

# CS4277 / CS5477

## 3D Computer Vision

### Lecture 10: Structure-from-Motion (SfM) and Bundle Adjustment

Assoc. Prof. Lee Gim Hee

AY 2022/23

Semester 2

# Course Schedule

Week	Date	Topic	Assignments
1	11 Jan	2D and 1D projective geometry	<b>Assignment 0:</b> Getting started with Python (Ungraded)
2	18 Jan	3D projective geometry, Circular points and Absolute conic	
3	25 Jan	Rigid body motion and Robust homography estimation	
4	01 Feb	Camera models and calibration	<b>Assignment 1:</b> Metric rectification and robust homography (10%) <b>Due:</b> 2359hrs, 07 Feb
5	08 Feb	Single view metrology	<b>Assignment 2:</b> Affine 3D measurement from vanishing line and point (10%) <b>Due:</b> 2359hrs, 14 Feb
6	15 Feb	The Fundamental and Essential matrices	
-	22 Feb	Semester Break	No lecture
7	01 Mar	<b>Mid-term Quiz (20%)</b> Lecture: Generalized cameras	<b>In-person Quiz (LT 15, 1900hrs – 2000hrs)</b> Lecture: 2000hrs – 2130hrs
8	08 Mar	Absolute pose estimation from points or lines	
9	15 Mar	Three-view geometry from points and/or lines	
10	22 Mar	Structure-from-Motion (SfM) and bundle adjustment	<b>Assignment 3:</b> SfM and Bundle adjustment (10%) <b>Due:</b> 2359hrs, 28 Mar
11	29 Mar	Two-view and multi-view stereo	<b>Assignment 4:</b> Dense 3D model from multi-view stereo (10%) <b>Due:</b> 2359hrs, 04 Apr
12	05 Apr	3D Point Cloud Processing	
13	12 Apr	Neural Field Representations	

**Final Exam: 03 MAY 2023**

# Learning Outcomes

- Students should be able to:
  1. Describe the pipeline of large-scale **3D reconstruction**: data association, structure-from-motion and dense stereo.
  2. Explain the use of **robust two-view geometry** and the **bag-of-words** algorithm for data association.
  3. Use two-view geometry, PnP and triangulation to **initialize** the 3D reconstruction.
  4. Apply the iterative methods: Newton, Gauss-Newton Gradient descent or Levenberg-Marquardt for **bundle adjustment**.

# Acknowledgements

- A lot of slides and content of this lecture are adopted from:
  1. Tutorial on large-scale 3D modeling, CVPR 2017  
<https://demuc.de/tutorials/cvpr2017/>
  2. R. Hartley, and A. Zisserman: “Multiple view geometry in computer vision”, Appendix 6.
  3. Bill Triggs et al, “Bundle Adjustment – A Modern Synthesis”, 1999
  4. Y. Ma, S. Soatto, J. Kosecka, S. S. Sastry, “ An invitation to 3-D vision”, Chapter 11.

# Large-Scale 3D Reconstruction

- **Given:** a set of images  $\{I_1, \dots, I_N\}$  taken from different viewpoints.
- **Find:**
  1. The motions of the cameras w.r.t. a world coordinate frame  $F_W$ , i.e. the camera projection matrices  $\{P_1, \dots, P_N\}$ .
  2. And then recover the 3D structures  $\{X_1, \dots, X_M\}$  of the scene.

# Large-Scale 3D Reconstruction: Pipeline

- Three key parts of the pipeline:
  1. **Data association:** check whether a pair of images are related using **image correspondences** and **robust two-view geometry**.
  2. **Structure-from-Motion (SfM):** initial reconstruction from **relative/absolute pose** estimation and **triangulation**, and then refine with **bundle adjustment**.
  3. **Multi-view stereo (MVS):** get the dense 3D model from the **plane sweeping algorithm** (next lecture).

# Large Scale 3D Reconstruction: Pipeline

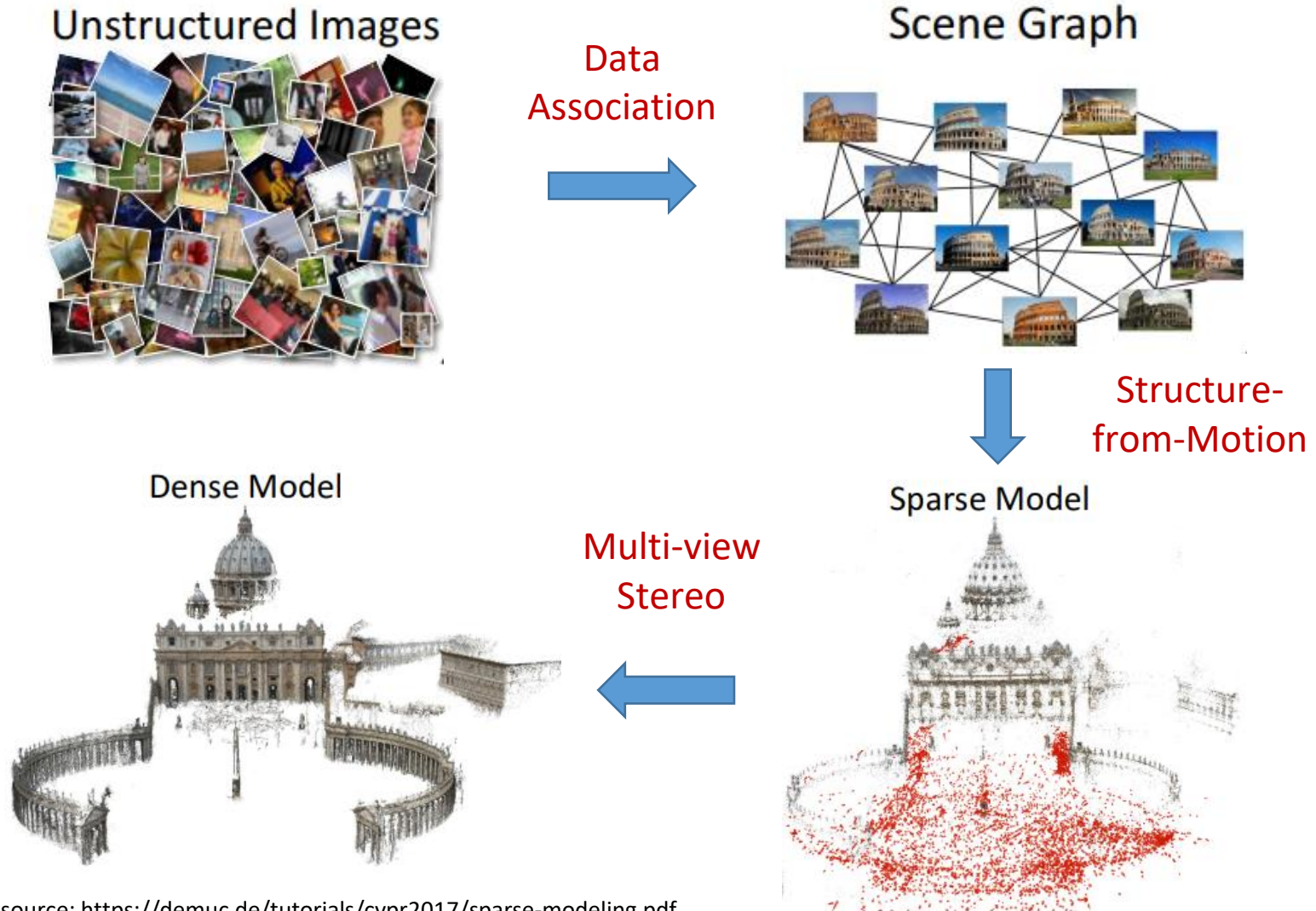


Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>

# Structure-from-Motion

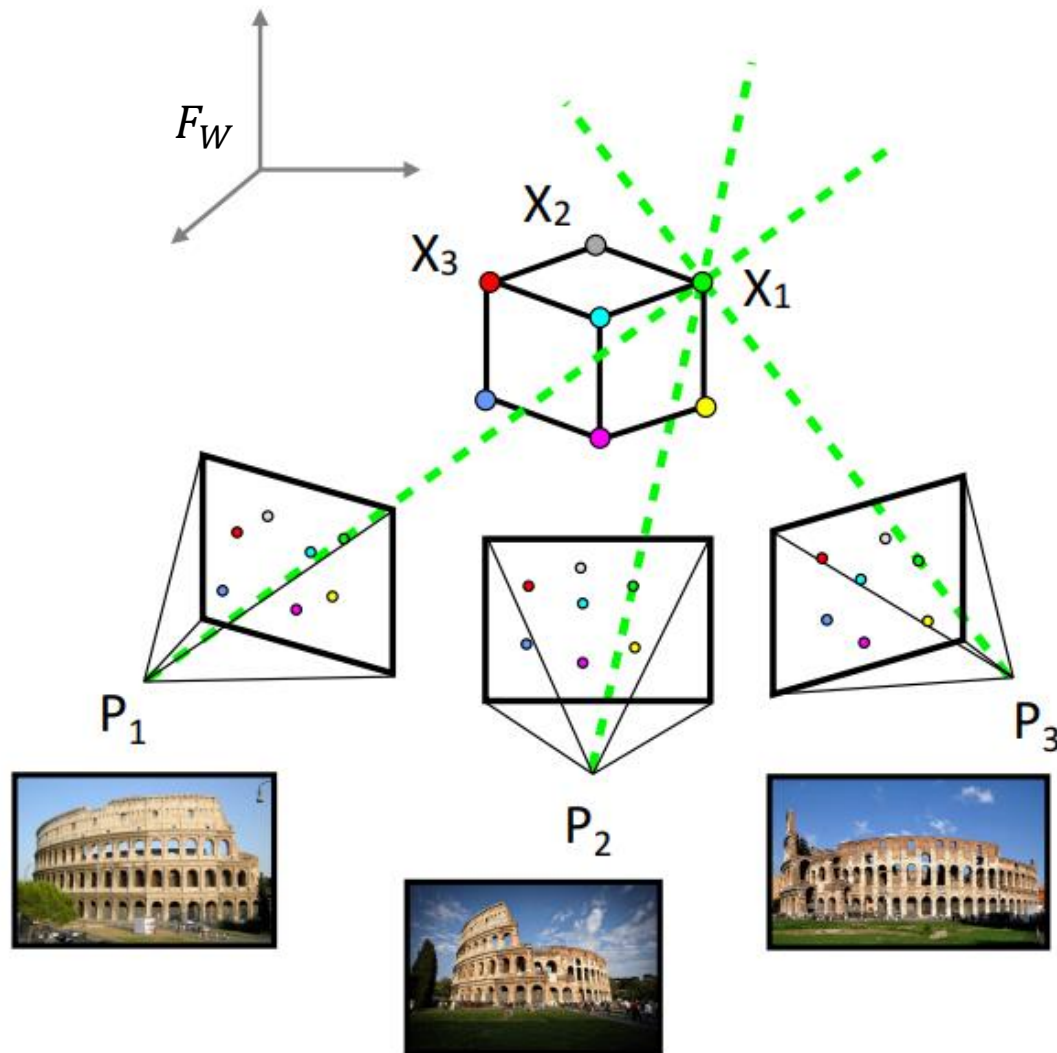


Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>



# Structure-from-Motion

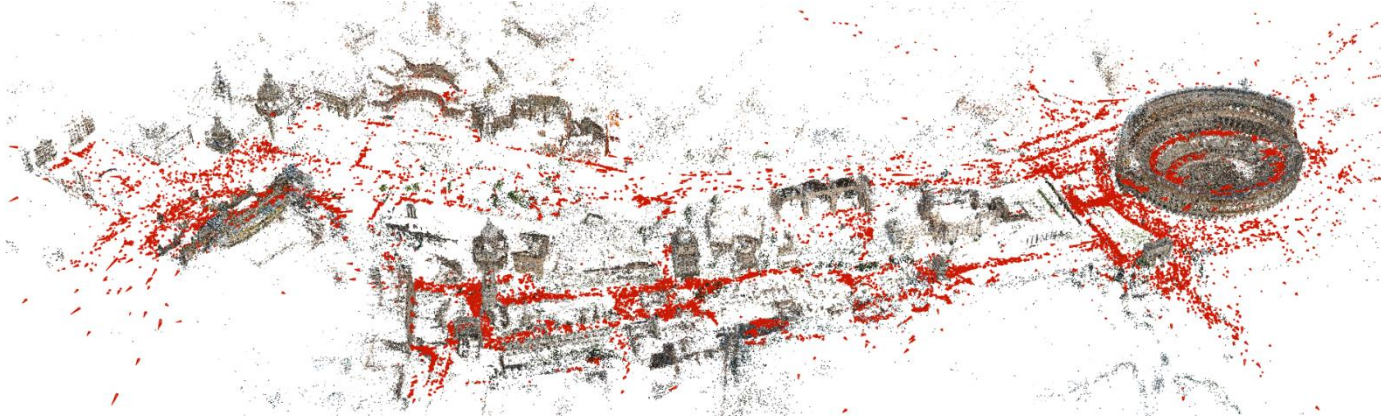
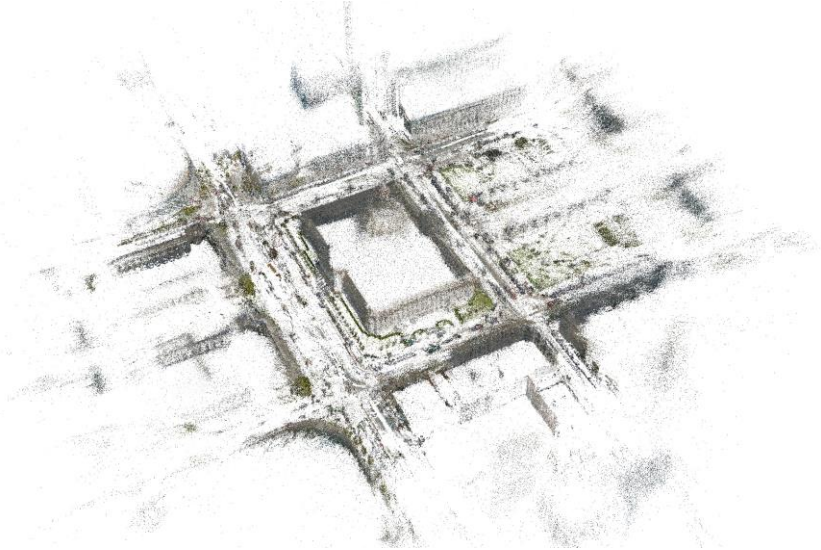
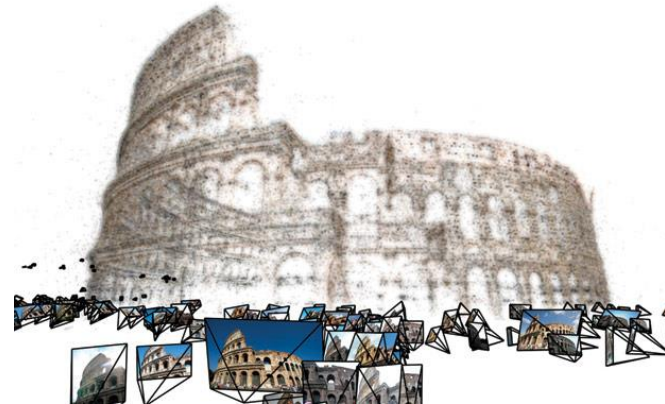


Image Source: COLMAP, <https://colmap.github.io/>



O. Saurer, M. Pollefeys, G. H. Lee, CVPR 2016



S. Agarwal, N. Snavely, I. Simon, S. Seitz and R. Szeliski, ICCV 2009

# Data Association

- We get the scene graph by finding images with **overlapping views**, i.e. connected components.

