

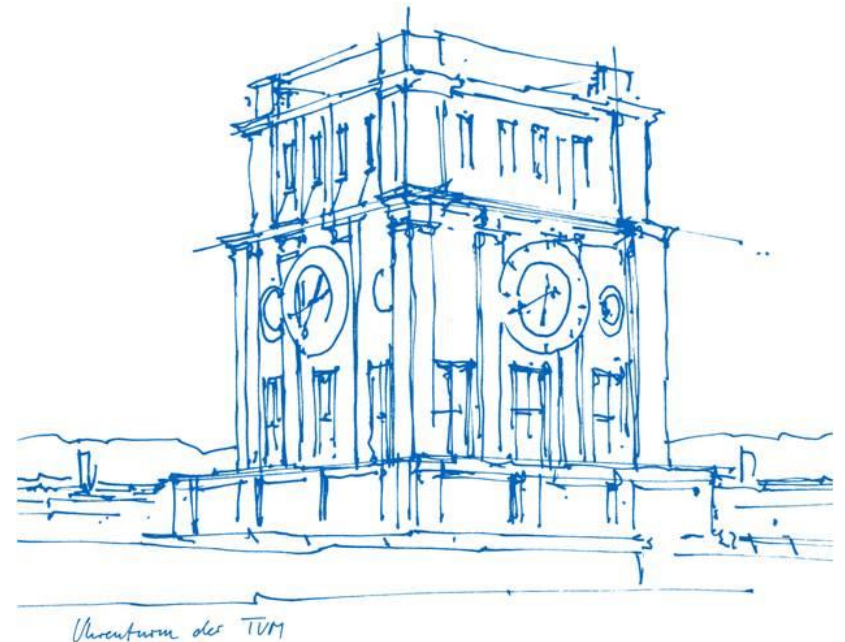
# Deep Virtual Stereo Odometry: Leveraging Deep Depth Prediction for Monocular Direct Sparse Odometry

