathbf{b}_{m_j} = \text{Sparse Vector}$$

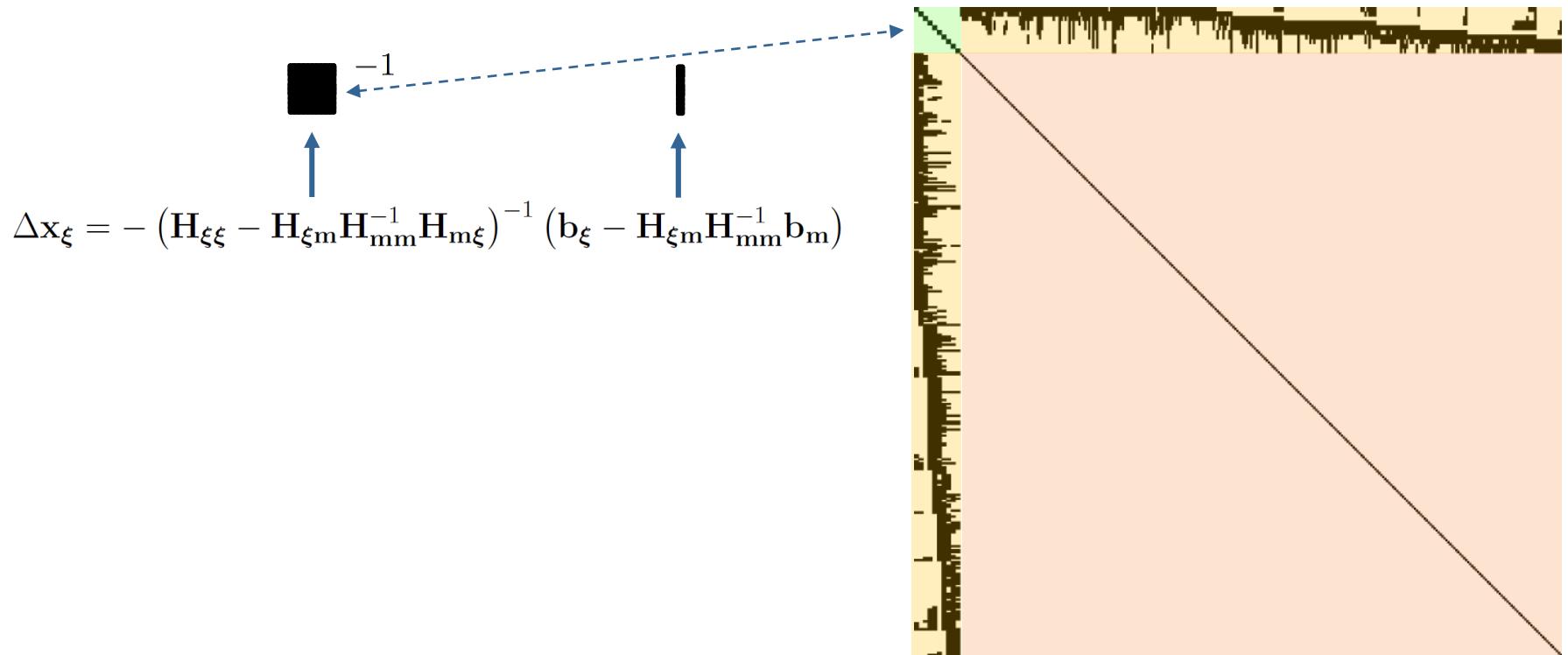
# Exploiting the Sparse Structure

- What is the structure of the two sub-problems ?
- Landmarks:  $\Delta \mathbf{x}_m = -H_{mm}^{-1} (b_m + H_{m\xi} \Delta \mathbf{x}_\xi)$