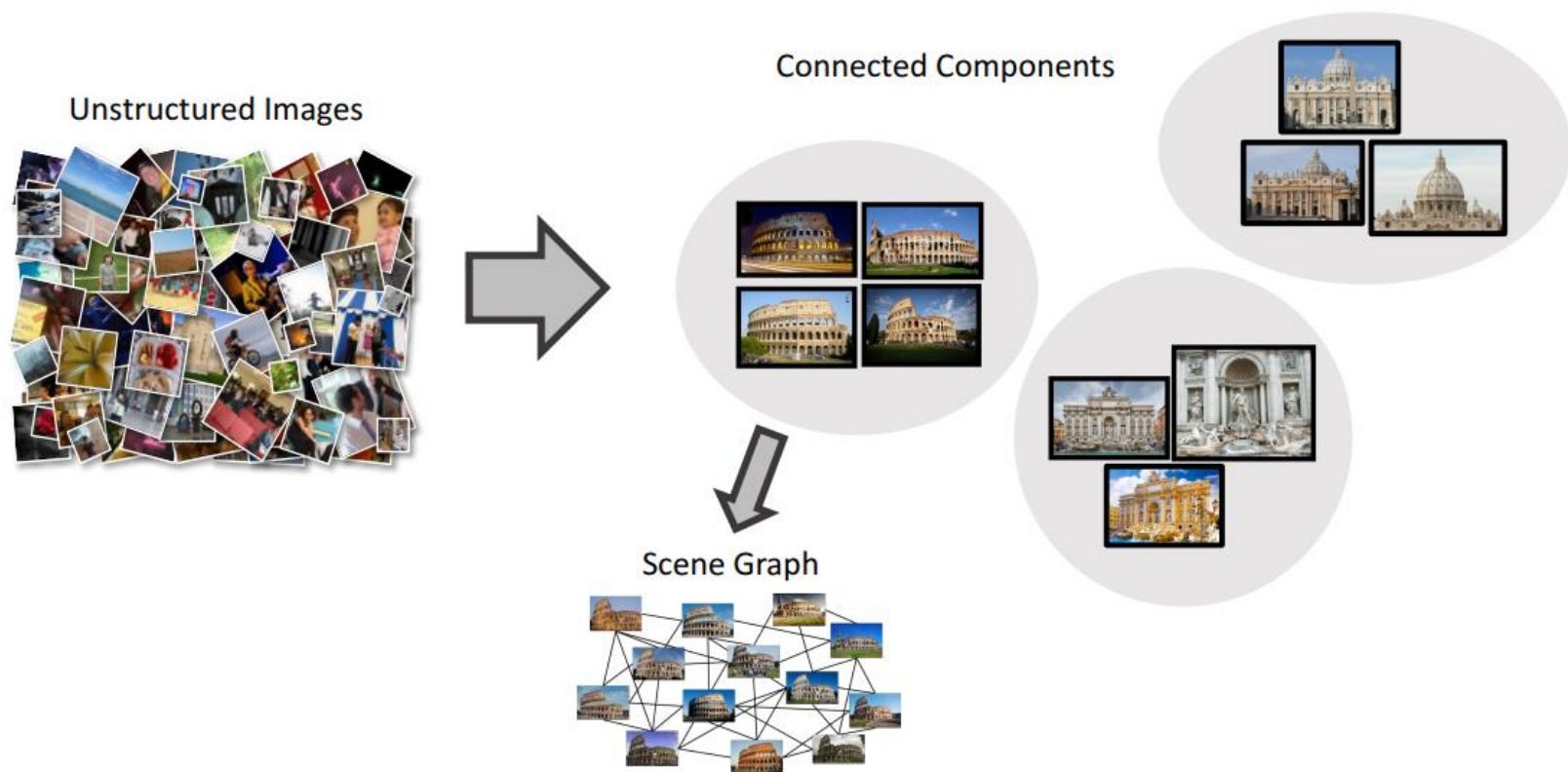


Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>

# Data Association

- Use **two-view geometry** to establish the connected components.

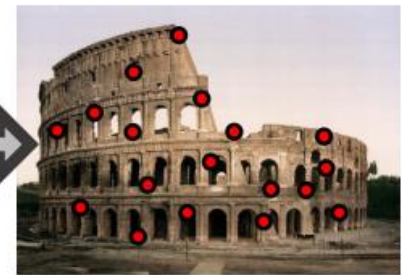
Unstructured Images



Two-View Geometry



?



Scene Graph



Image source:  
<https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>



# Data Association

**Step 1:** extract **image keypoints** e.g. SIFT [1] or ORB [2].

**Step 2:** establish the **putative correspondences**.

**Step 3:** **geometric verification** to get correct image pairs.

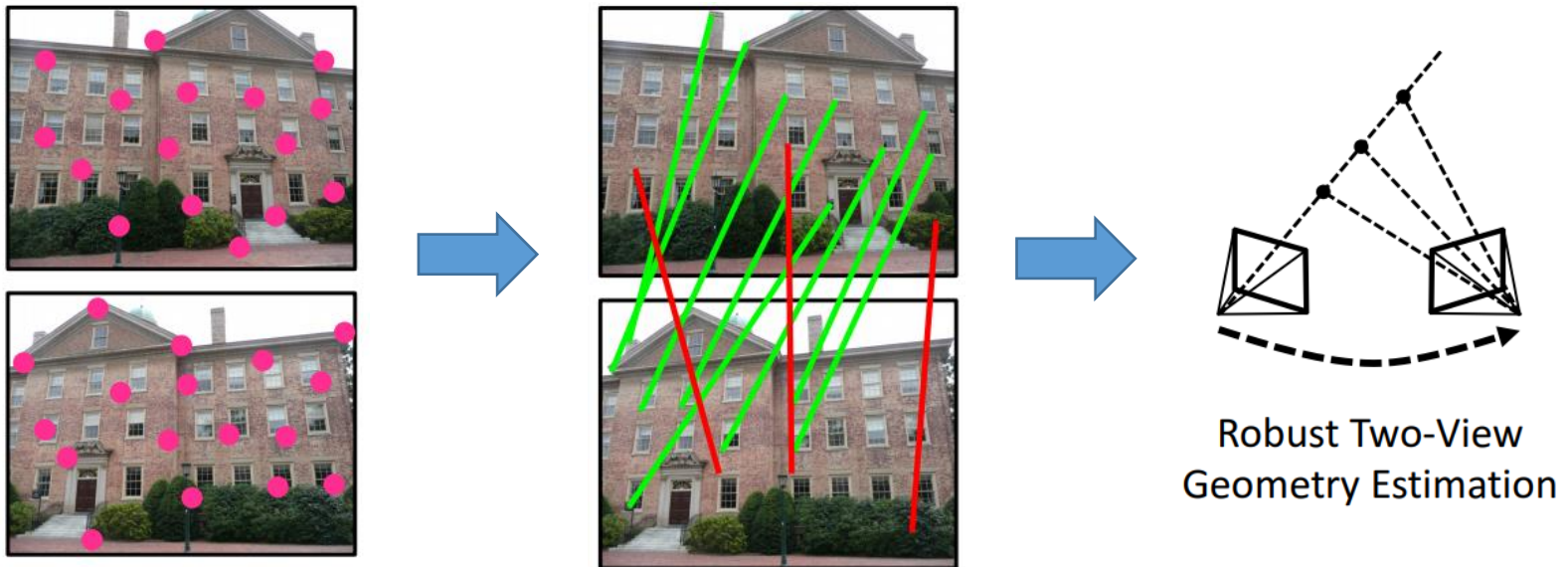


Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>

[1] David G. Lowe, "Distinctive Image Features from Scale-Invariant Keypoints", IJCV 2004.

[2] Ethan et. al., "ORB: An efficient alternative to SIFT or SURF", ICCV 2011.

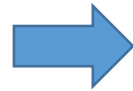
# Data Association

## Geometric Verification:

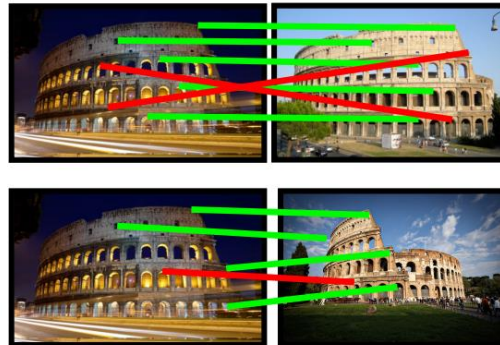
- The image keypoint matchings is only based on **appearance** and might be wrong due to **outliers**.
- We do a RANSAC-based two-view geometry, e.g. homography, fundamental or essential matrix to **detect outlier** matches **with geometry**.
- Image pairs with **inlier counts > threshold** are added to the scene graph.

# Data Association

Unstructured Images



Inlier/outlier correspondences



Scene Graph



**Problem:** Exhaustively searching through pairs of images in the set of  $N$  images **is intractable** when  $N$  is large!

**Complexity** of querying one image is  $\mathcal{O}(N \cdot K^2)$ , where  $K$  is the number of keypoints in each image.

Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>

# Data Association

## Example:

- Assume 1,000 SIFT features per image  $\Rightarrow K = 1,000$ .
- Assume  $N = 100,000,000$  images.
- $\Rightarrow N \cdot K^2 = 100,000,000,000,000$  feature comparisons!
- If we assume 0.1 ms per feature comparison  $\Rightarrow 1$  image query would take **317 years!**

# Place Recognition / Image Retrieval

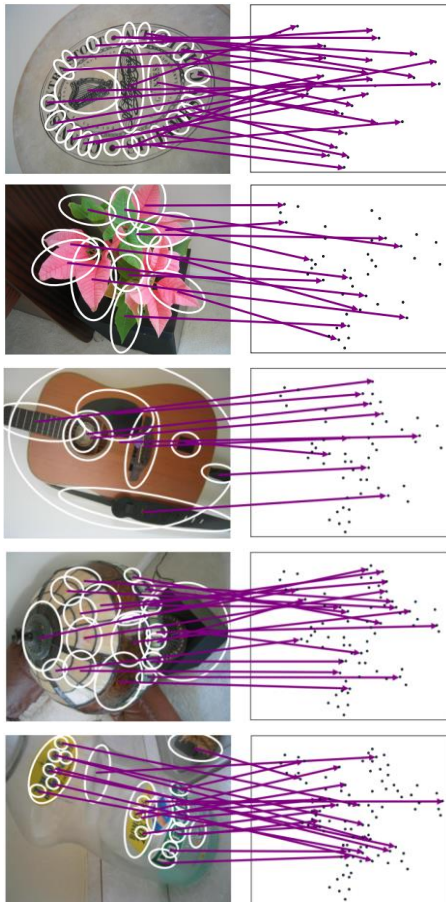
**Solution:** Use **Bag-of-Words** image retrieval!

- The goal is to build an **efficient tree-based search** algorithm for matching the query image features with the image features from the database.



# Place Recognition / Image Retrieval

**Step 1:** Extract **image keypoints** and **descriptors** from the training data.



For visualization we assume that the descriptor space has **2 dimensions**. In the case of SIFT, the descriptor space would have **128 dimensions**.



Image source: [http://rpg.ifi.uzh.ch/docs/teaching/2019/12b\\_recognition.pdf](http://rpg.ifi.uzh.ch/docs/teaching/2019/12b_recognition.pdf)

# Place Recognition / Image Retrieval

**Step 2:** Perform **hierarchical clustering** on the descriptors. Each leaf node is a “**visual word**” and the set forms a “**bag-of-words**”.

◆ Root node, ● 1<sup>st</sup> level node, ● Leaf node

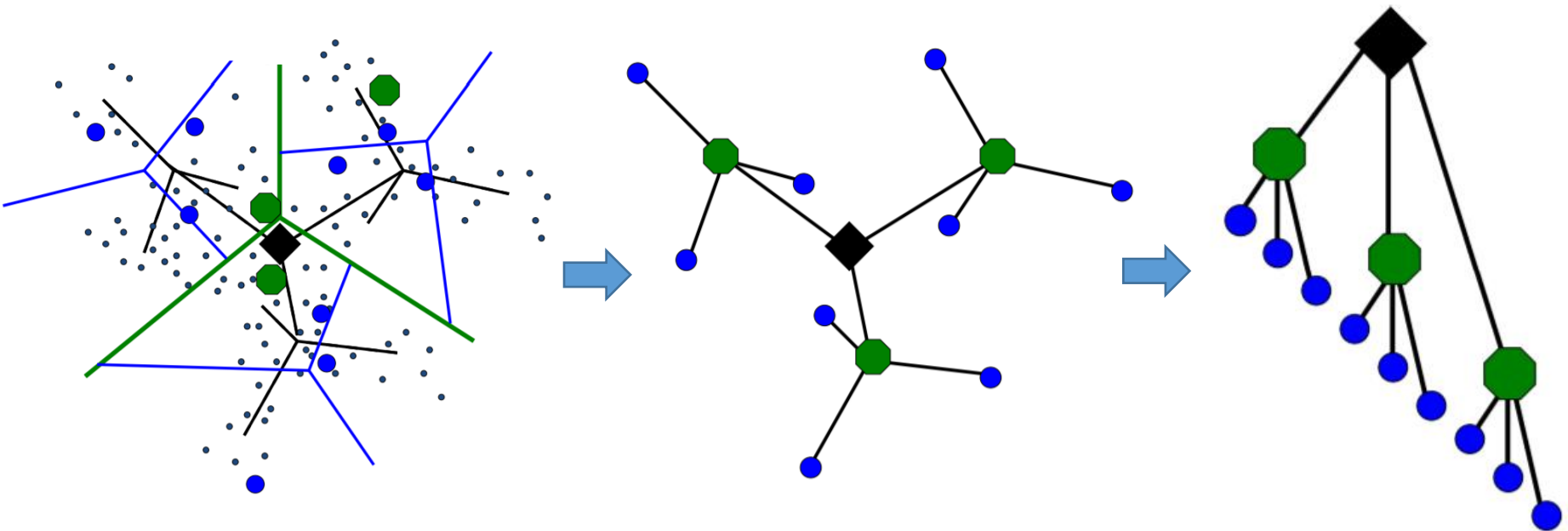


Image source: Nister et al, CVPR 2006

# Place Recognition / Image Retrieval

Step 3: Build the **inverted file index**.

- i. Extract the **image features** for each image in the database.
- ii. Do a tree search on each image feature for the **closest leaf nodes**.
- iii. Store the **frequency count** of the image features and the **image ID** on each leaf node as a global descriptor.

