# Bundle Adjustment

# Structure from Motion (SfM)

- Recovering the 3D geometry or "structure" of the scene and the camera motion from a set of 2D images when a camera is subject to "motion"
  - The structure refers to the 3D world coordinates of the captured points
  - The motion refers to the 3D world coordinates from which various images were captured from the camera(s)
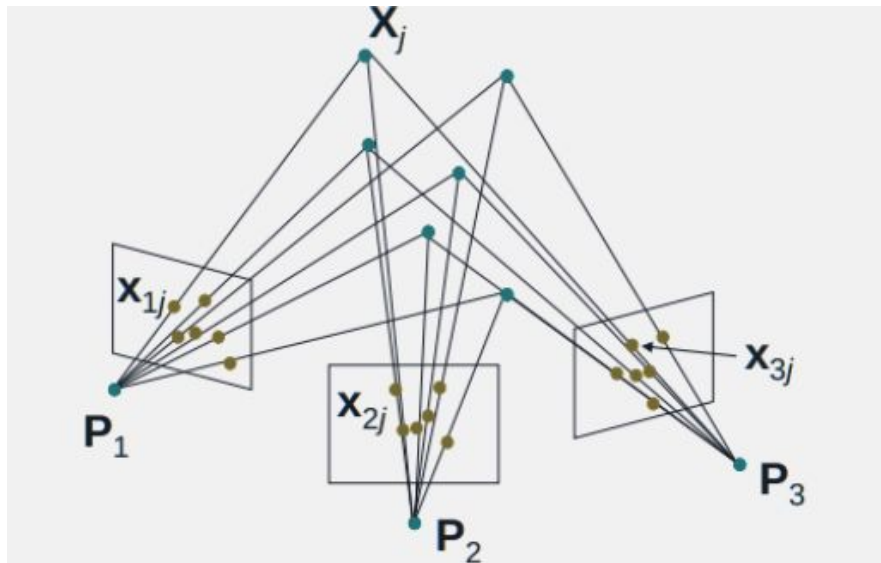
# Bundle Adjustment (BA)

- This is an optimisation algorithm to solve the problem of SfM
- The non-linear least squares constraint problem which we discussed previously is used to model the cost function in the Bundle Adjustment setting

# Initialisation for BA

- The following approaches help provide initialisation to BA: (brief overview)

  - Real world information like IMU, Odometry from which we obtain our initial guess for poses. Then by triangulation, obtain initial guess of 3D points as intial guess for BA

  - Fundamental Matrix estimation using the 8 point algorithm on the first pair of images to get the relative poses. Again, by triangulation we can obtain guesses for the 3D points from P3P

  - The relative poses and 3D point estimates can also be obtained by just pairwise fundamental matrix estimation

# Reprojection Error

- What we have:
  - m cameras
  - n points
  - Known correspondences across all these cameras and points
- Unknowns:
  - m matrices $P_i$
  - n coordinates $X_i$

- Known:
  - Data association, i.e., correspondences across images which relate the pixels that correspond to the same real world point taken from various cameras and hence, viewpoints
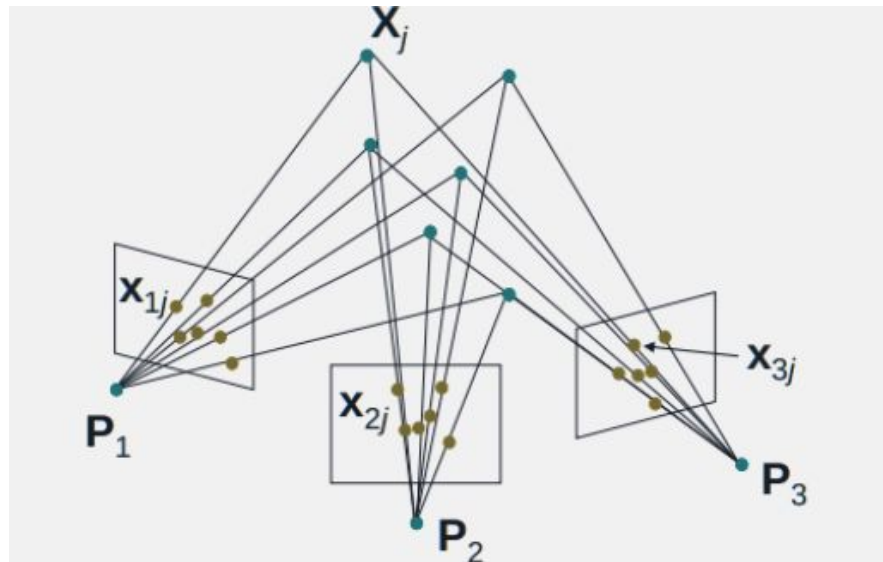
# Reprojection Error

