



Practical Course: Vision-based Navigation Summer Term 2015

Lecture 3: Visual Motion Estimation, Direct Dense Methods

Dr. Jörg Stückler

What we will cover today

- Direct, dense motion estimation
 - Motion representation using the $SE(3)$ Lie algebra
 - Non-linear least squares optimization
 - Direct RGB-D odometry

Direct Visual Odometry with RGB-D Cameras

Robust Odometry Estimation for RGB-D Cameras

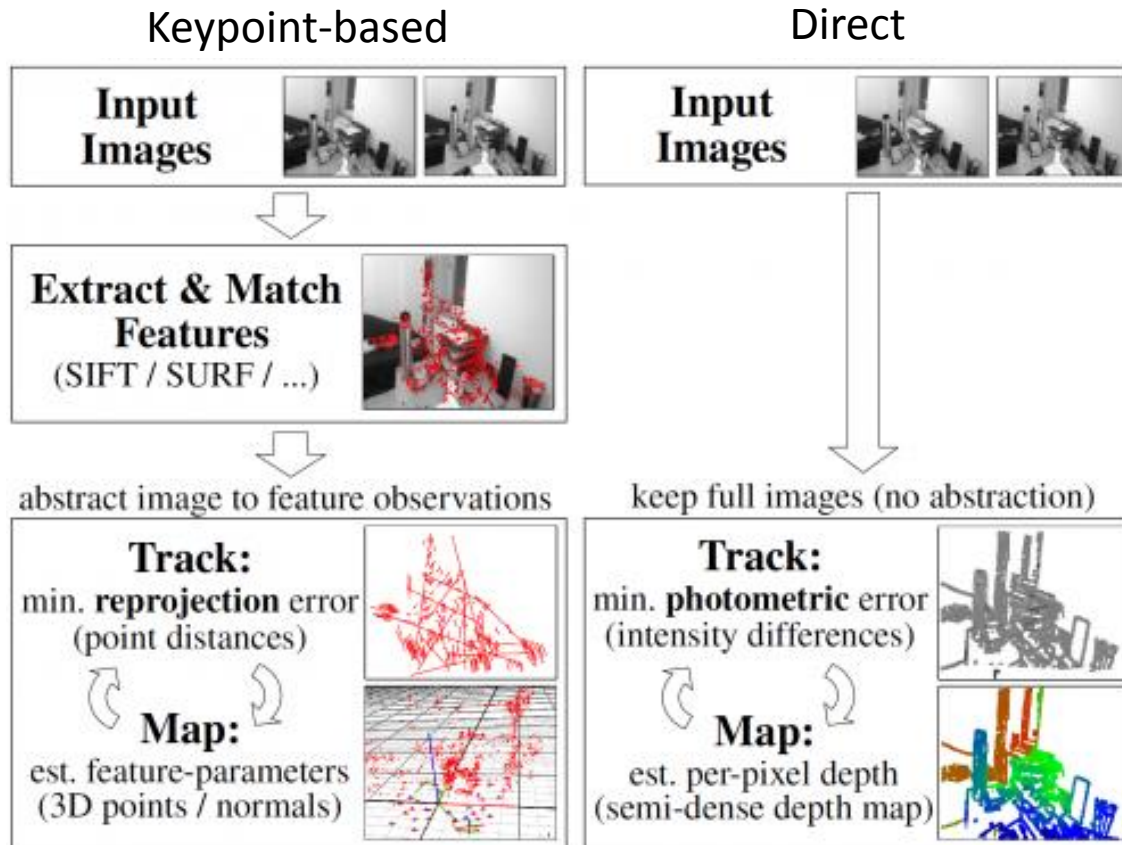
Christian Kerl, Jürgen Sturm, Daniel Cremers



Computer Vision and Pattern Recognition Group
Department of Computer Science
Technical University of Munich