# Place Recognition / Image Retrieval

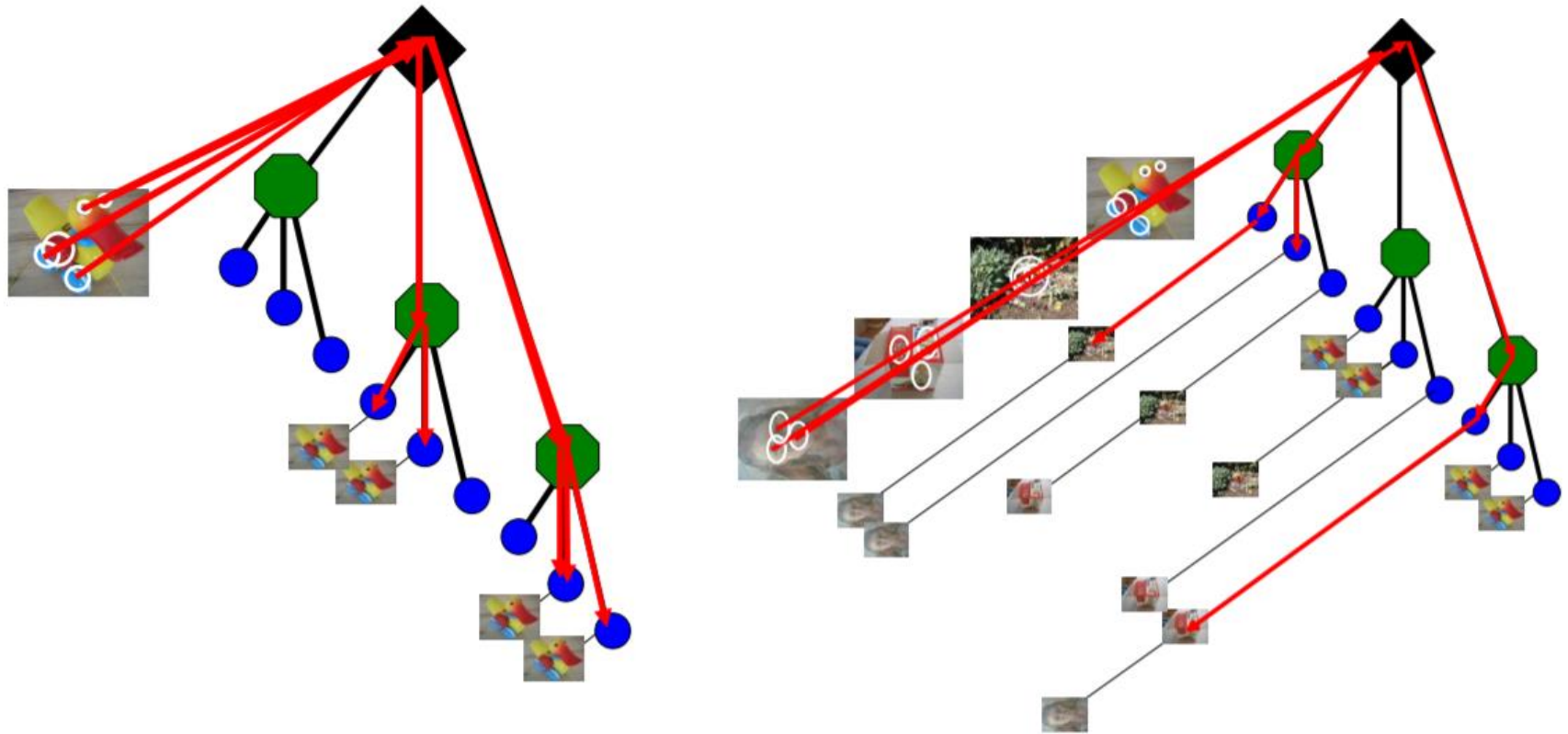


Image source: Nister et al, CVPR 2006

# Place Recognition / Image Retrieval

**Step 4:** Given a query image, retrieve the **visually most similar** image from the database.

- i. Extract **image features** from the query image.
- ii. For each image feature, find the **closest leaf node** in the hierarchical k-mean tree.
- iii. **Increase the counter** to the database image ID(s) with a feature in that leaf node.
- iv. Choose the database image with the **highest count** (normalized by the frequency count) as the visually most similar image.

# Place Recognition / Image Retrieval

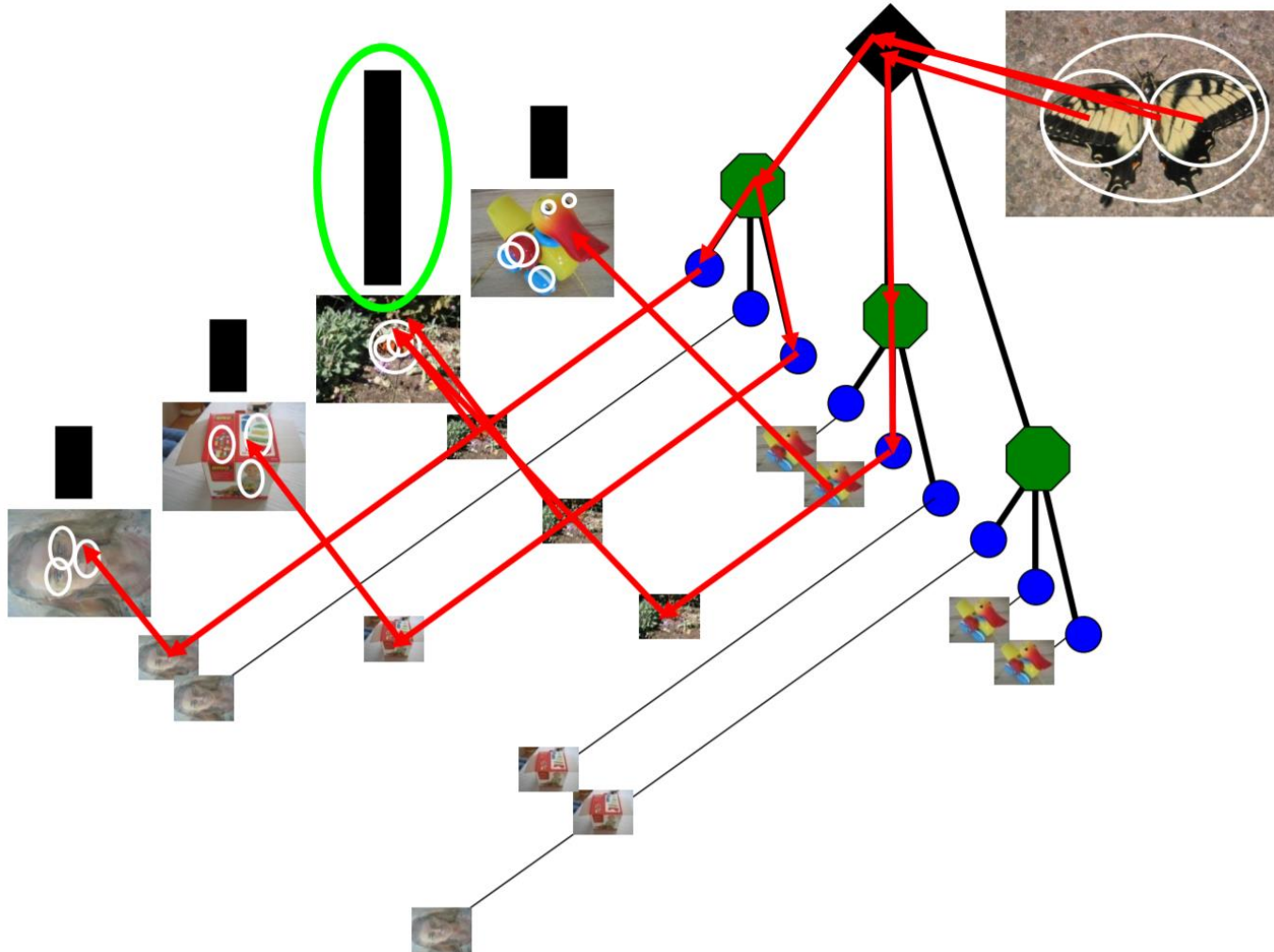


Image source: Nister et al, CVPR 2006

# Place Recognition / Image Retrieval

**Example:** Querying an image in a database of **100 million images**.

- Assume a query image with 1,000 SIFT features  $\Rightarrow K = 1,000$ .
- Assume 10 branches and 6 depth levels (i.e.,  $b^L = 1,000,000$  visual words)
- $\Rightarrow$  Number of feature comparisons  $= K \cdot b \cdot L = 1,000 \cdot 10 \cdot 6 = 60,000$
- If we assume 0.1 ms per feature comparison  $\rightarrow$  1 image query would take **6 seconds**!

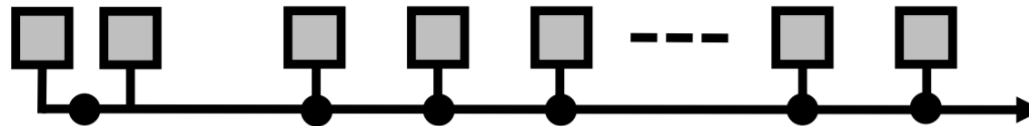
# Place Recognition / Image Retrieval

- **Reduced complexity:** For  $K$  features in the Query image, only  $\mathcal{O}(K \cdot b \cdot L)$  comparisons need to be made instead of  $\mathcal{O}(N \cdot K^2)$ .
- We first apply image retrieval to **eliminate the infeasible edges** in the scene graph.
- Since image retrieval is appearance-based, we then apply **geometric verification** on the remaining edges in the scene graph.

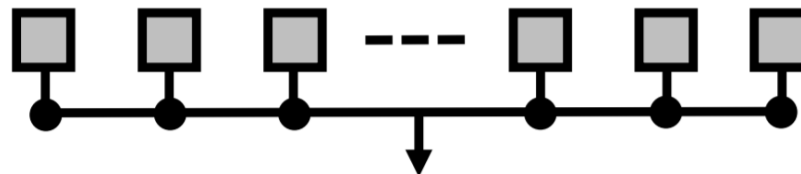


# Three Paradigms of SfM

## 1. Incremental:



## 2. Global:



## 3. Hierarchical:

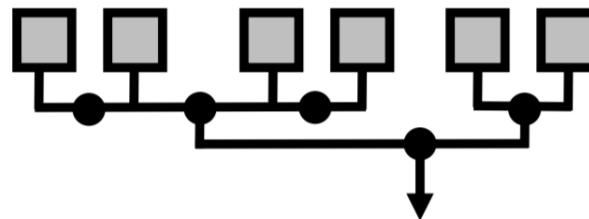
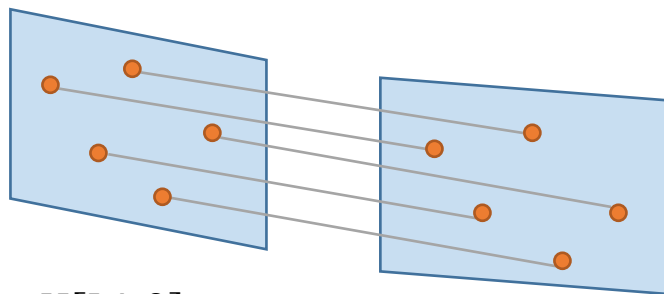


Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>

# Incremental SfM

## Step 1: Initialization

- i. Choose **two non-panoramic views** ( $\|\mathbf{t}\| \neq 0$ ). A good choice is the image pair with the **highest inlier count**.
- ii. Use **8-point algorithm** (F or E matrix) for non-planar scene or **4-point algorithm** (homography) for planar scene.



$$P = K[I \mid 0]$$

$$P' = K'[R' \mid \mathbf{t}']$$

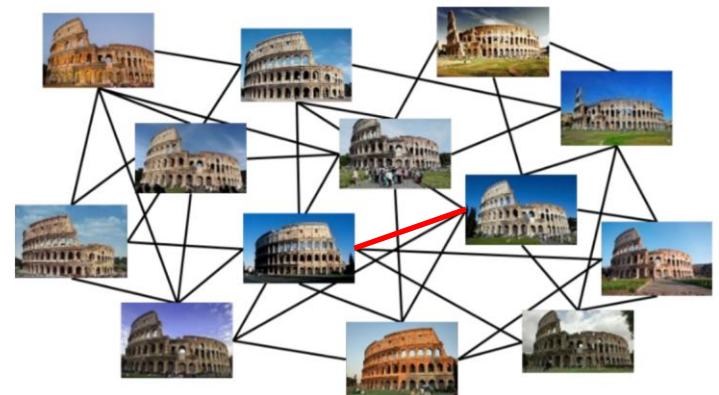
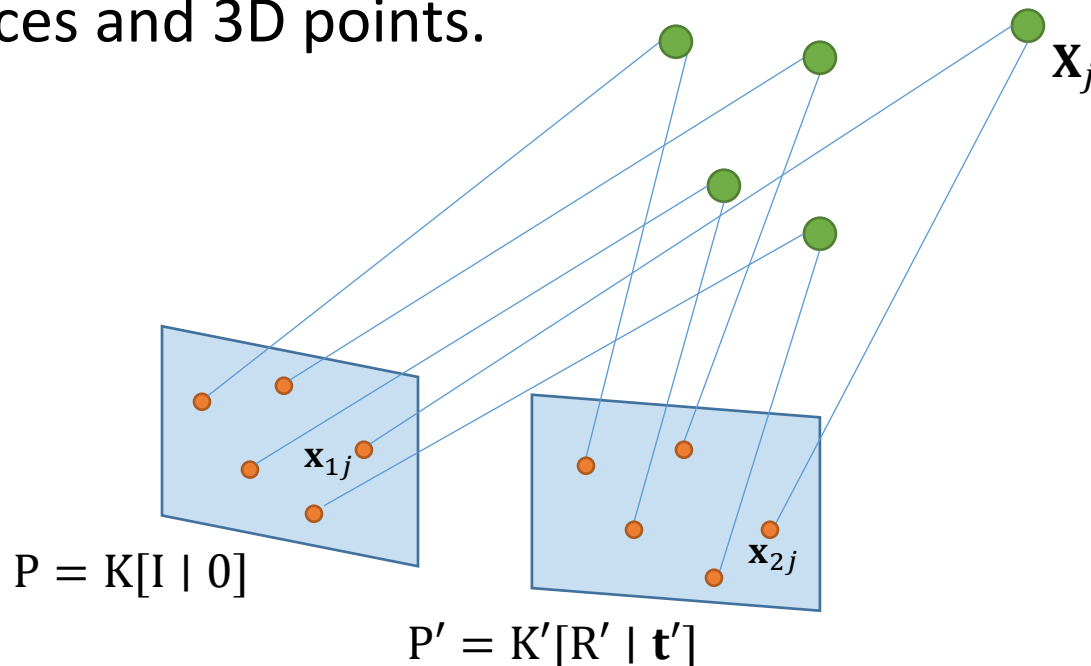


Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>

# Incremental SfM

## Step 1: Initialization

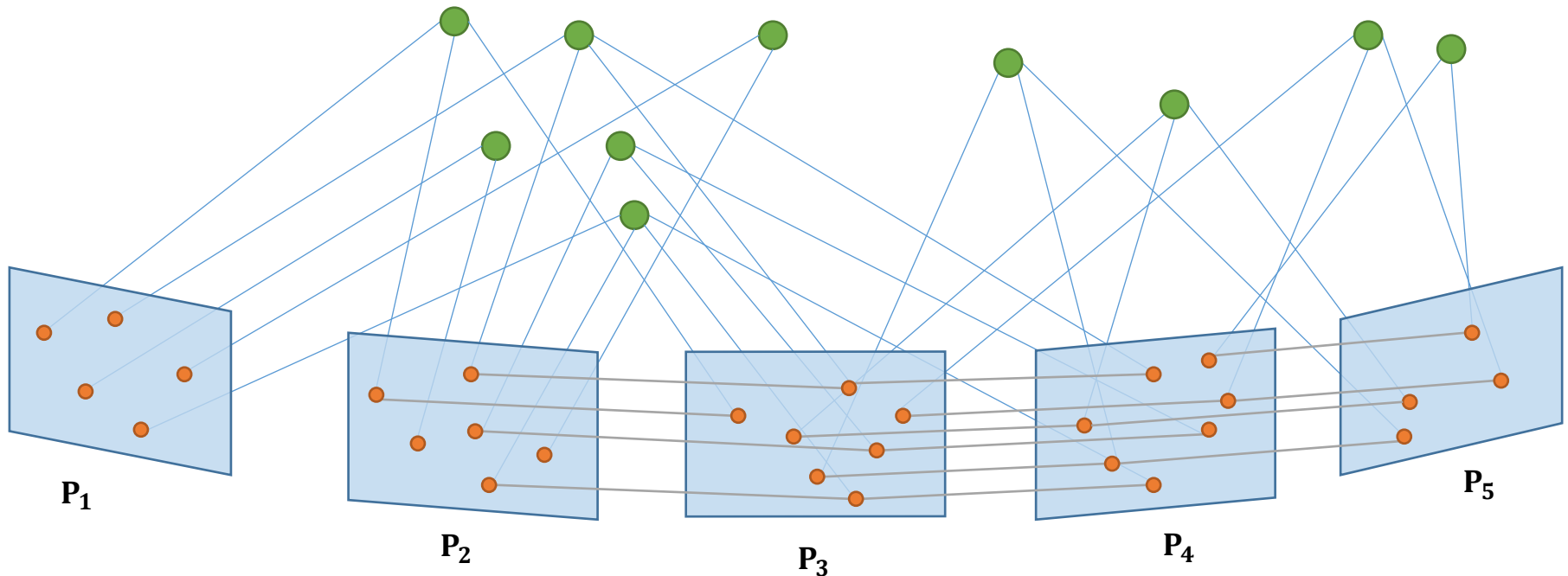
- iii. Set the **scale to 1**, i.e.  $\|\mathbf{t}\| = 1$ .
- iv. Triangulate inlier correspondences to get the **3D points**.
- v. Apply **bundle adjustment** (detail later) to refine camera matrices and 3D points.