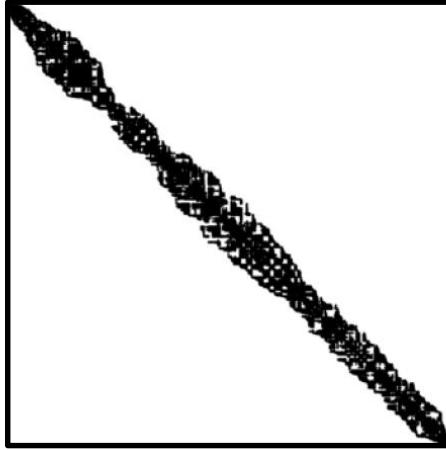
$$\rightarrow \Delta \mathbf{x}_{m_j} = -H_{m_j m_j}^{-1} (b_{m_j} + H_{m_j \xi} \Delta \mathbf{x}_\xi)$$
$$= - \quad \left( \quad + \quad \begin{array}{c|c|c|c} \xi_0 & \text{---} & \xi_t & \end{array} \quad \right)$$

- Landmark-wise solution
- Comparably small matrix operations
- Only involves poses that observe the landmark

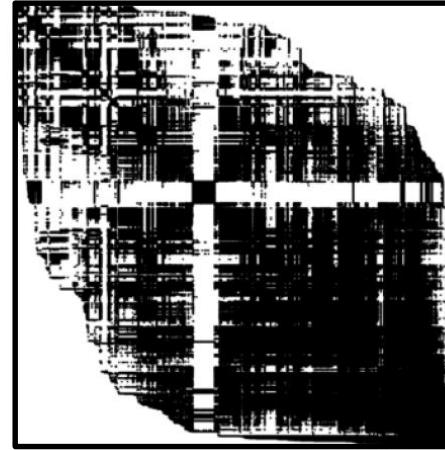
# Exploiting the Sparse Structure



# Exploiting the Sparse Structure



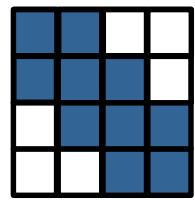
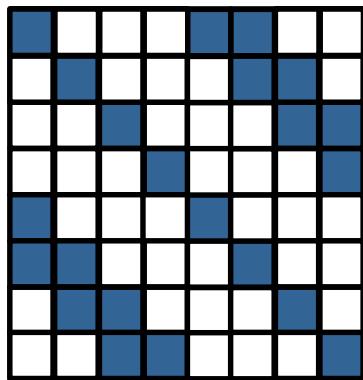
Camera on a moving vehicle  
(6375 images)



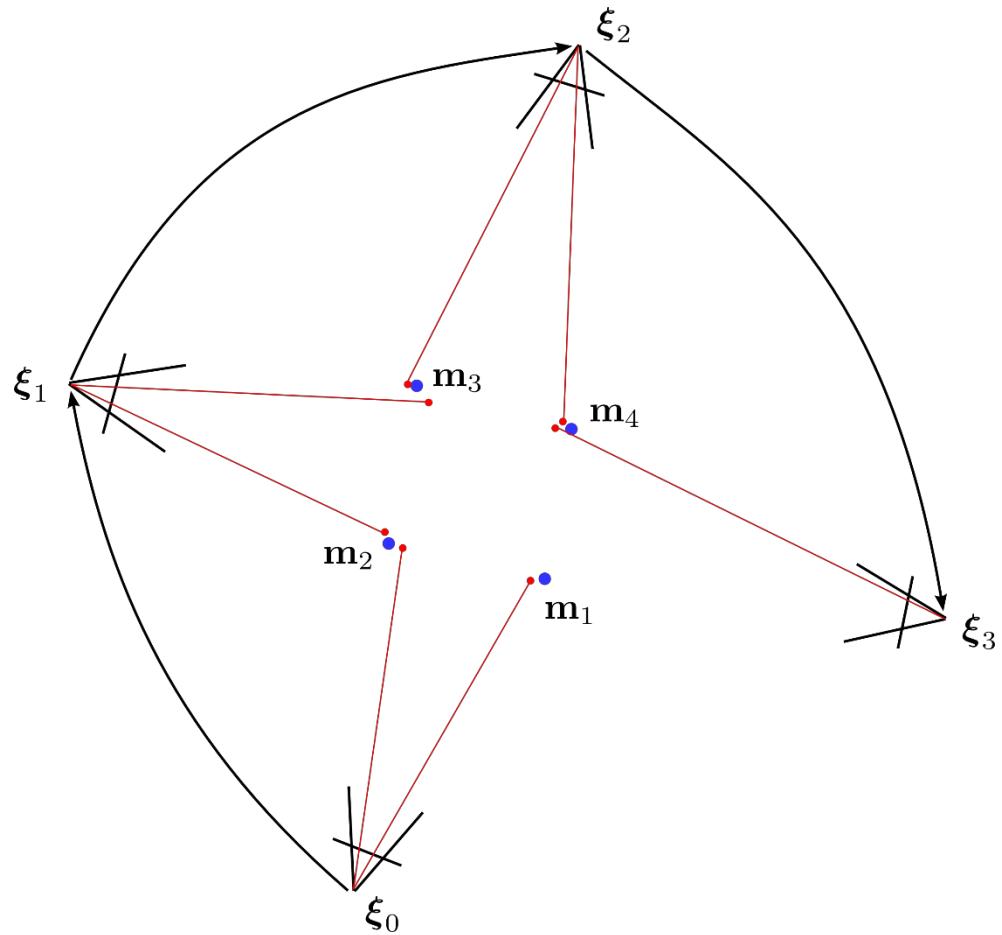
Flickr image search „Dubrovnik“  
(4585 images)