- What we have:
  - m cameras
  - n points
  - Known correspondences across all these cameras and points
- The camera model expressed in homogeneous coordinates is given by the following set of equations



$$\vec{x}_{ij} = \lambda_{ij} P_i \vec{X}_j, \quad 1 \leq i \leq m, 1 \leq j \leq n$$

- And finally, the reprojection error is given by $\|\lambda_{ij} P_i \vec{X}_j - \vec{x}_{ij}\|^2$, which denotes how closely the estimated 3D point resembles its actual projection
  - The lambdas are scaling factors to ensure that the third coordinate in the image coordinates in homogeneous system is 1

# An Example

- Given that we have 10k images with 1k points each. Each point is seen 10 times on avg

- How many known and unknown parameters do we have?

- Known parameters: 20M
  - These are the image coordinates: 2 * 10k * 1k = 20M known params
  - Why are we multiplying by 2? Pixel coordinates

- Unknown parameters: ~13M
  - No of world points = (10k * 1k) / 10 = 1M. Division by 10 as the same point is repeated in 10 images on average. Hence, we have 1M * 3 = 3M world point params
    - Why are we multiplying by 3? World coordinates
  - For each of the 10k images, we have 10k corresponding orientations of the camera that captured the image. Hence, 10k * 6 = 60k orientation params
  - For each of the image points, there is a single corresponding scaling factor. Hence, the scaling factors contribute to 10k * 1k = 10M scaling params

# Can we do better?

- Get rid of the scale parameters

- How? Convert from homogeneous coordinate system to Euclidean coordinate system

# Cost Function

$$\arg \min_{\vec{X}_j, P_i} \sum_{i=1}^{M} \sum_{j=1}^{N} \left( \left[ \frac{P_{i11} X_j + P_{i12} Y_j + P_{i13} Z_j + P_{i14}}{P_{i31} X_j + P_{i32} Y_j + P_{i33} Z_j + P_{i34}} - x_{ij} \right]^2 \right.$$

$$\left. + \left[ \frac{P_{i21} X_j + P_{i22} Y_j + P_{i23} Z_j + P_{i24}}{P_{i31} X_j + P_{i32} Y_j + P_{i33} Z_j + P_{i34}} - y_{ij} \right]^2 \right)$$

$$\arg \min_{\vec{X}_j, P_i} \sum_{i=1}^{M} \sum_{j=1}^{N} \left\| \frac{P_{(1:2)i} \vec{X}_j}{P_{3i} \vec{X}_j} - \vec{x}_{ij} \right\|^2$$

# Cost Function

$$\arg\min_{\vec{X}_j, P_i} \sum_{i=1}^{M} \sum_{j=1}^{N} \left( \left[ \frac{P_{i11}X_j + P_{i12}Y_j + P_{i13}Z_j + P_{i14}}{P_{i31}X_j + P_{i32}Y_j + P_{i33}Z_j + P_{i34}} - x_{ij} \right]^2 + \left[ \frac{P_{i21}X_j + P_{i22}Y_j + P_{i23}Z_j + P_{i24}}{P_{i31}X_j + P_{i32}Y_j + P_{i33}Z_j + P_{i34}} - y_{ij} \right]^2 \right)$$

$$\arg\min_{\vec{X}_j, P_i} \sum_{i=1}^{M} \sum_{j=1}^{N} \left\| P_i \vec{X}_j - \vec{x}_{ij} \right\|^2$$

$$\arg\min_{\vec{X}_j, P_i} \sum_{i=1}^{M} \sum_{j=1}^{N} \left\| \frac{P_{(1:2)i} \vec{X}_j}{P_{3i} \vec{X}_j} - \vec{x}_{ij} \right\|^2$$

$$\arg\min_{X_j, P_i} \sum_{i=1}^{M} \sum_{j=1}^{N} \left\| \widehat{\vec{x}}_{ij} - \vec{x}_{ij} \right\|^2$$