# Incremental SfM

## Step 2: Subsequent Views

- i. Find 2D-3D correspondences.
- ii. Solve the **Perspective-n-Point** (PnP) problem.



# Incremental SfM

**Step 3:** Refine the multi-view camera poses and structures with **bundle adjustment** (detail later).

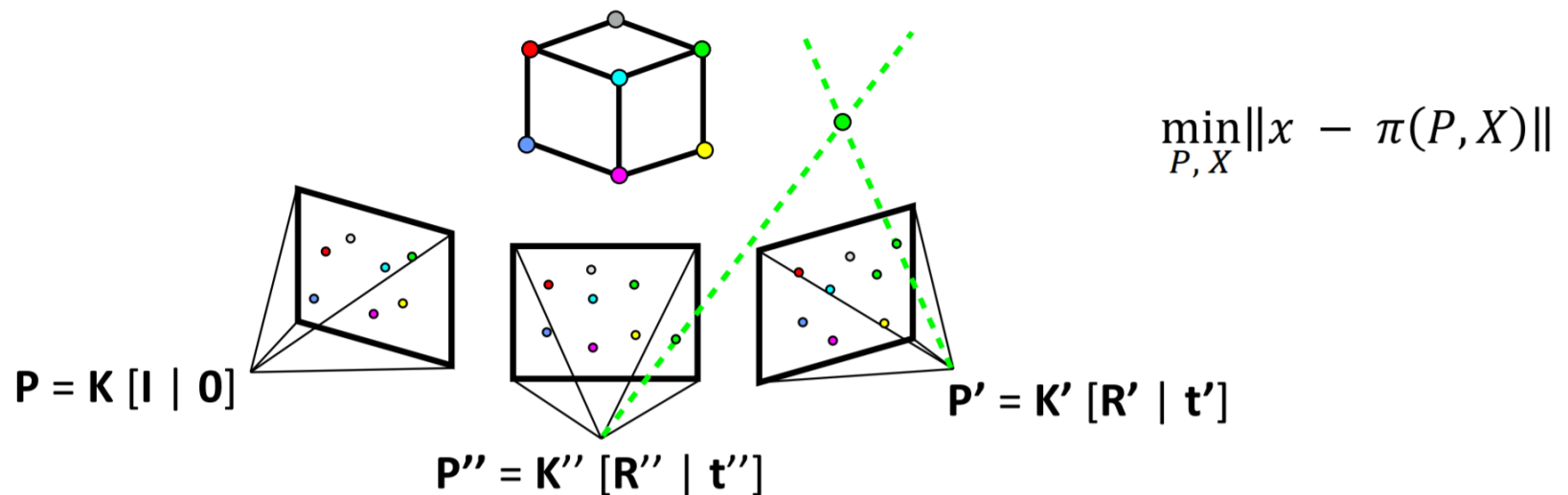


Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>

# Global SfM

1. Compute all **relative poses**  $(R, t)$  in the edges of the scene graph.

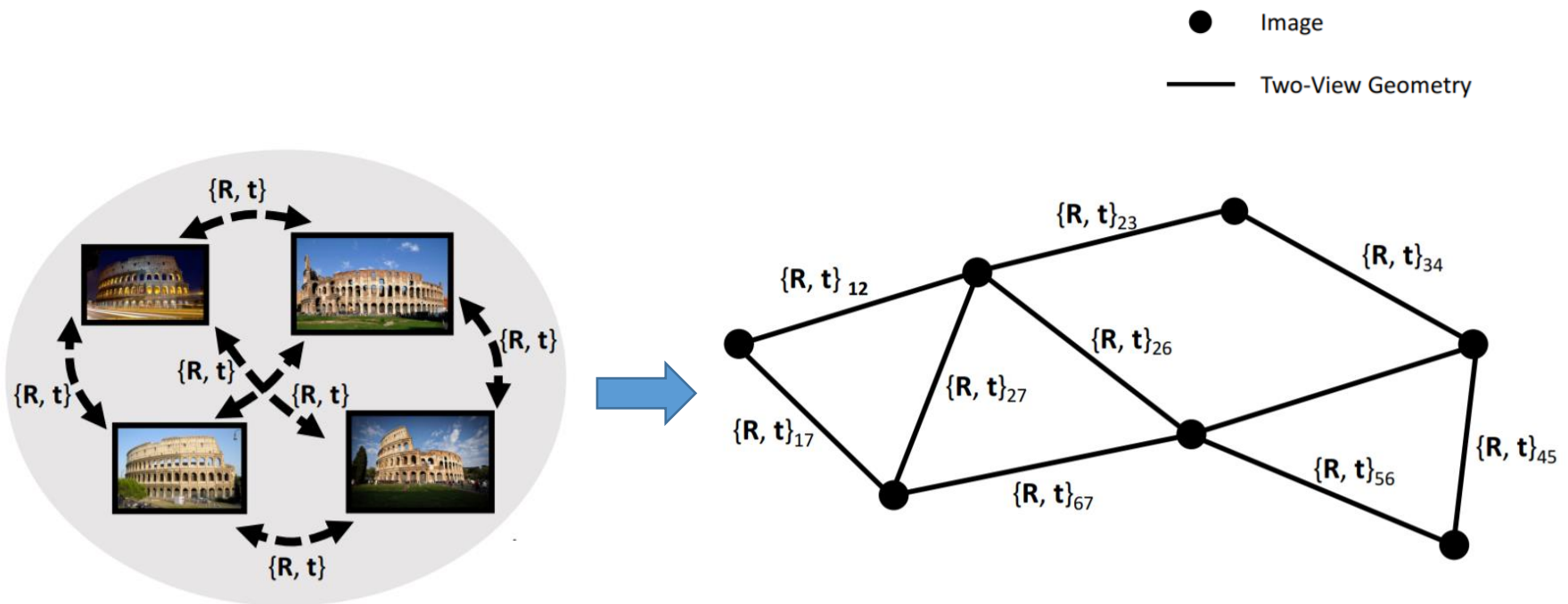


Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>

**Reference:** Chatterje and Govindu, "Efficient and Robust Large-Scale Rotation Averaging", 2013

# Global SfM

2. Estimate the **global rotations**:  $\min_{\mathbf{R}} \|\mathbf{R}_{ij} - \mathbf{R}_i \mathbf{R}_j^T\|$ ,  
where  $\mathbf{R}_{ij}$  are the relative rotations and  $\mathbf{R}_i$  are the global rotations.

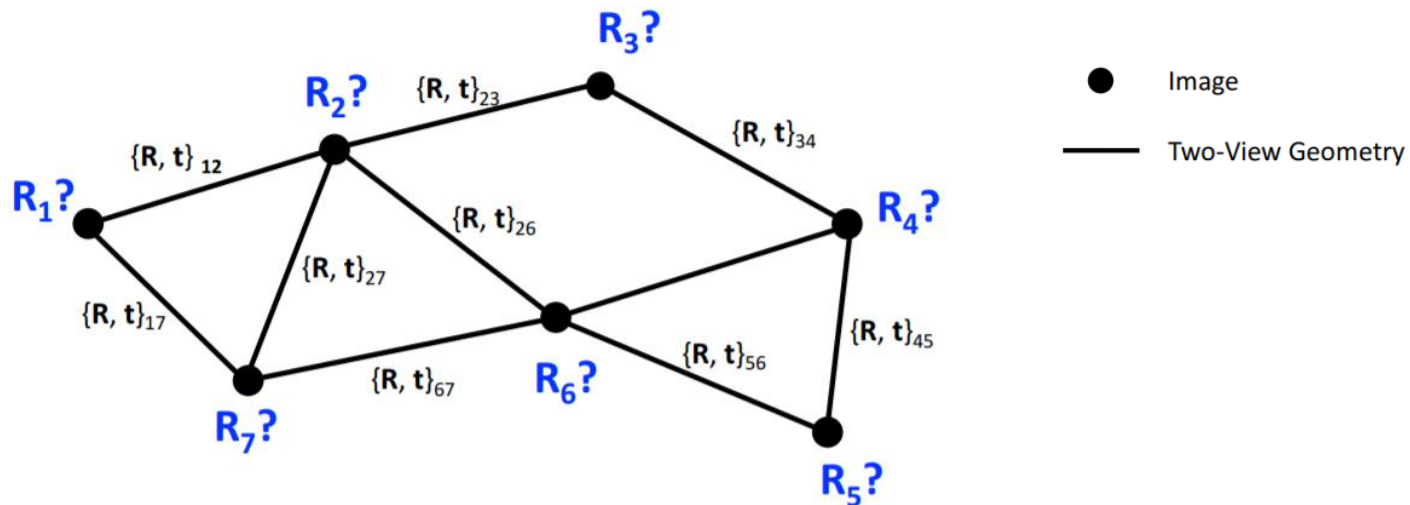


Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>

**Reference:** Chatterje and Govindu, "Efficient and Robust Large-Scale Rotation Averaging", 2013

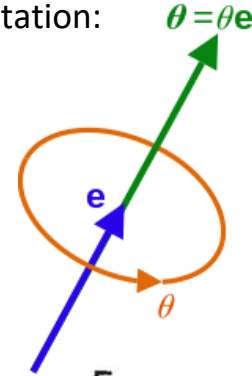
# Global SfM

- We define the **logarithmic map** as:  $\boldsymbol{\theta} = \hat{\boldsymbol{\omega}} = \log(R) \in so(3)$ .

- Specifically, for  $R = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix}$ ,

$$\theta = \|\boldsymbol{\omega}\| = \cos^{-1} \left( \frac{\text{trace}(R) - 1}{2} \right), \quad \mathbf{e} = \frac{\hat{\boldsymbol{\omega}}}{\|\boldsymbol{\omega}\|} = \frac{1}{2 \sin(\|\boldsymbol{\omega}\|)} \begin{bmatrix} r_{32} - r_{23} \\ r_{13} - r_{31} \\ r_{21} - r_{12} \end{bmatrix}.$$

Angle-axis representation  
of rotation:



- If  $R = I$ , then  $\|\boldsymbol{\omega}\| = 0$ , and  $\frac{\hat{\boldsymbol{\omega}}}{\|\boldsymbol{\omega}\|}$  is **not determined** and therefore can set to  $\frac{\hat{\boldsymbol{\omega}}}{\|\boldsymbol{\omega}\|} = [0, 0, 0]^T$ .

Image source: [https://en.wikipedia.org/wiki/Axis%E2%80%93angle\\_representation](https://en.wikipedia.org/wiki/Axis%E2%80%93angle_representation)



# Global SfM

- The inverse is given by the **exponential map**, i.e.  $R = \exp(\hat{\omega}) \in SO(3)$ .
- This is also known as the **Rodrigues' formula** for a rotation matrix:

$$R = e^{\hat{\omega}} = I + \frac{\hat{\omega}}{\|\omega\|} \sin(\|\omega\|) + \frac{\hat{\omega}^2}{\|\omega\|^2} (1 - \cos(\|\omega\|)) ,$$

where

$$K = [\mathbf{e}]_{\times},$$

$$K^2 = \mathbf{e}\mathbf{e}^T - I.$$

# Global SfM

- $R_{ij} = R_i R_j^T$  can be written as:

$$\hat{\omega}_{ij} = \hat{\omega}_j - \hat{\omega}_i = \underbrace{\left[ \cdots - \mathbf{I} \cdots \mathbf{I} \cdots \right]}_{\mathbf{A}_{ij}} \hat{\omega}_{global} .$$

- Then,  $\min_{\mathbf{R}} \|R_{ij} - R_i R_j^T\|$  can be solved as a **least squares problem**:  $\mathbf{A} \hat{\omega}_{global} = \hat{\omega}_{rel}$ .
- $\hat{\omega}_{rel}$  is the vector made by **stacking all** relative rotation observations  $\hat{\omega}_{ij}$  and  $\mathbf{A}$  is made by stacking the corresponding matrices  $\mathbf{A}_{ij}$

# Global SfM

## 2. Estimate **global translations**:

$$\min_{\mathbf{t}} \left\| \mathbf{t}_{ij} - \frac{\mathbf{t}_i - \mathbf{t}_j}{\|\mathbf{t}_i - \mathbf{t}_j\|} \right\| \quad \text{solve as} \quad \mathbf{Ax} = \mathbf{b}$$

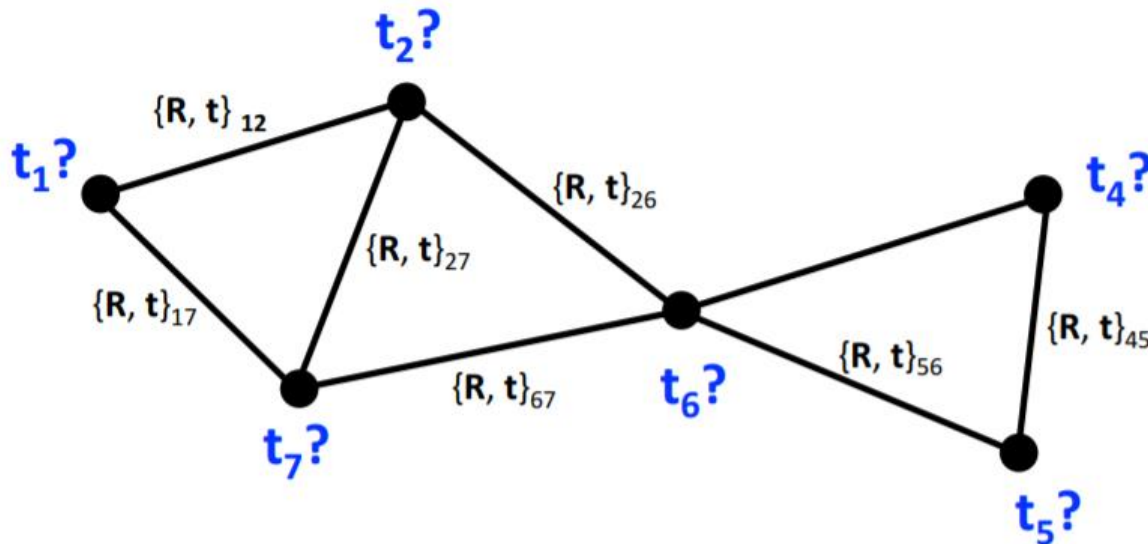


Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>

**Reference:** Chatterje and Govindu, "Efficient and Robust Large-Scale Rotation Averaging", 2013

# Global SfM

## 3. Triangulate and refine with **bundle adjustment** (detail later).

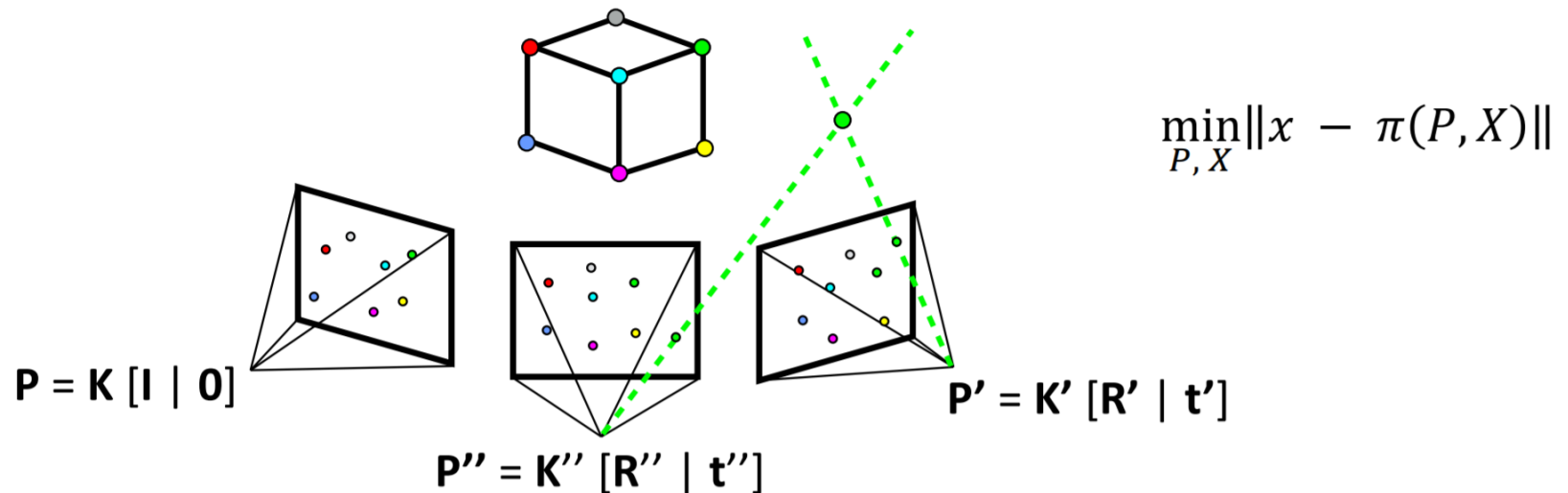


Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>

# Hierarchical SfM

1. **Hierarchical clustering** of scene graph
2. Reconstruct clusters independently (using incremental or global SfM)
3. **Merge clusters** using similarity transformations