- Reduced pose Hessian can still have sparse structure
- However: For many camera poses with many shared observations, the inversion of the reduced pose Hessian is still computationally expensive!
- Exploit further structure, e.g., using variable reordering or hierarchical decomposition

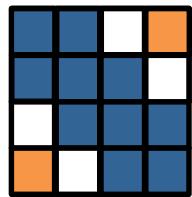
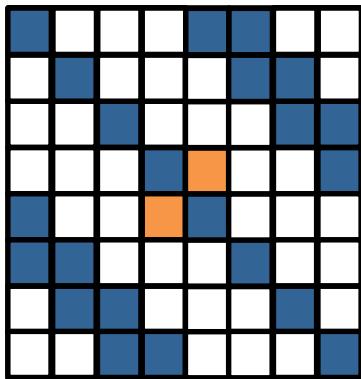
# Effect of Loop-Closures on the Hessian



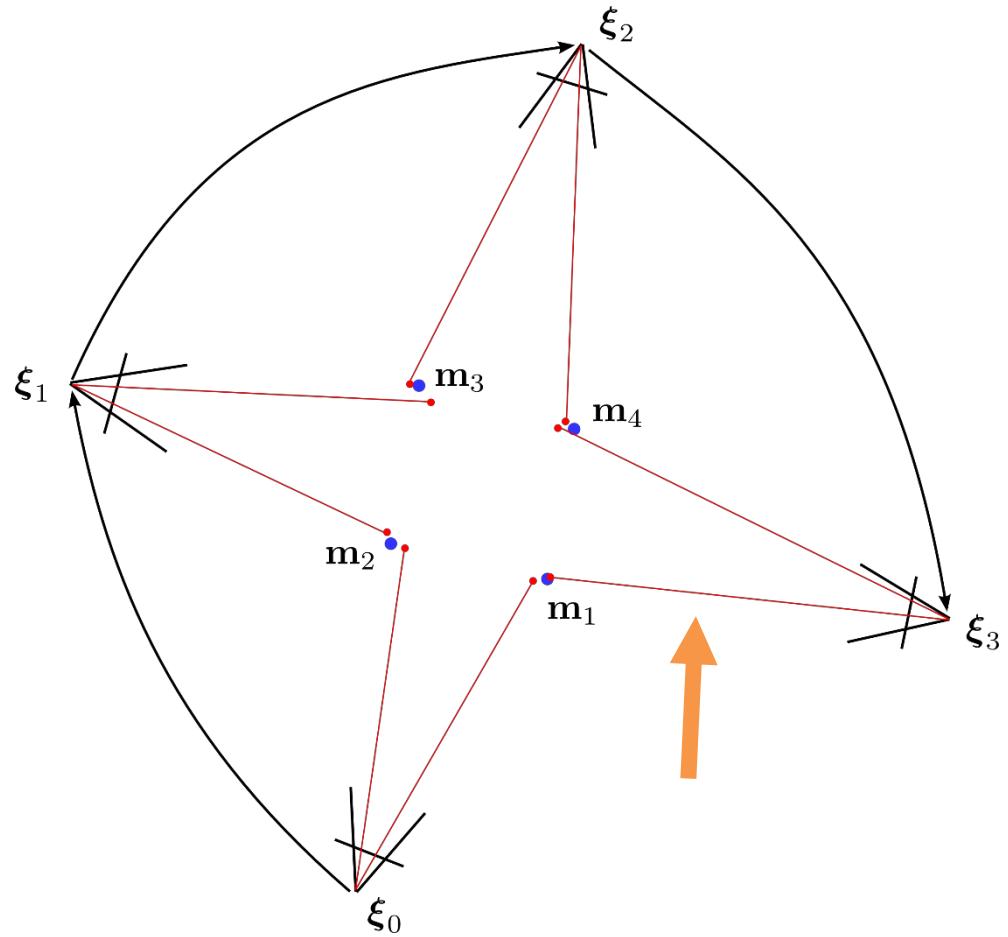
Band matrix



# Effect of Loop-Closures on the Hessian



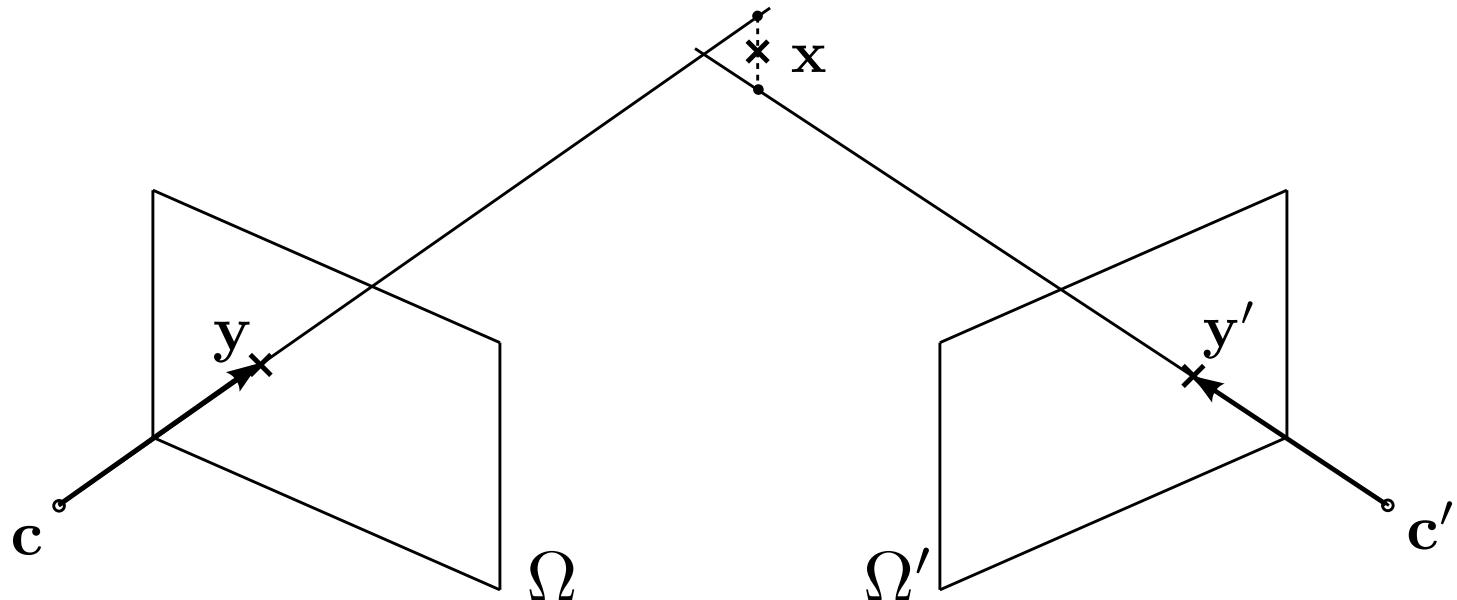
Not band matrix: **costlier** to solve



# Further Considerations

- Use matrix decompositions (f.e. Cholesky) to perform inversions
- Levenberg-Marquardt optimization improves basin of convergence
- Heavier-tail distributions / robust norms on the residuals can be implemented using Iteratively Reweighted Least Squares
- Twists are also a suitable pose parametrization for bundle adjustment: optimize increments on the twists