# Reprojection Error

- $P_i$ is the projection matrix of the ith view and $X_\square$ is the world coordinate of the jth world point
- $x_i\square$ is the known observation while $\widehat{x}_{ij}$ is the predicted projection of the jth world point in the ith view

# What do we finally have?

- We have 2mn equations in total (2 for each match)

- This is a non-linear optimisation problem wherein we minimise the sum of the squared reprojection errors of the reconstructed n 3D points over m images

- How do we solve this familiar looking formulation of a non-linear least squares cost function? **Levenberg-Marquardt Algorithm**

$$(\mathbf{J}^\top \mathbf{J} + \lambda \mathbf{I}) \, \Delta \mathbf{k} = -\mathbf{J}^\top \mathbf{r}(\mathbf{k})$$

- Here r is the residual vector, which is the difference vector between the observed image coordinate vector and the predicted projection of the corresponding world point through the view
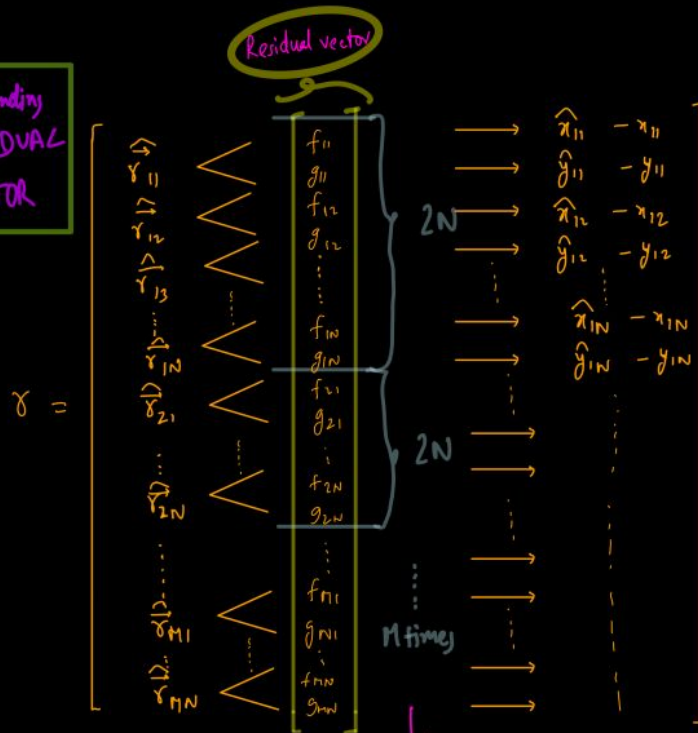
$$\widehat{\vec{x}}_{ij} = P_i \vec{X}_j$$

Objective Function:
$$F(k) = \gamma^T \gamma$$

$$\arg\min \sum_{i=1}^{M} \sum_{j=1}^{N} \left\| \widehat{\vec{x}}_{ij} - \vec{x}_{ij} \right\|^2$$

Residual vector

Understanding
RESIDUAL
VECTOR

$$\gamma = \begin{bmatrix} \overset{\rightharpoonup}{\gamma}_{11} \\ \overset{\rightharpoonup}{\gamma}_{12} \\ \widehat{\gamma}_{13} \\ \vdots \\ \widehat{\gamma}_{1N} \\ \widehat{\gamma}_{21} \\ \vdots \\ \widehat{\gamma}_{2N} \\ \vdots \\ \widehat{\gamma}_{M1} \\ \widehat{\gamma}_{MN} \end{bmatrix} \quad \begin{bmatrix} f_{11} \\ g_{11} \\ f_{12} \\ g_{12} \\ \vdots \\ f_{1N} \\ g_{1N} \\ f_{21} \\ g_{21} \\ \vdots \\ f_{2N} \\ g_{2N} \\ f_{M1} \\ g_{M1} \\ f_{MN} \\ g_{MN} \end{bmatrix} \begin{matrix} \\ \end{matrix}$$