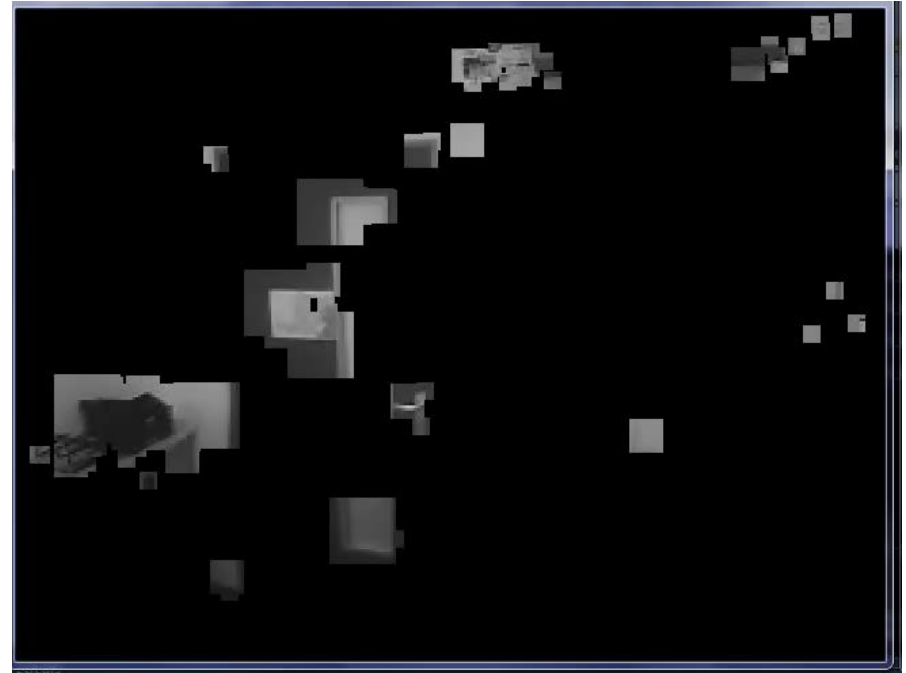


Keypoint-based vs. Direct VO Methods



- Sparse: use a small set of selected pixels (keypoints)
- Dense: use all (valid) pixels

Problem with Keypoint-based Methods



Special Euclidean Group SE(3)

- Not all matrices are transformation matrices: Transformation matrices have a special structure

$$\mathbf{T} = \begin{pmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbf{SE}(3) \subset \mathbb{R}^{4 \times 4}$$

- Translation \mathbf{t} has 3 degrees of freedom
- Rotation \mathbf{R} has 3 degrees of freedom
- They form a group which we call SE(3). The group operator is matrix multiplication:

$$\cdot : \mathbf{SE}(3) \times \mathbf{SE}(3) \rightarrow \mathbf{SE}(3)$$

$$\mathbf{T}_B^A \cdot \mathbf{T}_C^B \mapsto \mathbf{T}_C^A$$

- The operator is associative, but not commutative!
- There is also an inverse and a neutral element

Parametrizations of SE(3)

- Translation \mathbf{t} has 3 degrees of freedom
- Rotation \mathbf{R} has 3 degrees of freedom

$$\mathbf{T} = \begin{pmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{pmatrix} \in \mathbf{SE}(3) \subset \mathbb{R}^{4 \times 4}$$

- Different parametrizations θ of $\mathbf{T}(\theta)$
 - Direct matrix representation
 - Quaternion / translation
 - Axis,angle / translation
 - Later: Twist coordinates in Lie Algebra $\mathfrak{se}(3)$ of SE(3)

Pose Parametrization for Optimization

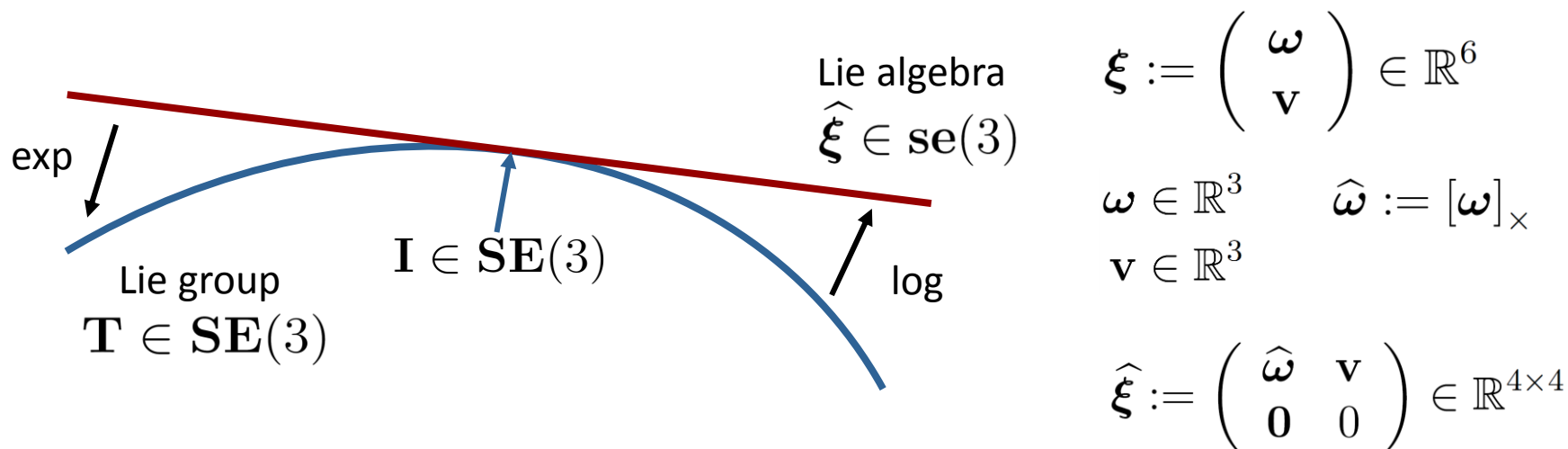
- Let's say we want to optimize a cost function $E(\theta)$ for the pose θ in some parametrization
- We need to set $\nabla_{\theta}E(\theta) = 0$