Nan Yang, Rui Wang, Jörg Stückler, Daniel Cremers

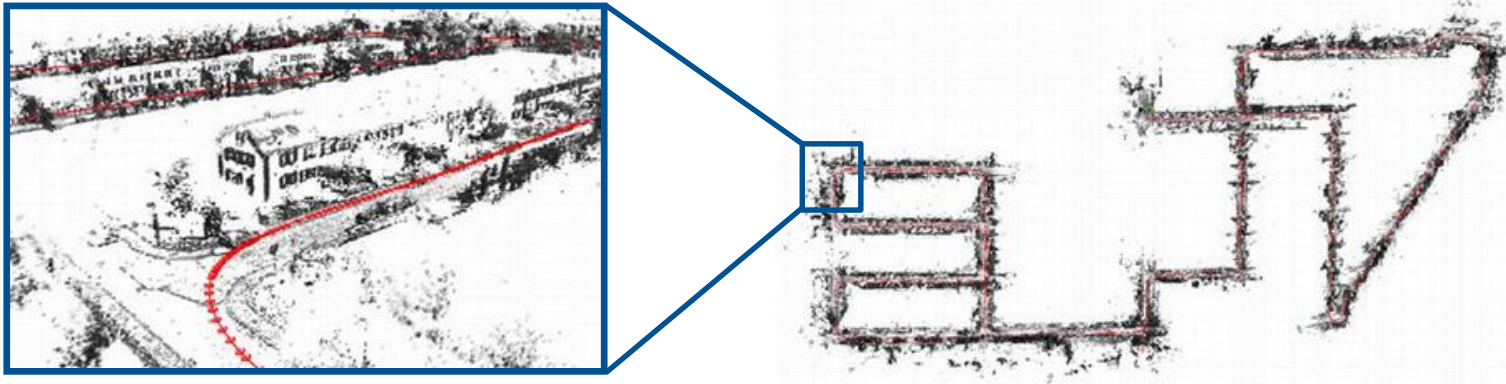
Lars Carius

Technical University of Munich

Munich, March 16th 2020



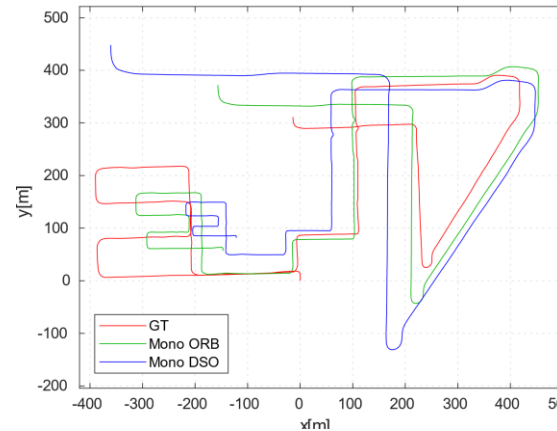
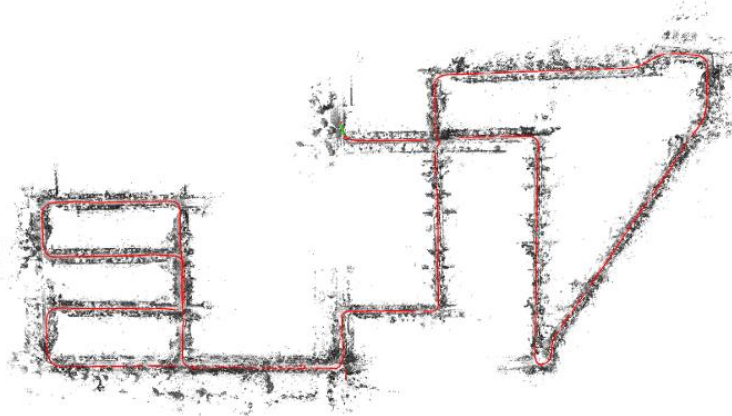
# Intro – Visual Odometry



- **VO:** Determine position and orientation of camera from visual input
- **VSLAM:** Globally optimize both map and camera pose

# Intro – Shortcomings of Existing Methods

- **Mono camera setups: Scale Drift**
  - Geometrically impossible to recover scale!



- **End-to-end trained Neural Networks: Performance**
- **Stereo Camera, RGB-D, LiDAR: Cost & Calibration**
- **This paper: Hybrid method (Single camera enhanced with Deep Learning)**

# Outline

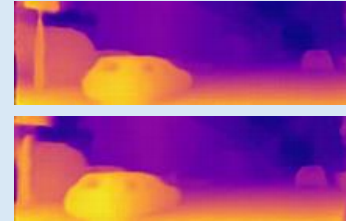
- Introduction
- Method
  - High-Level Concept
  - Depth Estimation
  - Bundle Adjustment
  - Deep Virtual Stereo Odometry
- Experiments & Results
  - Monocular Depth Estimation
  - Monocular Visual Odometry
- Personal Comments
- Summary