$2N$

$2N$

M times

$$\begin{matrix} \widehat{x}_{11} - x_{11} \\ \widehat{y}_{11} - y_{11} \\ \widehat{x}_{12} - x_{12} \\ \widehat{y}_{12} - y_{12} \\ \vdots \\ \widehat{x}_{1N} - x_{1N} \\ \widehat{y}_{1N} - y_{1N} \end{matrix}$$

$$\left( \underline{Side\ note}:\ Note\ that\ \frac{\partial f_{11}}{\partial k} = \frac{\partial \widehat{x}_{11}}{\partial k} \right)$$

2MN terms

$$\widehat{x}_{ij} = P_i \vec{X}_j$$

Residual vector

Objective Function:
$$F(k) = \gamma^T \gamma$$

$$\arg\min \sum_{i=1}^{M} \sum_{j=1}^{N} \left\| \widehat{\vec{x}}_{ij} - \vec{x}_{ij} \right\|^2$$

$$J = \frac{\partial \vec{r}}{\partial k}$$

$$\begin{bmatrix} P_1 & P_2 & \ldots & P_M & \vec{X}_1 & \vec{X}_2 & \ldots & \vec{X}_N \end{bmatrix}$$

$P_{1_1} P_{1_2} P_{1_3} \ldots P_{1_{34}}$ ... $P_{M_{11}} P_{M_{12}} \ldots P_{M_{34}}$    $X_1 Y_1 Z_1$  ......  $X_N Y_N Z_N$

**INDIVIDUAL JACOBIANS :** $J_{ij_m}, G_{ij_m}, J_{ij_s}, G_{ij_s}$ $\Rightarrow$

'm' otion
$$J_{ij_m} = \begin{bmatrix} \dfrac{\partial f_{ij}}{\partial P_i} \end{bmatrix} = \begin{bmatrix} \dfrac{\partial f_{ij}}{\partial P_{i_{11}}} & \ldots\ldots & \dfrac{\partial f_{ij}}{\partial P_{i_{34}}} \end{bmatrix}$$

$$G_{ij_m} = \begin{bmatrix} \dfrac{\partial g_{ij}}{\partial P_i} \end{bmatrix} = \begin{bmatrix} \dfrac{\partial g_{ij}}{\partial P_{i_{11}}} & \ldots\ldots & \dfrac{\partial g_{ij}}{\partial P_{i_{34}}} \end{bmatrix}$$

's'tructure
$$J_{ij_s} = \begin{bmatrix} \dfrac{\partial f_{ij}}{\partial \vec{X}_j} \end{bmatrix} = \begin{bmatrix} \dfrac{\partial f_{ij}}{\partial X_j} & \dfrac{\partial f_{ij}}{\partial Y_j} & \dfrac{\partial f_{ij}}{\partial Z_j} \end{bmatrix}$$

$$G_{ij_s} = \begin{bmatrix} \dfrac{\partial g_{ij}}{\partial \vec{X}_j} \end{bmatrix} = \begin{bmatrix} \dfrac{\partial g_{ij}}{\partial X_j} & \dfrac{\partial g_{ij}}{\partial Y_j} & \dfrac{\partial g_{ij}}{\partial Z_j} \end{bmatrix}$$

So, from above, notice when we write $J_{12_m}$ simply (Notice 'm') means derivative w.r.t $P_1$. And $J_{12_s}$ would mean derivative w.r.t $\vec{X}_2$.

$$J_{2HN, 12M+3N} = \frac{\delta r}{\delta k}$$

MOTION    STRUCTURE

$$= \left[ \frac{\partial r}{\partial P_1} \quad \frac{\partial r}{\partial P_2} \quad \cdots\cdots \quad \frac{\partial r}{\partial P_M} \quad \middle| \quad \frac{\partial r}{\partial \vec{X}_1} \quad \frac{\partial r}{\partial \vec{X}_2} \quad \cdots\cdots \quad \frac{\partial r}{\partial \vec{X}_N} \right]$$