which we can tackle using gradient descent (or higher-order methods) by making steps on θ

$$\theta \leftarrow \theta - \lambda \nabla_{\theta}E(\theta)$$

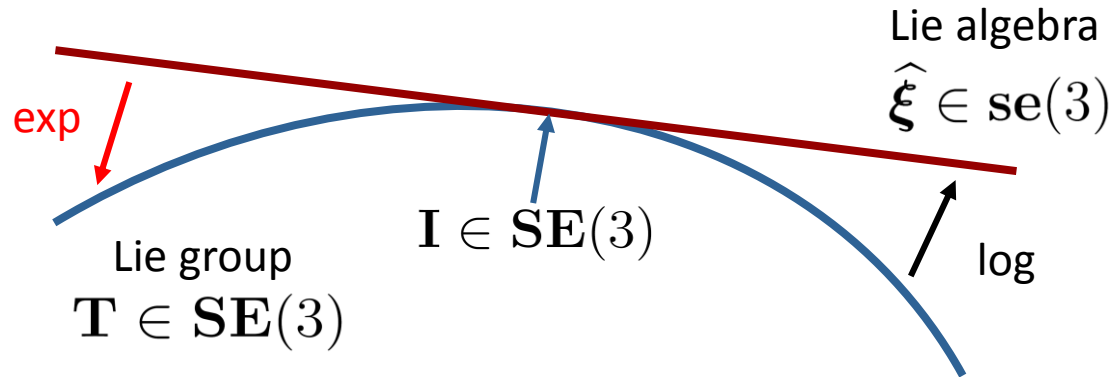
- When we determine the derivative of $E(\theta)$, we will require the derivative of $\mathbf{T}(\theta)$ for θ , which should have no singularities
- We also update the pose parametrization, which requires a minimal representation

SE(3) Lie Algebra for Representing Motion



- SE(3) is also a smooth manifold which makes it a Lie group
- The SE(3) Lie Algebra $\mathfrak{se}(3)$ provides an elegant way to parametrize poses for optimization
- Its elements $\hat{\xi} \in \mathfrak{se}(3)$ form the tangent space of SE(3) at its identity $\mathbf{I} \in \mathbf{SE}(3)$
- The $\mathfrak{se}(3)$ elements can be interpreted as rotational and translational velocities applied for some duration (twist) that explain the infinitesimal motion away from the identity transformation

Exponential Map of SE(3)

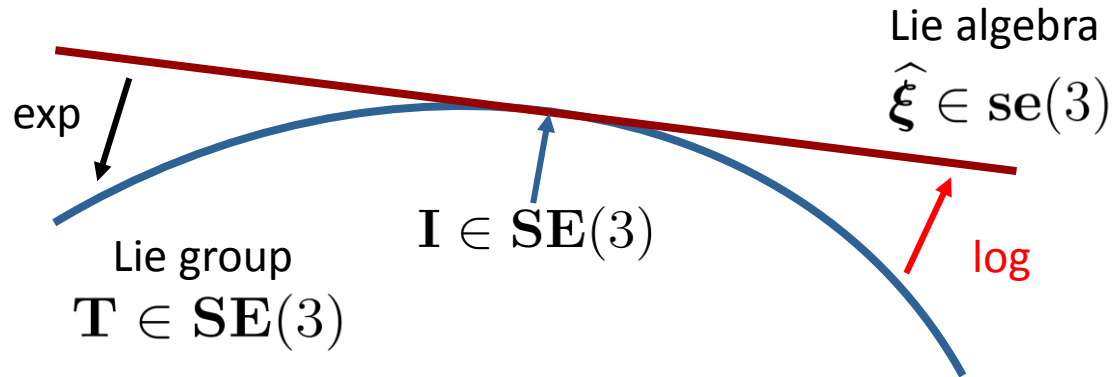


- The exponential map finds the transformation matrix for a twist:

$$\exp \left(\hat{\xi} \right) = \begin{pmatrix} \exp \left(\hat{\omega} \right) & \mathbf{A} \mathbf{v} \\ \mathbf{0} & 1 \end{pmatrix}$$

$$\exp \left(\hat{\omega} \right) = \mathbf{I} + \frac{\sin |\omega|}{|\omega|} \hat{\omega} + \frac{1 - \cos |\omega|}{|\omega|^2} \hat{\omega}^2 \quad \mathbf{A} = \mathbf{I} + \frac{1 - \cos |\omega|}{|\omega|^2} \hat{\omega} + \frac{|\omega| - \sin |\omega|}{|\omega|^3} \hat{\omega}^2$$

Logarithm Map of SE(3)



- The logarithm maps twists to transformation matrices:

$$\log(\mathbf{T}) = \begin{pmatrix} \log(\mathbf{R}) & \mathbf{A}^{-1}\mathbf{t} \\ \mathbf{0} & 0 \end{pmatrix}$$

$$|\omega| = \cos^{-1} \left(\frac{\text{tr}(\mathbf{R}) - 1}{2} \right) \quad \log(\mathbf{R}) = \frac{|\omega|}{2 \sin |\omega|} (\mathbf{R} - \mathbf{R}^T)$$

Optimization with Twist Coordinates

- How are twists useful in optimization?
- They provide a minimal representation without singularities close to identity
- Since $SE(3)$ is a smooth manifold, we can decompose $\mathbf{T}(\xi)$ in each optimization step into the transformation itself and a small increment (could be left or right-multiplied):

$$\mathbf{T}(\xi) := \mathbf{T}(\xi)\mathbf{T}(\delta\xi)$$

- Gradient descent operates on the auxiliary variable $\delta\xi$

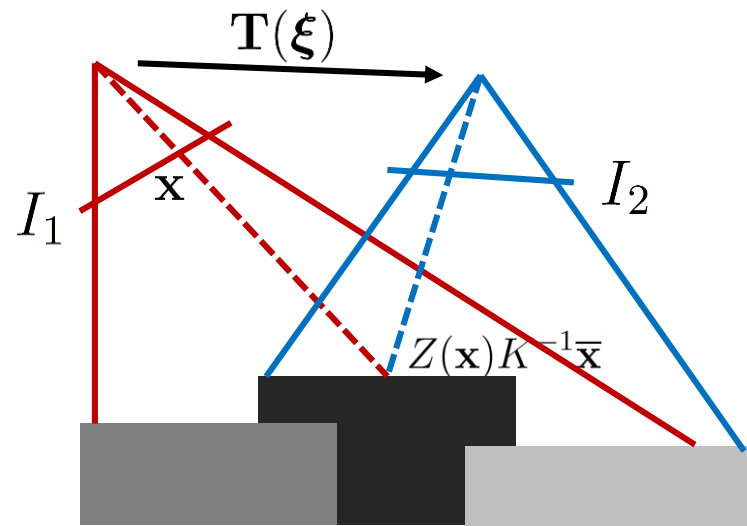
$$\delta\xi \leftarrow \mathbf{0} - \nabla_{\delta\xi} E(\delta\xi)$$

$$\hat{\xi} \leftarrow \log \left(\exp \left(\hat{\xi} \right) \exp \left(\delta\hat{\xi} \right) \right)$$

SE(3) Lie Algebra for Representing Motion

- C++ implementation: Sophus extension library for Eigen, by Hauke Strasdat, <https://github.com/strasdat/Sophus>
- Further reading on motion representation using the SE(3) Lie algebra:
 - Yi Ma, Stefano Soatto, Jana Kosecka, Shankar S. Sastry. An Invitation to 3-D Vision, Chapter 2: <http://vision.ucla.edu/MASKS/>
 - http://ingmec.ual.es/~jlblanco/papers/jlblanco2010geometry3D_techrep.pdf
 - <http://ethaneade.com/lie.pdf>

Dense Direct Image Alignment



- If we know pixel depth, we can „simulate“ an RGB-D image from a different view point
- Ideally, the warped image is the same like the image taken from that pose:

$$I_1(\mathbf{x}) = I_2(\pi(\mathbf{T}(\xi)Z(\mathbf{x})K^{-1}\bar{\mathbf{x}}))$$

- For RGB-D, we have the depth, but want to find the camera motion!