- Many further tricks to improve convergence/robustness/run-time efficiency, f.e.:
  - Preconditioning
  - Hierarchical optimization
  - Variable reordering
  - Delayed relinearization

# Triangulation



# Triangulation

- Goal: Reconstruct 3D point  $\tilde{\mathbf{x}} = (x, y, z, w)^\top \in \mathbb{P}^3$  from 2D image observations  $\{\mathbf{y}_1, \dots, \mathbf{y}_N\}$  for known camera poses  $\{\mathbf{T}_1, \dots, \mathbf{T}_N\}$
- Linear solution: Find 3D point such that reprojections equal its projections

$$\mathbf{y}'_i = \pi(\mathbf{T}_i \tilde{\mathbf{x}}) = \begin{pmatrix} r_{11}x + r_{12}y + r_{13}z + t_x w \\ r_{31}x + r_{32}y + r_{33}z + t_z w \\ r_{21}x + r_{22}y + r_{23}z + t_y w \\ r_{31}x + r_{32}y + r_{33}z + t_z w \end{pmatrix}$$

get initial values  
just optimize the landmarks,  
keep poses constant

- Each image provides one constraint  $\mathbf{y}_i - \mathbf{y}'_i = 0$
- Leads to system of linear equations  $\mathbf{A}\tilde{\mathbf{x}} = 0$ , two approaches:
  - Set  $w = 1$  and solve nonhomogeneous system
  - Find nullspace of  $\mathbf{A}$  using SVD
- Non-linear solution: Minimize least squares reprojection error (more accurate)

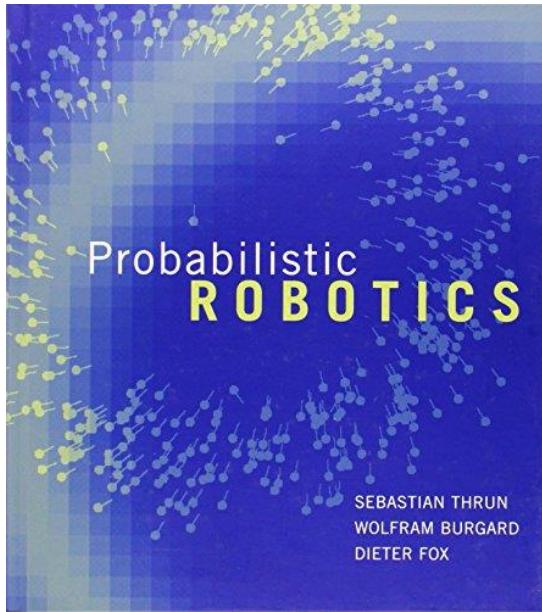
$$\min_{\mathbf{x}} \left\{ \sum_{i=1}^N \|\mathbf{y}_i - \mathbf{y}'_i\|_2^2 \right\}$$

# Lessons Learned Today

- SLAM is a chicken-or-egg problem:
  - Localization requires map
  - Mapping requires localization
  - Unknown association of measurements to map elements
- Bundle Adjustment has a sparse structure that can be exploited for efficient optimization
- Reduction of BA to pose optimization problem through marginalization of landmarks (using the Schur complement)
- Loop closure constraints make SLAM optimization problem less efficient to solve (but reduce drift!)

# Further Reading

- Probabilistic Robotics textbook



Probabilistic Robotics,  
S. Thrun, W. Burgard, D. Fox,  
MIT Press, 2005

- Triggs et al., Bundle Adjustment – A Modern Synthesis, 2002

Thanks for your attention!