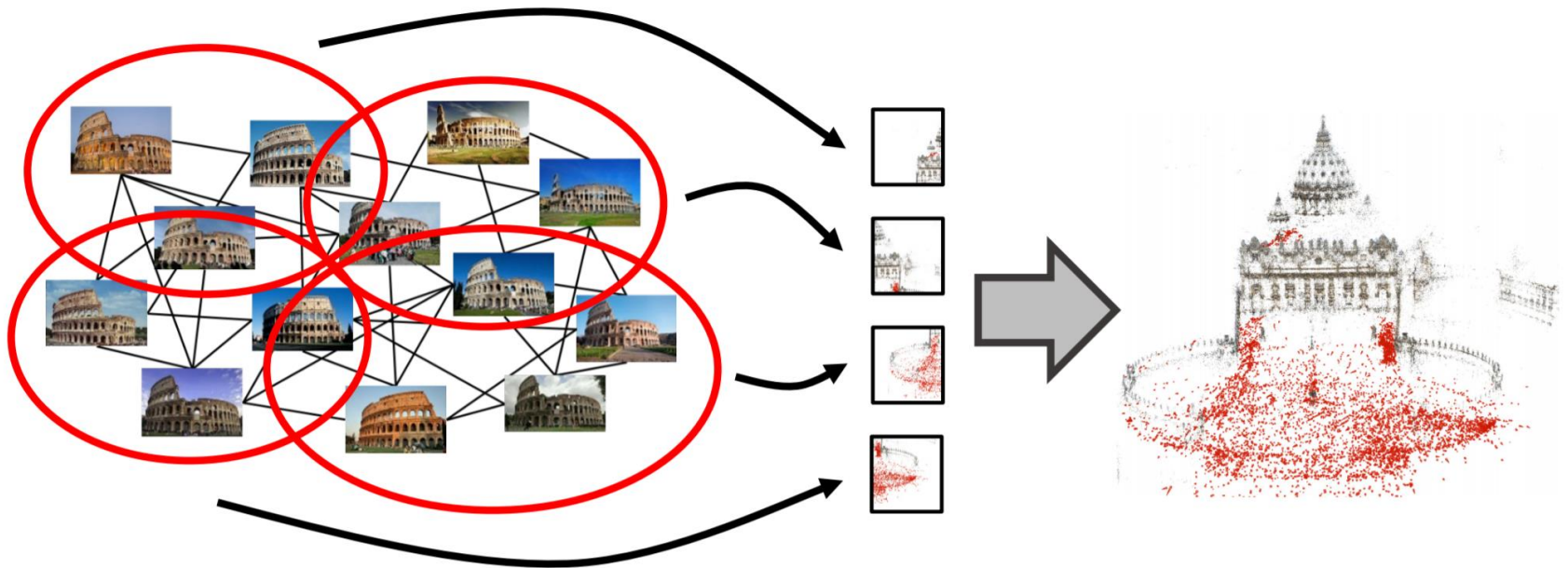
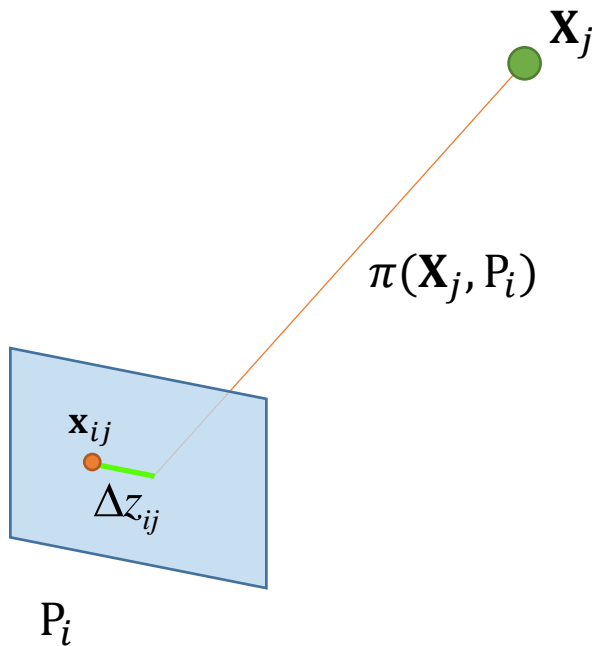


Image source: <https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf>

# Bundle Adjustment

- Refine a visual reconstruction to produce **jointly optimal** 3D structures  $\mathbf{X}_j$  and camera poses  $P_i$ .
- Minimize the total **reprojection errors**  $\Delta z_{ij}$ .



**Cost Function:**

$$\operatorname{argmin}_{\mathbf{P}, \mathbf{X}} \sum_i \sum_j \underbrace{\|\mathbf{x}_{ij} - \pi(P_i, \mathbf{X}_j)\|}_{\Delta z_{ij}}^2_{W_{ij}}$$

$W_{ij}$ : Measurement error covariance

$$\pi(\mathbf{X}_j, P_i) = \frac{P_i \mathbf{X}_j}{\hat{x}_3}, \text{ where } \hat{x}_3 \text{ is the 3rd element of } P_i \mathbf{X}_j$$

# Bundle Adjustment

- The **cost function** can be rewritten as:

$$\underset{\mathcal{P}}{\operatorname{argmin}} \underbrace{\sum_i \sum_j \Delta \mathbf{z}_{ij}^\top \mathbf{W}_{ij} \Delta \mathbf{z}_{ij}}_{g(\mathcal{P})}$$

- $\mathbf{W}_{ij}$ : Measurement error covariance
- $\mathcal{P}$  is a  $(12M \times 3N)$  vector with  **$12M$  camera matrix parameters** and  **$3N$  point coordinates**

# Iterative Estimation Methods

- Suppose we are given a hypothesized **nonlinear functional relation**  $\mathbf{X} = f(\mathbf{P})$ .
- $\mathbf{X} \in \mathbb{R}^N$  is a **measurement vector** and  $\mathbf{P} \in \mathbb{R}^M$  is a **parameter vector** in Euclidean spaces.
- A measured value of  $\mathbf{X}$  approximating the true value  $\bar{\mathbf{X}}$  is provided, and we wish to find the vector  $\hat{\mathbf{P}}$  that **most nearly satisfies** this functional relation.
- More precisely, we seek the vector  $\hat{\mathbf{P}}$  satisfying  $\mathbf{X} = f(\hat{\mathbf{P}}) - \epsilon$  for which  $g(\hat{\mathbf{P}}) = \frac{1}{2} \|\epsilon\|^2 = \frac{1}{2} \epsilon^T \epsilon$  **is minimized**.



# Iterative Estimation Methods

- We will look at **four iterative estimation methods** to minimize  $g(\mathbf{P})$ .
  1. Newton's Method
  2. Gauss-Newton
  3. Gradient Descent
  4. Levenberg-Marquardt

# Newton's Method and the Hessian

- Let  $\mathbf{P}_0$  be an **initial estimated value**, i.e.  $g(\mathbf{P}_0)$  .
- We may expand  $g(\mathbf{P})$  about  $\mathbf{P}_0$  in a **Taylor series** to get:

$$g(\mathbf{P}_0 + \Delta) = g + g_{\mathbf{P}}\Delta + \Delta^T g_{\mathbf{PP}}\Delta/2 + \dots$$

- where subscript  $\mathbf{P}$  denotes **differentiation**, and the right-hand side is evaluated at  $\mathbf{P}_0$ , i.e.

$$g_{\mathbf{P}} = \left. \frac{\partial g}{\partial \mathbf{P}} \right|_{\mathbf{p}=\mathbf{P}_0} \quad \text{and} \quad g_{\mathbf{PP}} = \left. \frac{\partial^2 g}{\partial \mathbf{P}^2} \right|_{\mathbf{p}=\mathbf{P}_0} .$$

# Jacobian and Hessian Matrices

- Suppose  $\mathbf{f}(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a function such that each of its 1<sup>st</sup> and 2<sup>nd</sup> order partial derivatives exist on  $\mathbb{R}^n$ , then the **Jacobian matrix** is given by an  $m \times n$  matrix:

$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \quad \text{or simply} \quad \mathbf{J}_{ij} = \frac{\partial f_i}{\partial x_j},$$

- And the **Hessian matrix** is an  $n \times n$  matrix:

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \quad \text{or simply} \quad \mathbf{H}_{i,j} = \frac{\partial^2 f}{\partial x_i \partial x_j}.$$

[https://en.wikipedia.org/wiki/Jacobian\\_matrix\\_and\\_determinant](https://en.wikipedia.org/wiki/Jacobian_matrix_and_determinant)  
[https://en.wikipedia.org/wiki/Hessian\\_matrix](https://en.wikipedia.org/wiki/Hessian_matrix)

# Newton's Method and the Hessian

- We seek a point  $g(\mathbf{P}_1)$ , with  $\mathbf{P}_1 = \mathbf{P}_0 + \Delta$ , which **minimizes**:

$$\begin{aligned} g(\mathbf{P}_1) &= g(\mathbf{P}_0) + g_{\mathbf{p}}\Delta + \Delta^T g_{\mathbf{PP}}\Delta/2 \\ &= g_0 + g_{\mathbf{p}}\Delta + \Delta^T g_{\mathbf{PP}}\Delta/2 \end{aligned}$$

with respect to  $\Delta$ .

- To this end, we differentiate with respect to  $\Delta$  and set the derivative to zero to get:

$$g_{\mathbf{p}} + g_{\mathbf{PP}}\Delta = 0 \quad \text{or} \quad g_{\mathbf{PP}}\Delta = -g_{\mathbf{p}}.$$

# Newton's Method and the Hessian

- The solution vector  $\hat{\mathbf{P}}$  is obtained by starting with an estimate  $\mathbf{P}_0$  and **computing successive approximations** according to the formula:

$$\mathbf{P}_{i+1} = \mathbf{P}_i + \Delta_i$$

- where  $\Delta_i$  is the solution to the **linear least-squares** problem

$$g_{\mathbf{P}\mathbf{P}}\Delta = -g_{\mathbf{P}}.$$

# Newton's Method and the Hessian

## Remarks:

- Newton iteration is based on the assumption of an **approximately quadratic cost** function near the minimum.
- And will show **rapid convergence** if this condition is met.
- The disadvantage of this approach is that the computation of the Hessian **may be difficult**.
- In addition, the assumption of quadratic behaviour is probably **invalid** when far from the minimum.

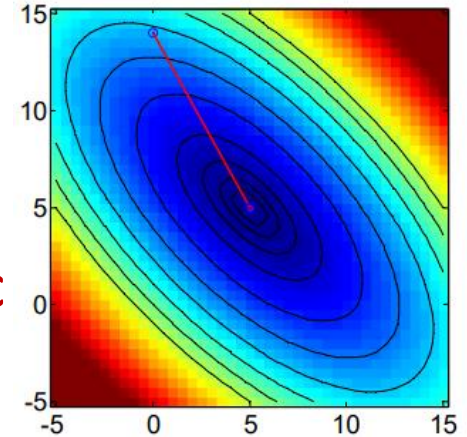


Image source: <http://www.robots.ox.ac.uk/~az/lectures/opt/lect1.pdf>

# Gauss-Newton Method

- Since  $g(\mathbf{P}) = \frac{1}{2} \|\boldsymbol{\epsilon}(\mathbf{P})\|^2 = \boldsymbol{\epsilon}(\mathbf{P})^T \boldsymbol{\epsilon}(\mathbf{P}) / 2$ , we can write:

$$g_{\mathbf{P}} = \boldsymbol{\epsilon}_{\mathbf{P}}^T \boldsymbol{\epsilon}, \quad \text{where} \quad \boldsymbol{\epsilon}_{\mathbf{P}} = \mathbf{f}_{\mathbf{P}} = \mathbf{J}, \quad \text{and}$$

$$g_{\mathbf{PP}} = \boldsymbol{\epsilon}_{\mathbf{P}}^T \boldsymbol{\epsilon}_{\mathbf{P}} + \boldsymbol{\epsilon}_{\mathbf{PP}}^T \boldsymbol{\epsilon}.$$

- In Gauss-Newton method, we **ignore the**  $\boldsymbol{\epsilon}_{\mathbf{PP}}^T \boldsymbol{\epsilon}$  term to get:

$$g_{\mathbf{PP}} = \boldsymbol{\epsilon}_{\mathbf{P}}^T \boldsymbol{\epsilon}_{\mathbf{P}} = \mathbf{J}^T \mathbf{J}.$$

# Gauss-Newton Method

- The **update equation**  $g_{PP}\Delta = -g_P$  becomes:

$$J^T J \Delta = -J^T \epsilon.$$

- This is also known as the **normal equation**.
- **Weighted iteration:** assume that the measurement  $\mathbf{X}$  satisfies a Gaussian distribution with **covariance matrix**  $\Sigma_X$ , the normal equation becomes

$$J^T \Sigma^{-1} J \Delta_i = -J^T \Sigma^{-1} \epsilon_i.$$

- $\Sigma_X$  is a **symmetric** and **positive definite** matrix.



# Gradient Descent

- The negative (or down-hill) gradient vector  $-g_{\mathbf{p}} = -\epsilon_{\mathbf{p}}^T \epsilon$  defines the direction of **most rapid decrease** of the cost function.
- A strategy for minimization of  $g$  is to move iteratively in the gradient direction; this is known as **gradient descent**.
- The **parameter increment**  $\Delta$  is computed from the equation  $\lambda \Delta = -g_{\mathbf{p}}$ , where  $\lambda$  controls the length of the step.

# Gradient Descent

## Remarks:

- We may consider this as related to Newton iteration where the Hessian is approximated (somewhat arbitrarily) by the **scalar matrix**  $\lambda I$ .
- Gradient descent by itself is not a very good minimization strategy, typically characterized by **slow convergence** due to zig-zagging.