Dense Direct Image Alignment

- Given a camera motion, we can find and compare corresponding pixels through projection.
- We measure in one image a noisy version of the intensity in the other image:

$$I_1(\mathbf{x}) = I_2(\pi(\mathbf{T}(\boldsymbol{\xi})Z(\mathbf{x})K^{-1}\bar{\mathbf{x}})) + \epsilon$$

- A simple assumption is Gaussian noise, e.g. if the noise only comes from pixel noise on the chip $\epsilon \sim \mathcal{N}(0, \sigma_I^2)$

- If we further assume that the measurements are stochastically independent at each pixel, we can formulate the joint probability

$$p(\boldsymbol{\xi} \mid I_1, I_2) \propto p(I_1 \mid \boldsymbol{\xi}, I_2)p(\boldsymbol{\xi})$$

$$p(\boldsymbol{\xi} \mid I_1, I_2) \propto \prod_{\mathbf{x} \in \Omega} \mathcal{N}(I_1(\mathbf{x}) - I_2(\pi(\mathbf{T}(\boldsymbol{\xi})Z(\mathbf{x})K^{-1}\bar{\mathbf{x}})); 0, \sigma_I^2)$$

Dense Direct Image Alignment

- Maximum-likelihood estimation problem
- Optimize negative log-likelihood
 - Product becomes a summation
 - Exponentials disappear
 - Normalizers are independent of the pose

$$E(\boldsymbol{\xi}) = \text{const.} + \frac{1}{2} \sum_{\mathbf{x} \in \Omega} \frac{r(\mathbf{x}, \boldsymbol{\xi})^2}{\sigma_I^2}$$

$$r(\mathbf{x}, \boldsymbol{\xi}) = I_1(\mathbf{x}) - I_2(\pi(\mathbf{T}(\boldsymbol{\xi})Z(\mathbf{x})K^{-1}\bar{\mathbf{x}}))$$

- This non-linear least squares error function can be efficiently optimized using standard methods (Gauss-Newton, Levenberg-Marquardt)

Least Squares Optimization

- If the residuals would be linear ξ , i.e., $r(\xi) = \mathbf{A}\xi + \mathbf{b}$, optimization would be simple, has a closed-form solution
- In this case, the error function and its derivatives are

$$E(\xi) = \frac{1}{2}r(\xi)^T \mathbf{W}r(\xi)$$

$$\nabla_{\xi} E(\xi) = \nabla_{\xi} r(\xi)^T \mathbf{W}r(\xi) = \mathbf{A}^T \mathbf{W}r(\xi)$$

$$\nabla_{\xi}^2 E(\xi) = \mathbf{A}^T \mathbf{W} \mathbf{A}$$

- Setting the first derivative to zero yields

$$\nabla_{\xi} E(\xi) = \nabla_{\xi} E(\xi_0) + \nabla_{\xi}^2 E(\xi_0)(\xi - \xi_0) = 0$$

$$\xi = \xi_0 - \nabla_{\xi}^2 E(\xi_0)^{-1} \nabla_{\xi} E(\xi_0)$$

$$\xi = \xi_0 - (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{W}r(\xi_0)$$

Non-linear Least Squares Optimization

- In direct image alignment, the residuals are non-linear in ξ
- Gauss-Newton method, iterate:

- Linearize residuals $\tilde{r}(\xi) = r(\xi_0) + \nabla_{\xi} r(\xi)(\xi - \xi_0)$

$$\tilde{E}(\xi) = \frac{1}{2} \tilde{r}(\xi)^T \mathbf{W} \tilde{r}(\xi)$$

$$\nabla_{\xi} \tilde{E}(\xi) = \nabla_{\xi} r(\xi)^T \mathbf{W} \tilde{r}(\xi)$$

$$\nabla_{\xi}^2 \tilde{E}(\xi) = \nabla_{\xi} r(\xi)^T \mathbf{W} \nabla_{\xi} r(\xi)$$