M times          N times

Update vector

$$\Delta k$$

$$\begin{bmatrix}
J_{11m} & O & O & \cdots & O & J_{11s} & O & \cdots & O \\
1\times12 & 1\times12 & 1\times12 & & 1\times12 & 1\times3 & 1\times3 & & 1\times3 \\
G_{11m} & O & & \cdots & O & G_{11s} & O & & \\
1\times12 & & & & & O & & & \\
J_{12m} & & & & O & O & J_{12s} & \cdots & O \\
G_{12m} & & & & & & G_{12s} & \cdots & O \\
\vdots & & & & O & O & & & \\
J_{1Nm} & & & & O & O & \cdots & & J_{1Ns} \\
G_{1Nm} & & & & O & O & & & G_{1Ns} \\
O & J_{21m} & \cdots & & O & J_{21s} & O & & O \\
O & G_{21m} & & \cdots & O & G_{21s} & O & & O \\
\vdots & & & & & & J_{22s} & \cdots & \\
O & G_{2Nm} & \cdots & & O & O & G_{22s} & \cdots & J_{2Ns} \\
& & & & & & & & G_{2Ns} \\
\vdots & \vdots & \vdots & & & & & & \\
O & O & O & J_{M1m} & J_{M1s} & O & & O \\
& & & G_{M1m} & G_{M1s} & O & & \\
& & & & O & & & \\
O & & & \vdots & O & & & \\
O & & & G_{MNm} & O & & & J_{MNs} \\
& & & & & & & G_{MNs}
\end{bmatrix}$$

2nd

2N

$(12M+3N, 12M+3N)$

$$\delta P_{11} \Big\} \ 1^{st} \text{ camera}$$
$$\delta P_{1\,34}$$
$$\delta P_{2\,11} \Big\} \ 2^{nd} \text{ camera}$$
$$\vdots$$
$$\delta P_{2\,34}$$
$$\vdots$$
$$\delta P_{M\,11} \Big\} \ M^{th} \text{ cam}$$
$$\vdots$$
$$\delta P_{M\,34}$$
$$\delta X_1 \atop \delta Y_1 \atop \delta Z_1 \Big\} \ 1^{st} \text{ 3D point}$$
$$\delta X_N \atop \delta Y_N \atop \delta Z_N \Big\} \ N^{th} \text{ 3D point}$$

$$\left[ J^T J + \lambda I \right] \quad \Delta k = - \quad J^T \quad r(k)$$

$$(12M+3N, 2HN) \times (2HN, 12M+3N) \qquad (12M+3N, 1) \qquad (12M+3N, 2HN) \qquad (2HN, 1)$$

# A Short Detour - Lie Algebra

- For Transformation Matrices, it's better to work in the Lie Space than in the R12 or R3 space due to the specific properties of Rotation Matrices, more details in the corresponding chapters in the pdfs provided in the folder. (Note: This is not compulsory, you can do it if you are interested - the derivatives given are mostly for ICP-SLAM and not SFM, might need to change it a little bit to adapt it for this task)

**A few resources:**

- ICP-SLAM-ShubhodhFin.pdf : More ICP SLAM than SFM, but a good read nonetheless to understand Lie Algebra
- https://karnikram.info/blog/lie/
- slambook-en.pdf : Chapter 3
- OptimizationOverManifolds.pdf: Page 39 onwards (Go in this particular order, would make things easier to understand)
- https://github.com/princeton-vl/lietorch

# Demos

1.  COLMAP
    https://github.com/colmap/colmap

2.  HLoc
    https://github.com/cvg/Hierarchical-Localization
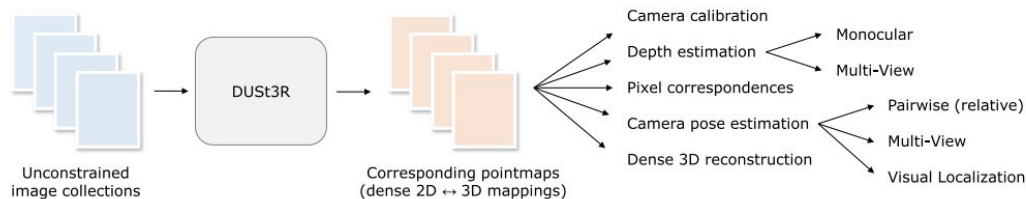
# What if we trained a Network for this task?