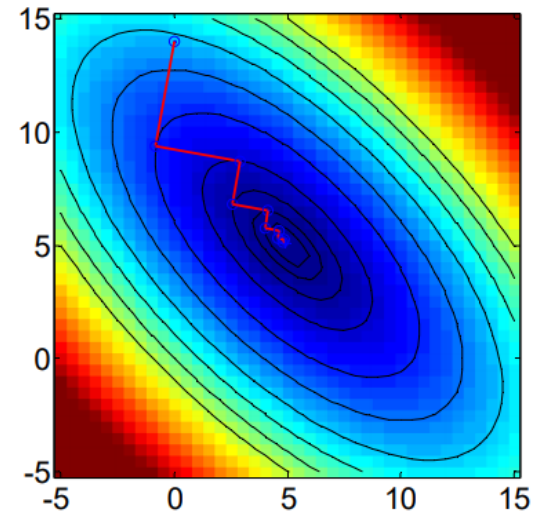


Image source: <http://www.robots.ox.ac.uk/~az/lectures/opt/lect1.pdf>

# Levenberg–Marquardt Method

- The Levenberg–Marquardt (abbreviated LM) iteration method is a **slight variation** on the Gauss-Newton iteration method.
- The normal equations  $J^T J \Delta = -J^T \epsilon$  are replaced by the **augmented normal equations**

$$(J^T J + \lambda I) \Delta = -J^T \epsilon,$$

- For some value of  $\lambda$  that **varies** from iteration to iteration, and  $I$  is the identity matrix.

# Levenberg–Marquardt Method

- A typical **initial value** of  $\lambda$  is  $10^{-3}$  times the average of the diagonal elements of  $N = J^T J$ .
- If the value of  $\Delta$  obtained by solving the augmented normal equations leads to a:
  - i. **Reduction in the error**, then the increment is accepted and  $\lambda$  is divided by a factor (typically 10) before the next iteration.
  - ii. **Increased error**, then the increment is rejected;  $\lambda$  is multiplied by the same factor and the augmented normal equations are solved again.

# Justification of LM

- When  $\lambda$  is very small, the method is essentially the same as Gauss-Newton iteration, i.e.

$$(J^T J + \cancel{\lambda I}) \Delta = -J^T \epsilon$$



$$J^T J \Delta = -J^T \epsilon$$

- **Remarks:** The error function  $\|\epsilon\|^2 = f(\mathbf{P}) - \mathbf{X}^2$  is quadratic in  $\mathbf{P}$  when near the minimum, hence will converge fast to the minimum value.

# Justification of LM

- When  $\lambda$  is large, the normal equation matrix is approximated by  $\lambda I$ , and the normal equations become  $\lambda \Delta = -J^T \epsilon$ .
- Recalling that  $J^T \epsilon$  is simply the gradient vector of  $\|\epsilon\|^2$ , we see that the direction of the parameter increment  $\Delta$  approaches that given by gradient descent.
- **Remarks:** Gauss-Newton which does not work well due to the violation of the quadratic assumption when far from the minimum is replaced by Gradient descent.

# Justification of LM

- The LM algorithm moves seamlessly between:
  1. **Gauss-Newton iteration**, which will cause rapid convergence in the neighbourhood of the solution, and
  2. A **gradient descent approach**, which will guarantee a decrease in the cost function when the going is difficult.
- Indeed, as  $\lambda$  becomes **increasingly larger**, the length of the increment step  $\Delta$  **decreases** and eventually it will lead to a **decrease** of the cost function  $\|\epsilon\|^2$ .

# A Sparse Levenberg–Marquardt Algorithm

- The bare LM algorithm is **not very suitable** for minimizing cost functions w.r.t. **large numbers of parameters**, e.g. bundle adjustment.
- This is because the solution of the normal equations has **complexity  $\mathcal{O}(N^3)$**  in the number of parameters, and this step is repeated many times.
- However, the normal equation matrix has a certain **sparse block structure** that we may take advantage of to realize very great time savings.



# A Sparse Levenberg–Marquardt Algorithm

- More specifically, the set of parameters  $\mathbf{P} \in \mathbb{R}^M$  may be **partitioned into** parameter vectors  $\mathbf{a}$  and  $\mathbf{b}$  so that  $\mathbf{P} = (\mathbf{a}^\top, \mathbf{b}^\top)^\top$ .
- Now the Jacobian matrix  $\mathbf{J} = [\partial \hat{\mathbf{X}} / \partial \mathbf{P}]$  has a block structure of the form  $\mathbf{J} = [\mathbf{A} \mid \mathbf{B}]$ , where **Jacobian submatrices** are defined by

$$\mathbf{A} = \left[ \partial \hat{\mathbf{X}} / \partial \mathbf{a} \right] \quad \text{and} \quad \mathbf{B} = \left[ \partial \hat{\mathbf{X}} / \partial \mathbf{b} \right].$$

# A Sparse Levenberg–Marquardt Algorithm

- The **normal equation** to be solved at each step of the LM algorithm now becomes:

$$\left[ \begin{array}{c|c} A^T \Sigma_X^{-1} A & A^T \Sigma_X^{-1} B \\ \hline B^T \Sigma_X^{-1} A & B^T \Sigma_X^{-1} B \end{array} \right] \begin{pmatrix} \delta_a \\ \delta_b \end{pmatrix} = \begin{pmatrix} A^T \Sigma_X^{-1} \epsilon \\ B^T \Sigma_X^{-1} \epsilon \end{pmatrix}.$$

- And may now be written in **block form** as:

$$\begin{bmatrix} U^* & W \\ W^T & V^* \end{bmatrix} \begin{pmatrix} \delta_a \\ \delta_b \end{pmatrix} = \begin{pmatrix} \epsilon_A \\ \epsilon_B \end{pmatrix}.$$

# A Sparse Levenberg–Marquardt Algorithm

- As a first step to solving these equations, both sides are now multiplied on the left by

$$\begin{bmatrix} \mathbf{I} & -\mathbf{W}\mathbf{V}^{*-1} \\ 0 & \mathbf{I} \end{bmatrix}$$

resulting in

$$\begin{bmatrix} \mathbf{U}^* - \mathbf{W}\mathbf{V}^{*-1}\mathbf{W}^T & 0 \\ \mathbf{W}^T & \mathbf{V}^* \end{bmatrix} \begin{pmatrix} \delta_a \\ \delta_b \end{pmatrix} = \begin{pmatrix} \epsilon_A - \mathbf{W}\mathbf{V}^{*-1}\epsilon_B \\ \epsilon_B \end{pmatrix}.$$

- This results in the **elimination of** the top right hand block.

# A Sparse Levenberg–Marquardt Algorithm

- The top half of this set of equations is

$$(U^* - WV^{*-1}W^T)\delta_a = \epsilon_A - WV^{*-1}\epsilon_B.$$

- This is also known as the **Schur's complement**; and these equations may be solved to find  $\delta_a$ .
- Subsequently, the value of  $\delta_b$  may be found by **back-substitution**, giving

$$V^*\delta_b = W^T\delta_a - \epsilon_B$$

# A Sparse Levenberg–Marquardt Algorithm

- If the newly computed value of the parameter vector  $\mathbf{P} = ((\mathbf{a} + \delta\mathbf{a})^\top, (\mathbf{b} + \delta\mathbf{b})^\top)^\top$  results in a **diminished value** of the error function.
- Then **accept** the new parameter vector  $\mathbf{P}$ , diminishes the value of  $\lambda$  by a factor of 10, and proceeds to the next iteration.
- On the other hand, if the **error value is increased**, then **reject** the new  $\mathbf{P}$  and tries again with a new value of  $\lambda$ , increased by a factor of 10.

# A Sparse Levenberg–Marquardt Algorithm

Given A vector of measurements  $\mathbf{X}$  with covariance matrix  $\Sigma_{\mathbf{X}}$ , an initial estimate of a set of parameters  $\mathbf{P} = (\mathbf{a}^T, \mathbf{b}^T)^T$  and a function  $f : \mathbf{P} \mapsto \hat{\mathbf{X}}$  taking the parameter vector to an estimate of the measurement vector.

Objective Find the set of parameters  $\mathbf{P}$  that minimizes  $\epsilon^T \Sigma_{\mathbf{X}}^{-1} \epsilon$  where  $\epsilon = \mathbf{X} - \hat{\mathbf{X}}$ .

Algorithm

- (i) Initialize a constant  $\lambda = 0.001$  (typical value).
- (ii) Compute the derivative matrices  $\mathbf{A} = [\partial \hat{\mathbf{X}} / \partial \mathbf{a}]$  and  $\mathbf{B} = [\partial \hat{\mathbf{X}} / \partial \mathbf{b}]$  and the error vector  $\epsilon$ .
- (iii) Compute intermediate expressions

$$\mathbf{U} = \mathbf{A}^T \Sigma_{\mathbf{X}}^{-1} \mathbf{A} \quad \mathbf{V} = \mathbf{B}^T \Sigma_{\mathbf{X}}^{-1} \mathbf{B} \quad \mathbf{W} = \mathbf{A}^T \Sigma_{\mathbf{X}}^{-1} \mathbf{B}$$

$$\epsilon_{\mathbf{A}} = \mathbf{A}^T \Sigma_{\mathbf{X}}^{-1} \epsilon \quad \epsilon_{\mathbf{B}} = \mathbf{B}^T \Sigma_{\mathbf{X}}^{-1} \epsilon$$

- (iv) Augment  $\mathbf{U}$  and  $\mathbf{V}$  by multiplying their diagonal elements by  $1 + \lambda$ .
- (v) Compute the inverse  $\mathbf{V}^{*-1}$ , and define  $\mathbf{Y} = \mathbf{W} \mathbf{V}^{*-1}$ . The inverse may overwrite the value of  $\mathbf{V}^*$  which will not be needed again.
- (vi) Find  $\delta_{\mathbf{a}}$  by solving  $(\mathbf{U}^* - \mathbf{Y} \mathbf{W}^T) \delta_{\mathbf{a}} = \epsilon_{\mathbf{A}} - \mathbf{Y} \epsilon_{\mathbf{B}}$ .
- (vii) Find  $\delta_{\mathbf{b}}$  by back-substitution:  $\delta_{\mathbf{b}} = \mathbf{V}^{*-1} (\epsilon_{\mathbf{B}} - \mathbf{W}^T \delta_{\mathbf{a}})$ .
- (viii) Update the parameter vector by adding the incremental vector  $(\delta_{\mathbf{a}}^T, \delta_{\mathbf{b}}^T)^T$  and compute the new error vector.
- (ix) If the new error is less than the old error, then accept the new values of the parameters, diminish the value of  $\lambda$  by a factor of 10, and start again at step (ii), or else terminate.
- (x) If the new error is greater than the old error, then revert to the old parameter values, increase the value of  $\lambda$  by a factor of 10, and try again from step (iv).

Table Source: Page 605, R. Hartley and A. Zisserman, "Multiple View Geometry in Computer Vision"

# Sparse LM on Multiple Image Bundle Adjustment

- We take advantage of the **lack of interaction** between parameters of the different cameras.

- **Measurement Data:**

$\mathbf{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_i\}$ , where  $\mathbf{X}_i = (\mathbf{x}_{i1}^\top, \mathbf{x}_{i2}^\top, \dots, \mathbf{x}_{ij}^\top)^\top$  and  $\mathbf{x}_{ij}$  is the **image** of the  $i$ -th 3D point in the  $j$ -th image .

- **The parameter vector:**

$\mathbf{a} = (\mathbf{a}_1^\top, \mathbf{a}_2^\top, \dots, \mathbf{a}_j^\top)^\top$ , where  $\mathbf{a}_j$  are the parameters of **the  $j$ -th camera**.

$\mathbf{b} = (\mathbf{b}_1^\top, \mathbf{b}_2^\top, \dots, \mathbf{b}_i^\top)^\top$ , where  $\mathbf{b}_i$  are the parameters of **the  $i$ -th 3D point**.

# Sparse LM on Multiple Image Bundle Adjustment

- Since the image point  $\mathbf{x}_{ij}$  does not depend on the parameters of any but the  **$j$ -th camera**,

$$\partial \hat{\mathbf{x}}_{ij} / \partial \mathbf{a}_k = 0 \quad \text{unless} \quad j = k.$$

- In a similar way for derivatives with respect to the parameters  $\mathbf{b}_k$  of the  **$k$ -th 3D point**,

$$\partial \hat{\mathbf{x}}_{ij} / \partial \mathbf{b}_k = 0 \quad \text{unless} \quad i = k.$$



# Sparse LM on Multiple Image Bundle Adjustment

- The form of the **Jacobian matrix**  $J$  for this problem and the resulting **normal equations**  $J^T J \delta = -J^T \epsilon$  are shown in the figures.

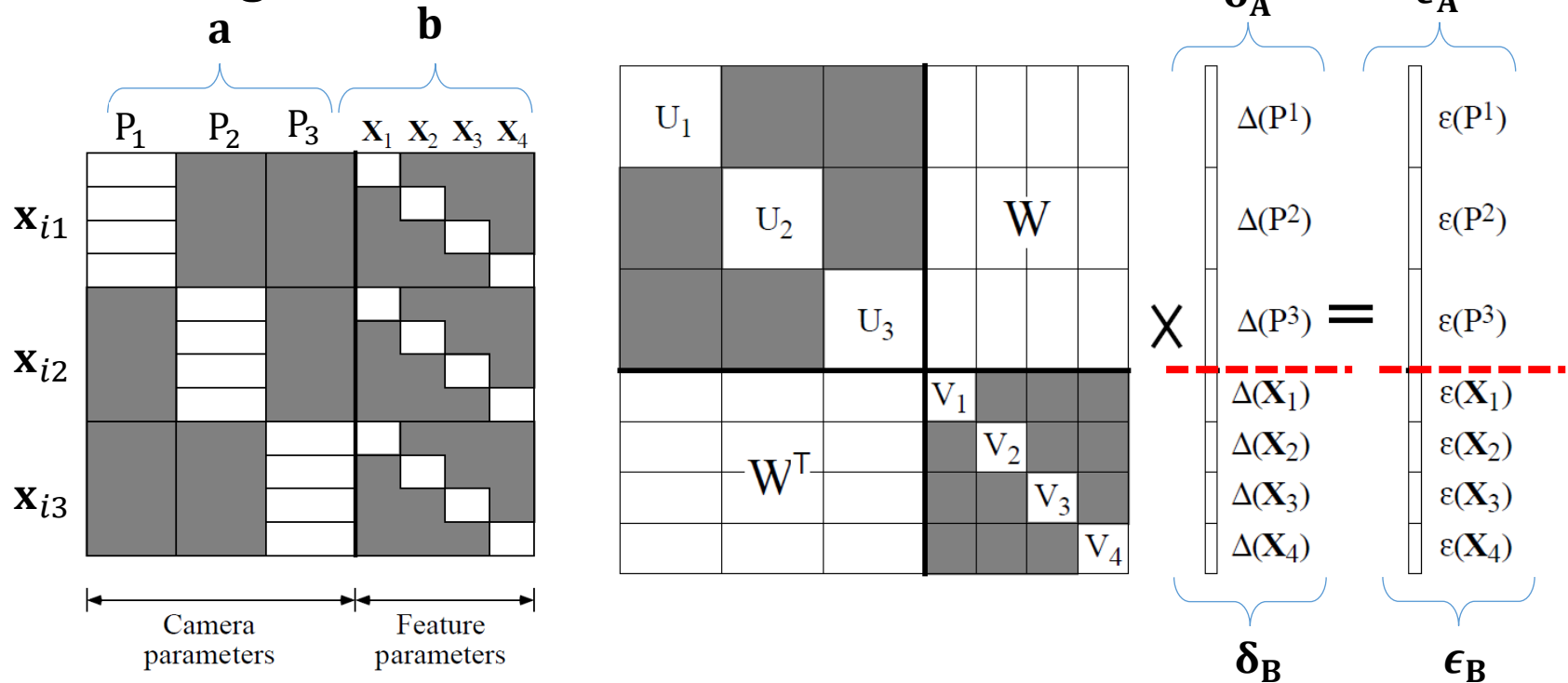


Image Source: R. Hartley and A. Zisserman, "Multiple View Geometry in Computer Vision"

# Sparse LM on Multiple Image Bundle Adjustment

- We can see that  $A_i = [\partial \hat{\mathbf{x}}_i / \partial \mathbf{a}]$  is a **block diagonal** matrix

$$A_i = \text{diag}(A_{i1}, \dots, A_{im}),$$

where  $A_{ij} = \partial \hat{\mathbf{x}}_{ij} / \partial \mathbf{a}_j$ .