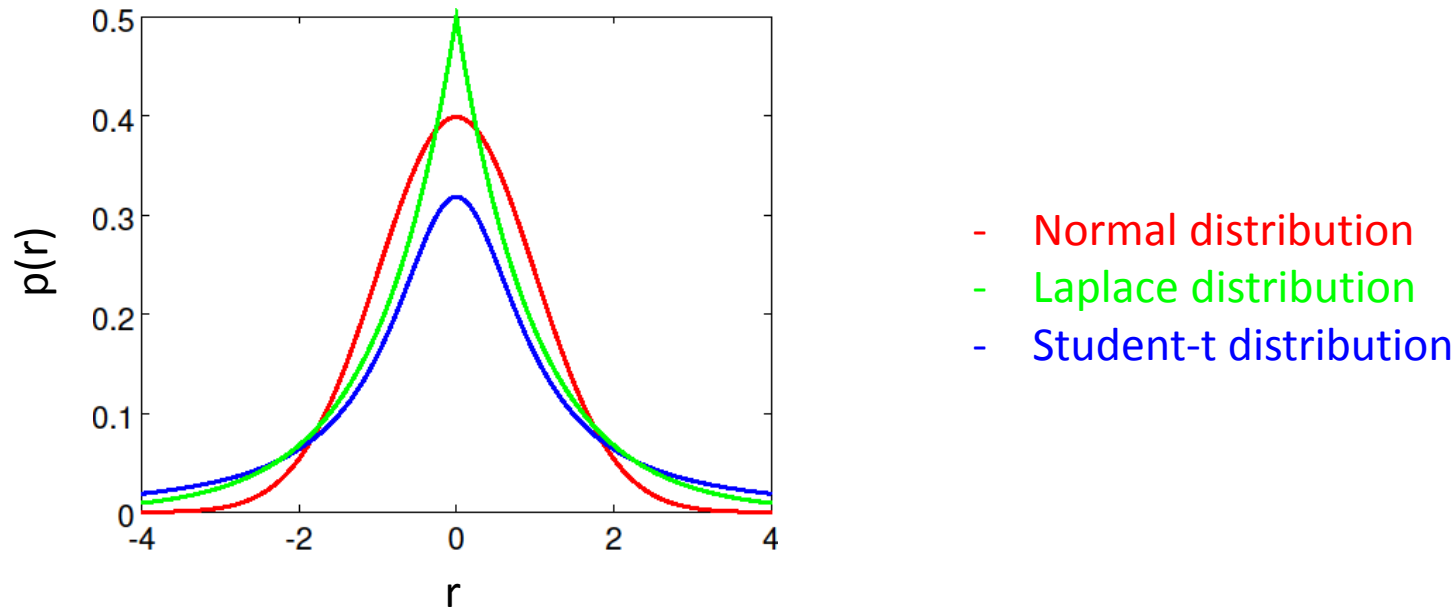
- Solve linearized system

$$\nabla_{\xi} \tilde{E}(\xi) = \nabla_{\xi} \tilde{E}(\xi_0) + \nabla_{\xi}^2 \tilde{E}(\xi_0)(\xi - \xi_0) = 0$$

$$\xi \leftarrow \xi - \nabla_{\xi}^2 \tilde{E}(\xi)^{-1} \nabla_{\xi} \tilde{E}(\xi)$$

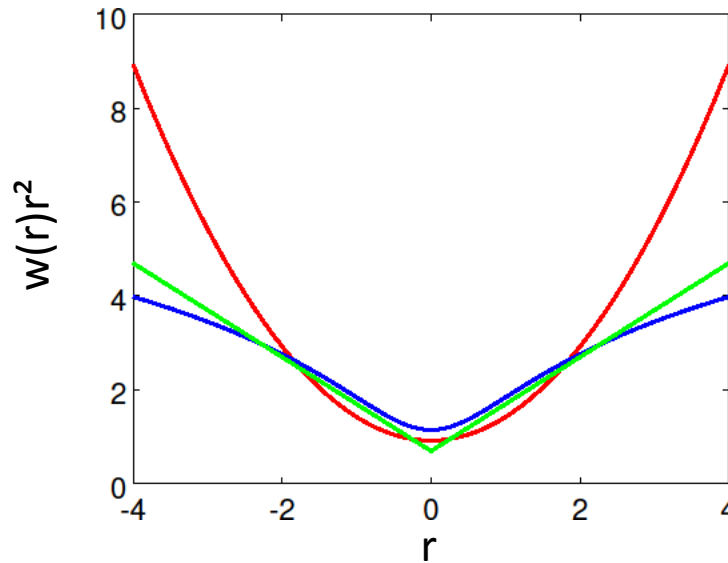
$$\xi \leftarrow \xi - (\nabla_{\xi} r(\xi)^T \mathbf{W} \nabla_{\xi} r(\xi))^{-1} \nabla_{\xi} r(\xi)^T \mathbf{W} r(\xi)$$

Actual Residual Distribution



- The Gaussian noise assumption is not valid
- Many outliers (occlusions, motion, etc.)
- Residuals are distributed with more mass on the larger values