1. Dust3r
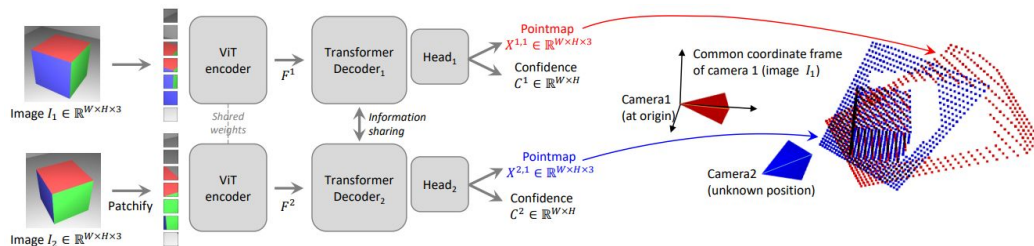   https://europe.naverlabs.com/research/publications/dust3r-geometric-3d-vision-made-easy/

2. Mast3r
   https://europe.naverlabs.com/research/publications/grounding-image-matching-in-3d-with-mast3r/

# Dust3r



Trained a transformer network to predict a joint pointmap given a pair of images.



The outputs are the 3D points in the first image's frame for each pixel, along with the confidence of the 3D point for each pixel

# Mast3r

1. Trains a new decoder Head to predict pixel correspondences.
2. Dust3r's decoder was originally scale invariant, it trains it with metric data now to predict metric depth, so that it can be used for pose estimation as well.
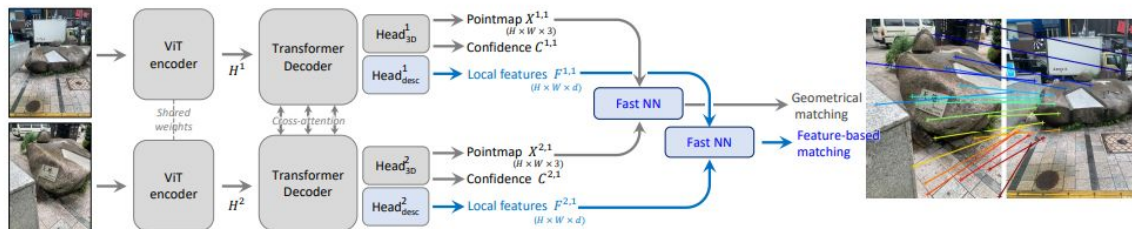


Figure 2: Overview of the proposed approach. Given two input images to match, our network regresses for each image and each input pixel a 3D point, a confidence value and a local feature. Plugging either 3D points or local features into our fast reciprocal NN matcher (3.3) yields robust correspondences. Compared to the DUSt3R framework which we build upon, our contributions are highlighted in blue.

# Questions?

# Acknowledgements

- Borrowed most of the slides from the SfM notion doc by Sai Shubodh

# References

- [Muti view geometry](#) by Hartley and Zisserman

- [RRC Summer School 2023](#)

- [Bundle Adjustment – Part I Introduction & Application](#) - Cyril Stachniss

- [Bundle Adjustment – Part II Numerics of BA](#) - Cyril Stachniss

- [Lecture: Photogrammetry I & II (2021, Uni Bonn, Cyrill Stachniss)](#) - Lectures 53, 54, 55

# Further Readings

- Incremental SfM, Failure cases of SfM, Multi View Stereo, and much more

  - [Structure from motion, L. Lazebnik, N. Snavely, M. Herbert](#)

  - [Structure from Motion, Derek Hoiem](#)

# Thank You!