**Example:**  $i^{\text{th}}$  3D Point in the  $2^{\text{nd}}$  camera, i.e.  $\mathbf{x}_{i2}$

$$\frac{\partial \mathbf{x}_{i2}}{\partial P_1} = 0, \quad \frac{\partial \mathbf{x}_{i2}}{\partial P_2} \neq 0, \quad \frac{\partial \mathbf{x}_{i2}}{\partial P_3} = 0$$

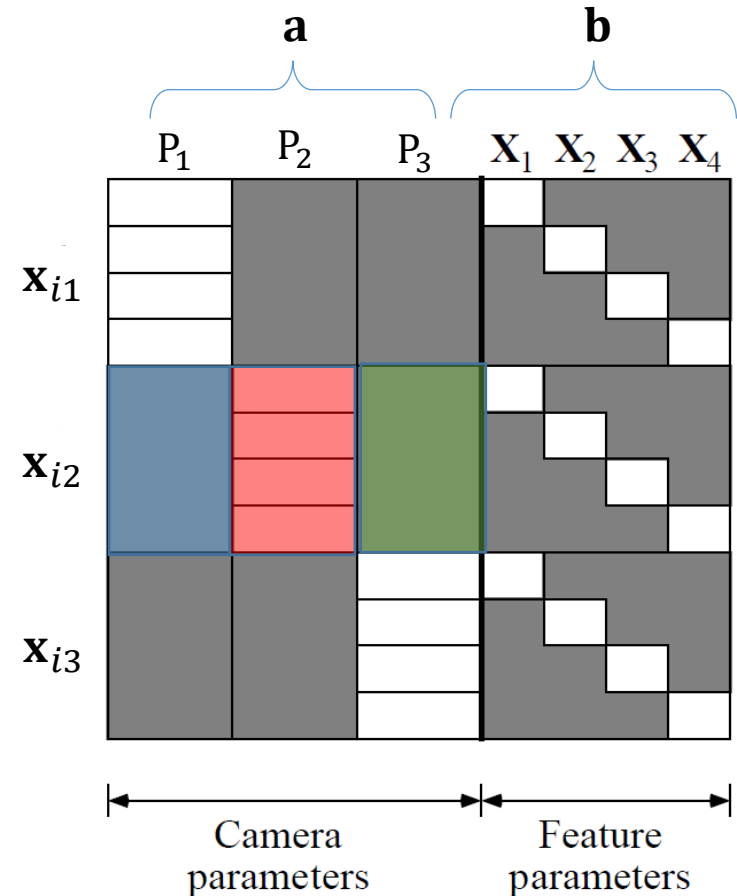


Image Source: R. Hartley and A. Zisserman, "Multiple View Geometry in Computer Vision"

# Sparse LM on Multiple Image Bundle Adjustment

- Similarly, matrix  $B_i = [\partial \hat{\mathbf{X}}_i / \partial \mathbf{b}_i]$  decomposes as:

$$B_i = [B_{i1}^T, \dots, B_{im}^T]^T,$$

where  $B_{ij} = \partial \hat{\mathbf{x}}_{ij} / \partial \mathbf{b}_i$ .

**Example:**  $i^{\text{th}}$  3D Point in the  $2^{\text{nd}}$  camera, i.e.  $\mathbf{x}_{i2}$

$$\frac{\partial \mathbf{x}_{12}}{\partial \mathbf{X}_1} \neq 0, \quad \frac{\partial \mathbf{x}_{22}}{\partial \mathbf{X}_2} \neq 0, \quad \frac{\partial \mathbf{x}_{32}}{\partial \mathbf{X}_3} \neq 0$$

$$\frac{\partial \mathbf{x}_{42}}{\partial \mathbf{X}_4} \neq 0, \quad \frac{\partial \mathbf{x}_{i2}}{\partial \mathbf{X}_k} = 0, k \neq i$$

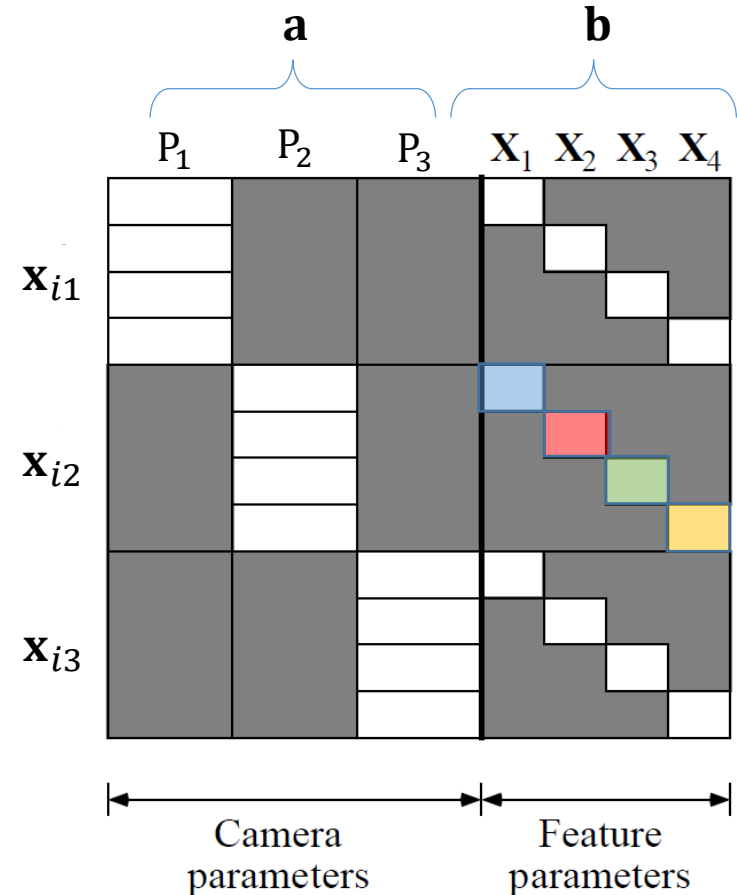
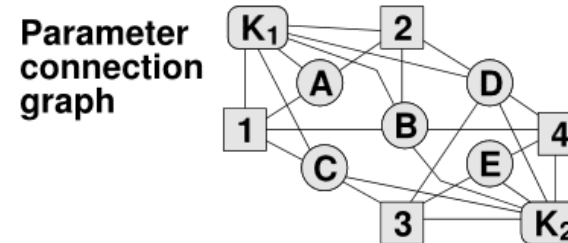
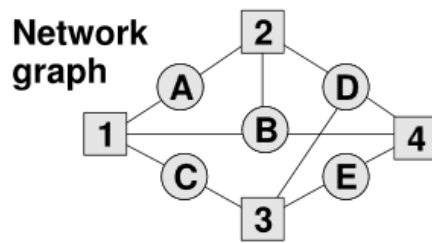


Image Source: R. Hartley and A. Zisserman, "Multiple View Geometry in Computer Vision"

# Hessian and Adjacency Relation

- The Hessian matrix  $H = J^T J$  is also the **adjacency matrix** of the graph that relates the parameters.



$J =$

	A	B	C	D	E	1	2	3	4	$K_1$	$K_2$
A1	■					■				■	
A2	■						■			■	
B1		■				■				■	
B2		■					■			■	
B4		■						■		■	
C1			■			■				■	
C3			■					■		■	
D2				■		■				■	
D3				■			■			■	
D4				■				■		■	
E3					■			■		■	
E4					■				■	■	

$H =$

	A	B	C	D	E	1	2	3	4	$K_1$	$K_2$
A	■					■				■	
B		■				■	■			■	
C			■			■		■		■	
D				■		■			■	■	
E					■	■			■	■	
1	■	■	■			■				■	
2	■	■		■			■			■	
3			■	■				■		■	
4				■	■				■	■	
$K_1$	■	■	■	■		■	■			■	
$K_2$		■	■	■	■			■	■	■	■

Image source: Bill Triggs et al, "Bundle Adjustment – A Modern Synthesis", 1999

# Solving the Normal Equation

- As mentioned earlier, we use the **Schur's complement** to solve the normal equations.

$$H_{LM}\delta = -J^T \Sigma \epsilon \quad \rightarrow \quad \left[ \begin{array}{c|c} A^T \Sigma_X^{-1} A & A^T \Sigma_X^{-1} B \\ \hline B^T \Sigma_X^{-1} A & B^T \Sigma_X^{-1} B \end{array} \right] \begin{pmatrix} \delta_a \\ \delta_b \end{pmatrix} = \begin{pmatrix} A^T \Sigma_X^{-1} \epsilon \\ B^T \Sigma_X^{-1} \epsilon \end{pmatrix}$$

The diagram illustrates the block structure of the normal equations. On the left, a block matrix  $H_{LM}$  is shown, partitioned into four main blocks:  $U$  (top-left),  $W$  (top-right),  $W^T$  (bottom-left), and  $V$  (bottom-right). The matrix is symmetric, with  $U$  and  $V$  being square and  $W$  and  $W^T$  being transposes of each other. The matrix is partitioned into four main blocks:  $U$  (top-left),  $W$  (top-right),  $W^T$  (bottom-left), and  $V$  (bottom-right). The matrix is symmetric, with  $U$  and  $V$  being square and  $W$  and  $W^T$  being transposes of each other. The matrix is partitioned into four main blocks:  $U$  (top-left),  $W$  (top-right),  $W^T$  (bottom-left), and  $V$  (bottom-right). The matrix is symmetric, with  $U$  and  $V$  being square and  $W$  and  $W^T$  being transposes of each other.

Next to the matrix is a vector of residuals  $J^T \Sigma \epsilon$ , partitioned into two parts:  $\Delta(P)$  (top) and  $\Delta(X)$  (bottom). The vector is partitioned into two parts:  $\Delta(P)$  (top) and  $\Delta(X)$  (bottom). The vector is partitioned into two parts:  $\Delta(P)$  (top) and  $\Delta(X)$  (bottom). The vector is partitioned into two parts:  $\Delta(P)$  (top) and  $\Delta(X)$  (bottom).

The equation is then written as:

$$\begin{bmatrix} U^* & W \\ W^T & V^* \end{bmatrix} \begin{pmatrix} \delta_a \\ \delta_b \end{pmatrix} = \begin{pmatrix} \epsilon_A \\ \epsilon_B \end{pmatrix}.$$

Image Source: R. Hartley and A. Zisserman, "Multiple View Geometry in Computer Vision"

# Solving the Normal Equation

- As mentioned earlier, we use the **Schur's complement** to solve the normal equations.

$$\begin{bmatrix} U^* & W \\ W^T & V^* \end{bmatrix} \begin{pmatrix} \delta_a \\ \delta_b \end{pmatrix} = \begin{pmatrix} \epsilon_A \\ \epsilon_B \end{pmatrix} \quad \begin{array}{l} \leftarrow \text{Camera Parameters} \\ \leftarrow \text{3D Structures} \end{array}$$

$$\Downarrow \quad \times \begin{bmatrix} I & -WV^{*-1} \\ 0 & I \end{bmatrix}$$

$$\begin{bmatrix} U^* - WV^{*-1}W^T & 0 \\ W^T & V^* \end{bmatrix} \begin{pmatrix} \delta_a \\ \delta_b \end{pmatrix} = \begin{pmatrix} \epsilon_A - WV^{*-1}\epsilon_B \\ \epsilon_B \end{pmatrix}.$$

# Solving the Normal Equation

- As mentioned earlier, we use the **Schur's complement** to solve the normal equations.

$$\begin{bmatrix} U^* - WV^{*-1}W^T & 0 \\ W^T & V^* \end{bmatrix} \begin{pmatrix} \delta_a \\ \delta_b \end{pmatrix} = \begin{pmatrix} \epsilon_A - WV^{*-1}\epsilon_B \\ \epsilon_B \end{pmatrix}$$

- First solve for  $\delta_a$  from:

$$\underbrace{(U^* - WV^{*-1}W^T)}_{\text{Schur Complement}} \delta_a = \epsilon_A - WV^{*-1}\epsilon_B$$

Easy to invert a block diagonal matrix

Schur Complement

(Sparse and Symmetric Positive Definite Matrix)

- Then Solve for  $\delta_b$  by backward substitution.

# Solving the Normal Equation

- The **non-homogeneous linear system** of equations can be solved without inverting  $A$  since it is a sparse matrix!

$$(U^* - WV^{*-1}W^T)\delta_a = \epsilon_A - WV^{*-1}\epsilon_B \quad \Rightarrow \quad A\mathbf{x} = \mathbf{b}$$

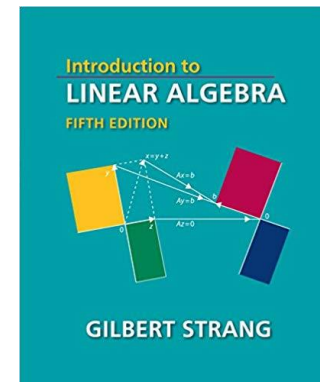
- Sparse matrix factorization

1. LU Factorization  $\longrightarrow A = LU$
2. QR factorization  $\longrightarrow A = QR$
3. Cholesky Factorization  $\longrightarrow A = LL^T$

Solve for  $\mathbf{x}$  by forward backward substitutions.

- Iterative methods

1. Conjugate gradient
2. Gauss-Seidel



Linear Algebra reference:

Gilbert Strang, "Introduction to Linear Algebra", Fifth Edition, 2016



# Problem of Fill-In

- Sparse matrix factorization can lead to the problem of fill-in, where the factorized matrix **becomes dense**.