Iteratively Reweighted Least Squares



- Normal distribution
- Laplace distribution
- Student-t distribution

- Can we change the residual distribution in the least squares optimization?
- We can reweight the residuals in each iteration to adapt residual distribution

$$E(\boldsymbol{\xi}) = \frac{1}{2} \sum_{\mathbf{x} \in \Omega} w(r(\mathbf{x}, \boldsymbol{\xi})) \frac{r(\mathbf{x}, \boldsymbol{\xi})^2}{\sigma_I^2}$$

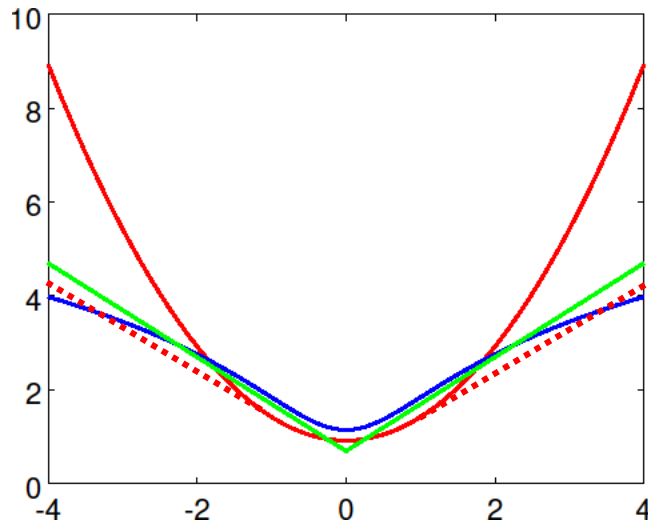
E.g., for Laplace distribution:

$$w(r(\mathbf{x}, \boldsymbol{\xi})) = |r(\mathbf{x}, \boldsymbol{\xi})|^{-1}$$

Huber-Loss

- Huber-loss „switches“ between normal (locally at mean) and Laplace distribution

$$\|r\|_{\delta} = \begin{cases} \frac{1}{2} \|r\|_2^2 & \text{if } \|r\|_2 \leq \delta \\ \delta (\|r\|_1 - \frac{1}{2}\delta) & \text{otherwise} \end{cases}$$



..... Huber-loss for $\delta = 1$

Linearization of Image Alignment Residuals

- In our direct image alignment case, the linearized residuals are

$$\nabla_{\xi} r(\mathbf{x}, \xi) = -\nabla_{\pi} I_2(\pi(\mathbf{p}(\mathbf{x}, \xi))) \cdot \nabla_{\xi} \pi(\mathbf{p}(\mathbf{x}, \xi))$$

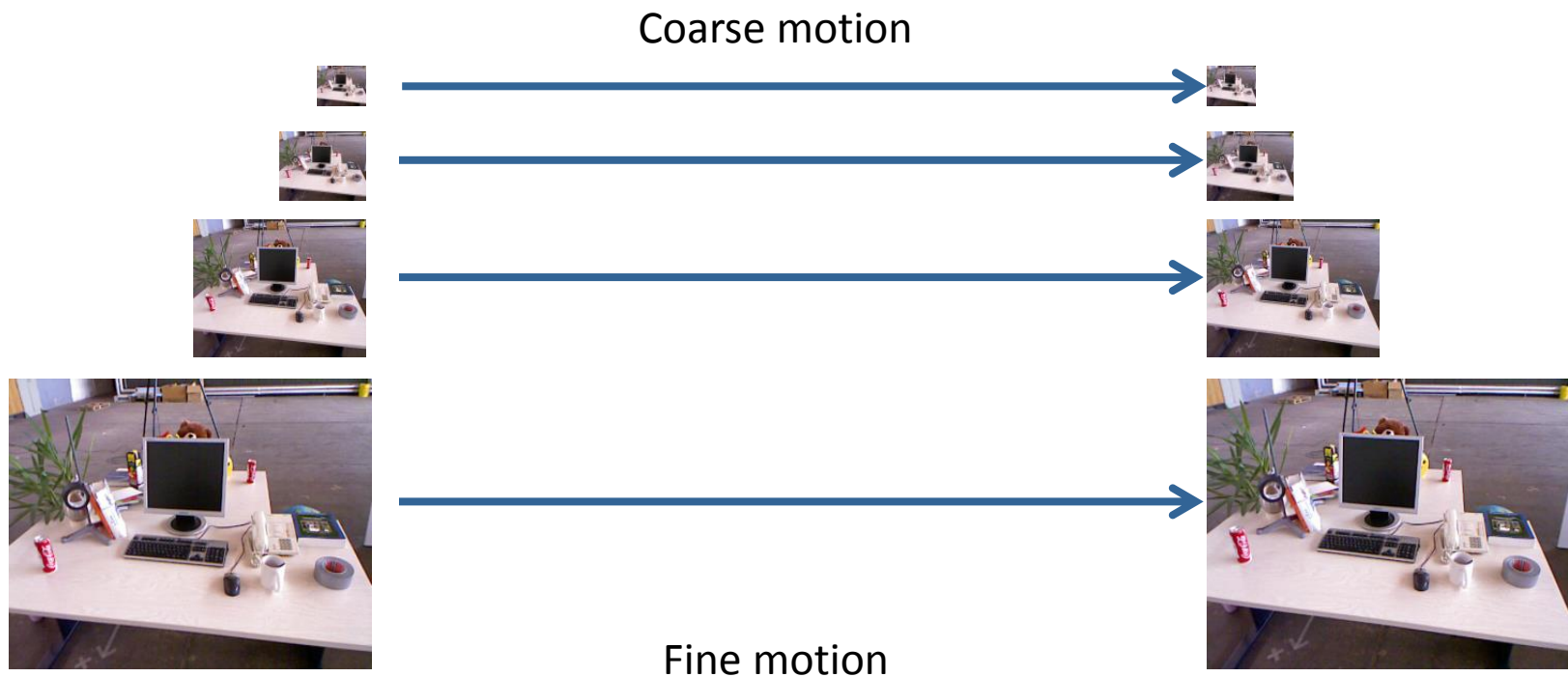
with $\mathbf{p}(\mathbf{x}, \xi) = \mathbf{T}(\xi)Z(\mathbf{x})K^{-1}\bar{\mathbf{x}}$

$$r(\mathbf{x}, \xi) = I_1(\mathbf{x}) - I_2(\pi(\mathbf{p}(\mathbf{x}, \xi)))$$

- Linearization is only valid for motions that change the projection in a small image neighborhood (where the gradient hints into the direction)

Coarse-To-Fine

- Adapt size of the neighborhood from coarse to fine

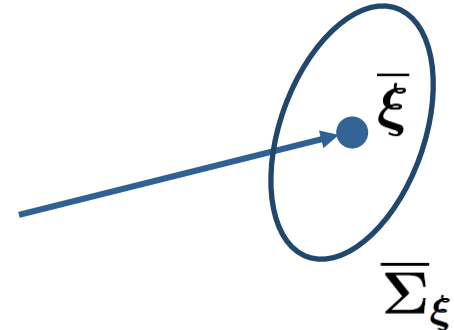


Covariance of the Pose Estimate

- Non-linear least squares determines a Gaussian estimate

$$p(\xi \mid I_1, I_2) = \mathcal{N}(\bar{\xi}, \bar{\Sigma}_\xi)$$

$$\bar{\Sigma}_\xi = (\nabla_\xi r(\bar{\xi})^T \mathbf{W} \nabla_\xi r(\bar{\xi}))^{-1}$$



- Due to pose decomposition, we have to change the coordinate frame of the covariance using the adjoint in SE(3)

$$p(\xi \mid I_1, I_2) = \mathcal{N}(\bar{\xi}, \text{ad}_{\mathbf{T}(\bar{\xi})} \bar{\Sigma}_{\delta\xi} \text{ad}_{\mathbf{T}(\bar{\xi})}^T)$$

$$\bar{\Sigma}_{\delta\xi} = (\nabla_{\delta\xi} r(\delta\xi = 0, \bar{\xi})^T \mathbf{W} \nabla_{\delta\xi} r(\delta\xi = 0, \bar{\xi}))^{-1}$$

$$\text{ad}_{\mathbf{T}} = \begin{pmatrix} \mathbf{R} & [\mathbf{t}]_{\times} \mathbf{R} \\ \mathbf{0} & \mathbf{R} \end{pmatrix} \in \mathbb{R}^{6 \times 6}$$

Lessons Learned

- The SE(3) Lie algebra is an elegant way of motion representation, especially for gradient-based optimization of motion parameters
- Non-linear least squares optimization is a versatile tool that can be applied for direct image alignment
- Iteratively Reweighted Least Squares allows for overcoming the limitation of basic least squares on the Gaussian residual distribution/L2 loss on the residuals
- Dense RGB-D odometry through direct image alignment can be implemented in a non-linear least squares framework.
 - The linear approximation of the residuals requires a coarse-to-fine optimization scheme
 - Non-linear least squares also provides the pose covariance

Questions ?