## Example:

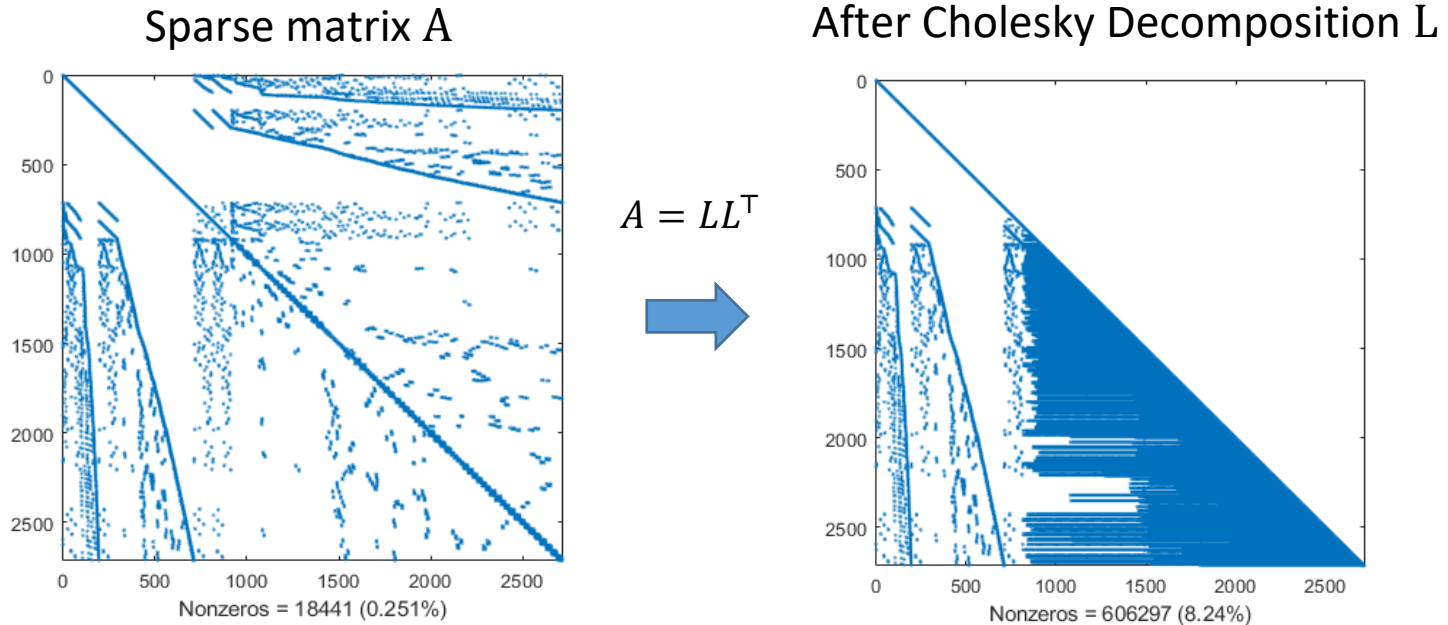


Image source: <https://www.mathworks.com/help/matlab/math/sparse-matrix-reordering.html>

# Problem of Fill-In

- **Reorder sparse matrix** to minimize fill-in.

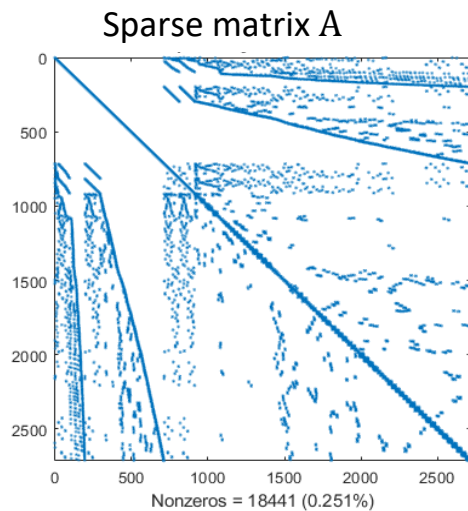
$$(P^T A P)(P^T \mathbf{x}) = P^T \mathbf{b}$$

Permutation matrix to reorder  $A$

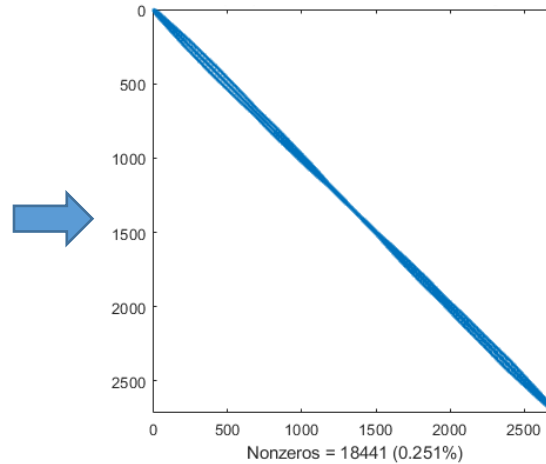
- NP-Complete problem.
- Approximate solutions:
  1. Minimum degree
  2. Column approximate minimum degree permutation
  3. Reverse Cuthill-McKee
  4. Nested Dissection ...

Reference: <https://www.mathworks.com/help/matlab/math/sparse-matrix-reordering.html>

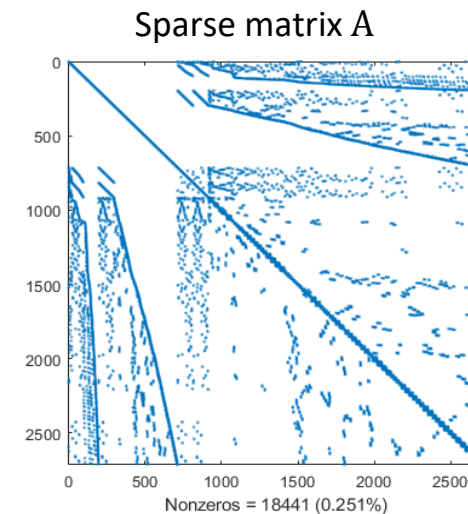
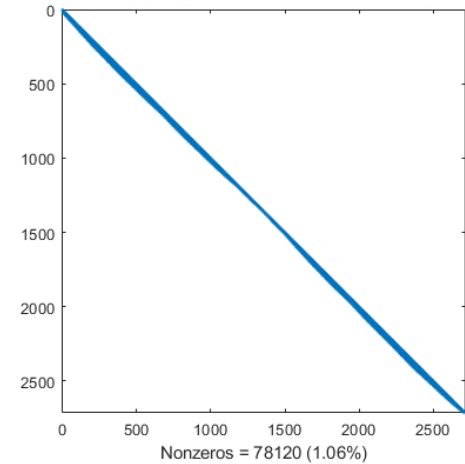
# Problem of Fill-In: Examples



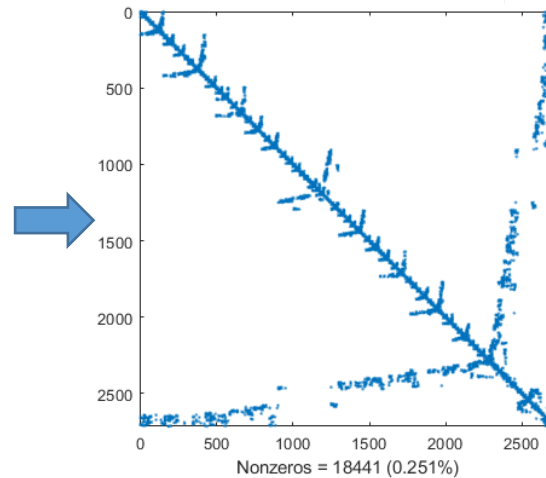
After **Cuthill-McKee** Ordering  $P^T A P$



After Cholesky Decomposition L



After **Minimum Degree** Ordering  $P^T A P$



After Cholesky Decomposition L

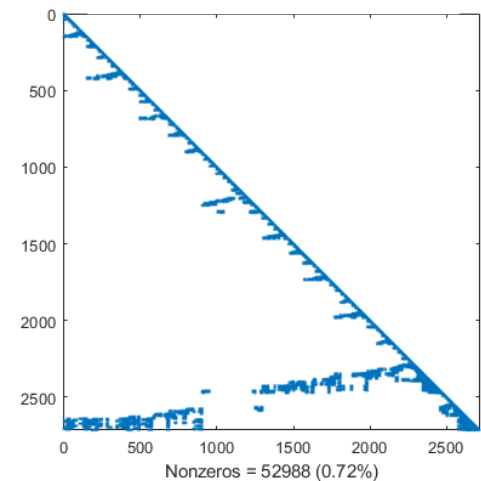
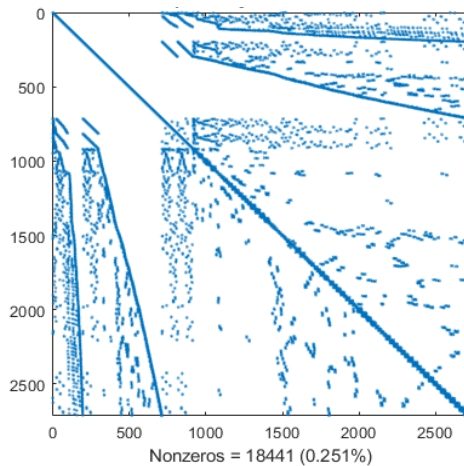


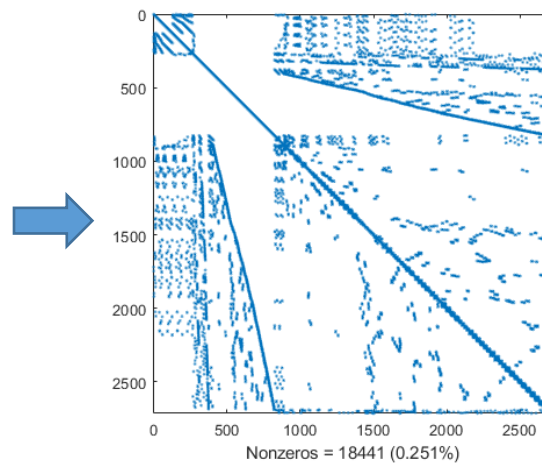
Image source: <https://www.mathworks.com/help/matlab/math/sparse-matrix-reordering.html>

# Problem of Fill-In: Examples

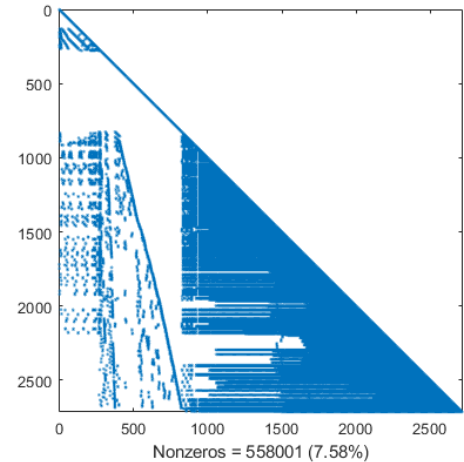
Sparse matrix A



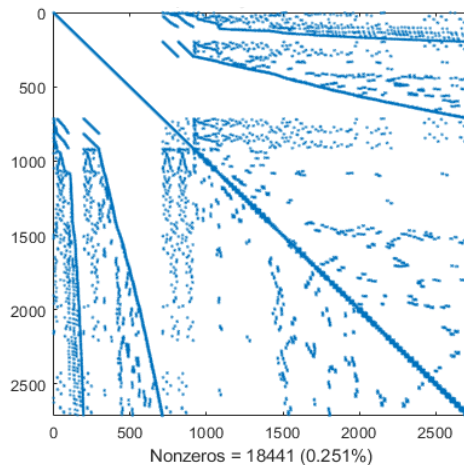
After **Column Count** Ordering  $P^T A P$



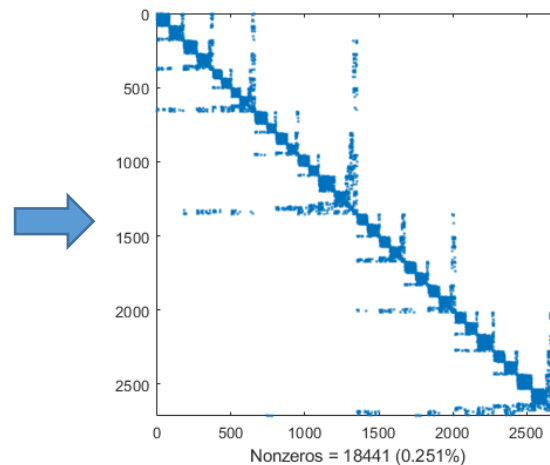
After Cholesky Decomposition L



Sparse matrix A



After **Nested Dissection** Ordering  $P^T A P$



After Cholesky Decomposition L

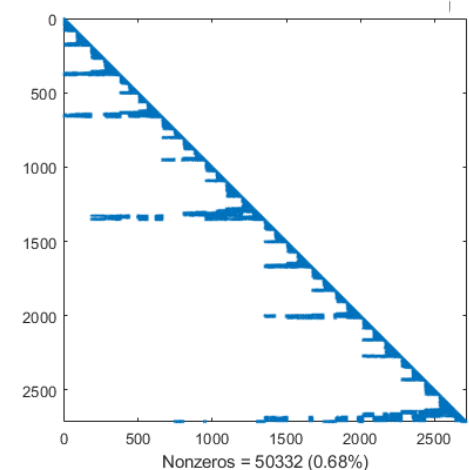


Image source: <https://www.mathworks.com/help/matlab/math/sparse-matrix-reordering.html>

# Problem of Fill-In: Examples

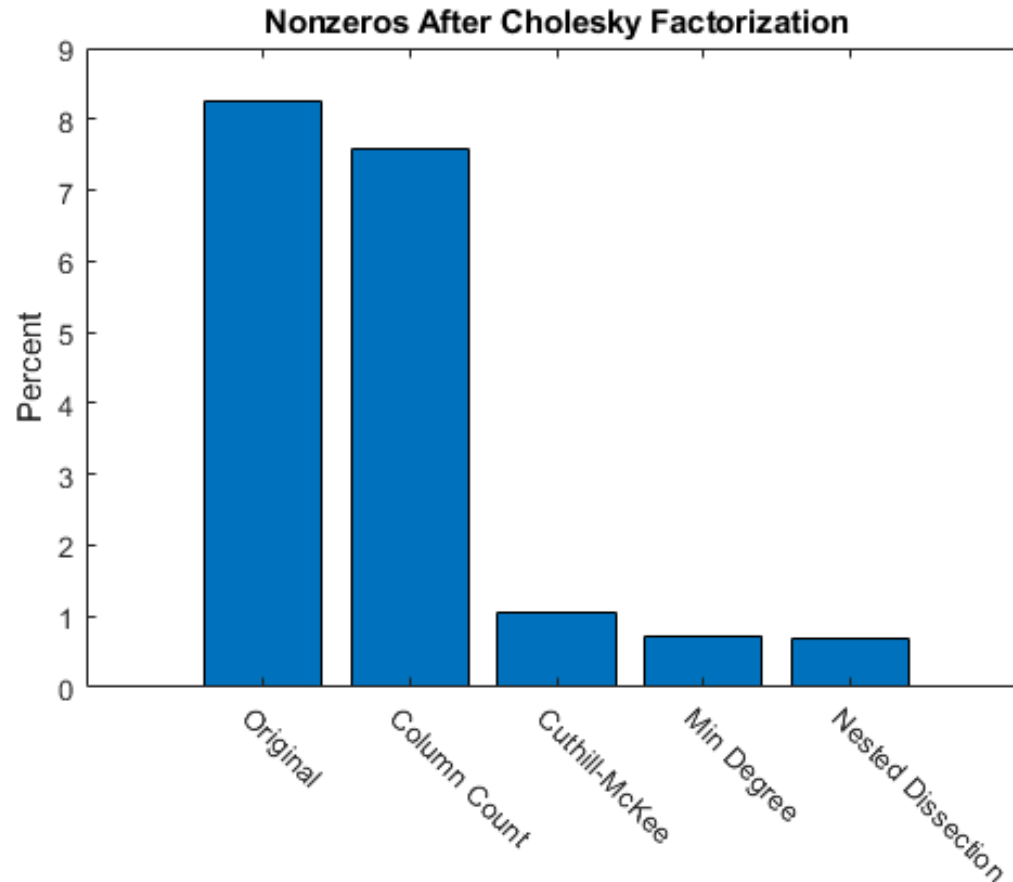


Image source: <https://www.mathworks.com/help/matlab/math/sparse-matrix-reordering.html>

# Open-Source Software

- Google Ceres: <http://ceres-solver.org/>
- g2o: <https://github.com/RainerKuemmerle/g2o>
- gtsam: <https://bitbucket.org/gtborg/gtsam/>
- Colmap: <https://colmap.github.io/>



Image source: <https://colmap.github.io/>

# Summary

- We have looked at how to:
  1. Describe the pipeline of large-scale **3D reconstruction**: data association, structure-from-motion and dense stereo.
  2. Explain the use of **robust two-view geometry** and the **bag-of-words** algorithm for data association.
  3. Use two-view geometry, PnP and triangulation to **initialize** the 3D reconstruction.
  4. Apply the iterative methods: Newton, Gauss-Newton Gradient descent or Levenberg-Marquardt for **bundle adjustment**.