# Method – High-Level Concept

## Deep Learning Pipeline



Mono camera image



Left disparity map

Right disparity map

Consistent metric  
scale initialization

Virtual stereo  
photometric error

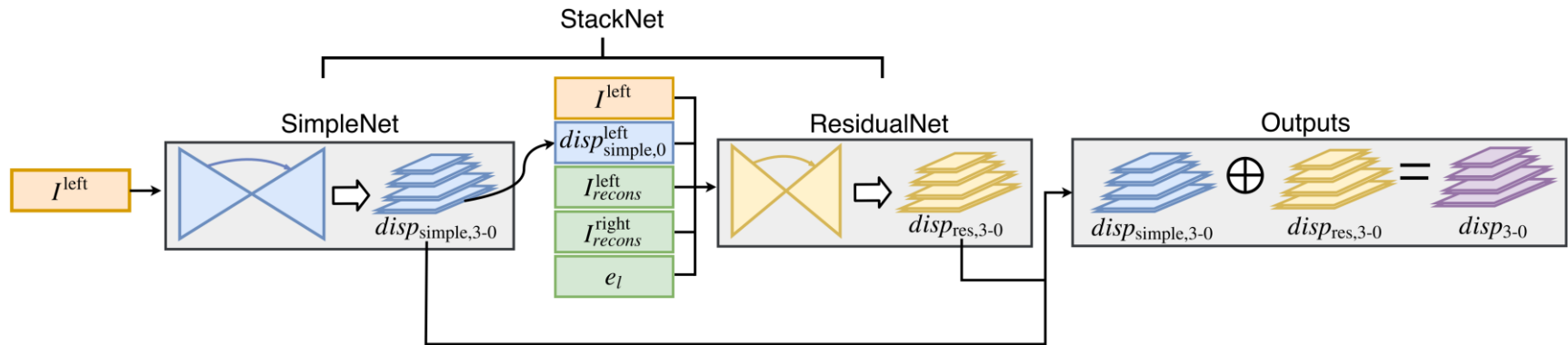
## Classic Visual Odometry Pipeline



Mono camera image



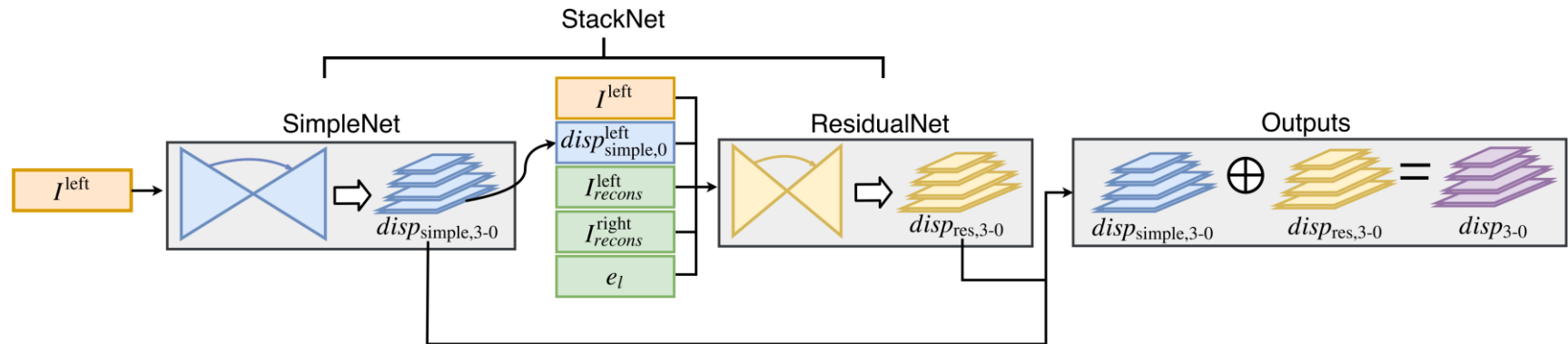
# Method – Depth Estimation



## Architecture:

- 2 fully-convolutional encoder-decoder networks with skip connections
- ResidualNet refines disparity maps output by SimpleNet

# Method – Depth Estimation



## Architecture:

- 2 fully-convolutional encoder-decoder networks with skip connections
- ResidualNet refines disparity maps output by SimpleNet

## Loss Function:

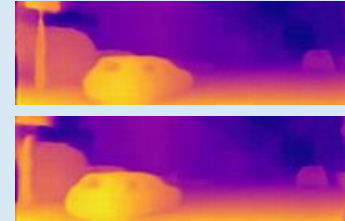
$$\mathcal{L}_s = \underbrace{\alpha_U (\mathcal{L}_U^{left} + \mathcal{L}_U^{right})}_{\text{Self-supervised loss}} + \underbrace{\alpha_S (\mathcal{L}_S^{left} + \mathcal{L}_S^{right})}_{\text{Supervised loss}} + \underbrace{\alpha_{lr} (\mathcal{L}_{lr}^{left} + \mathcal{L}_{lr}^{right})}_{\text{Left-right consistency loss}} + \underbrace{\alpha_{smooth} (\mathcal{L}_{smooth}^{left} + \mathcal{L}_{smooth}^{right})}_{\text{Disparity smoothness regularization}} + \underbrace{\alpha_{occ} (\mathcal{L}_{occ}^{left} + \mathcal{L}_{occ}^{right})}_{\text{Occlusion regularization}}$$

# Method – Deep Virtual Stereo Odometry

## Deep Learning Pipeline



Mono camera image



Left disparity map

Right disparity map

Consistent metric  
scale initialization

Virtual stereo  
photometric error

## Classic Visual Odometry Pipeline



Mono camera image



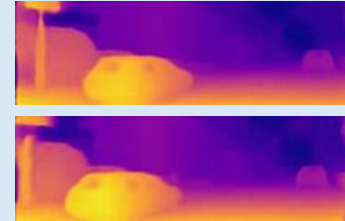


# Method – Deep Virtual Stereo Odometry

## Deep Learning Pipeline



Mono camera image



Left disparity map

Right disparity map



Consistent metric  
scale initialization  
→ mitigates scale  
drift

Virtual stereo  
photometric error

## Classic Visual Odometry Pipeline



Mono camera image

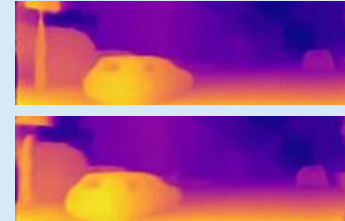


# Method – Deep Virtual Stereo Odometry

## Deep Learning Pipeline



Mono camera image



Left disparity map

Right disparity map

Consistent metric  
scale initialization

→ mitigates scale  
drift

Virtual stereo  
photometric error

→ increases accuracy  
& mitigates scale drift

## Classic Visual Odometry Pipeline



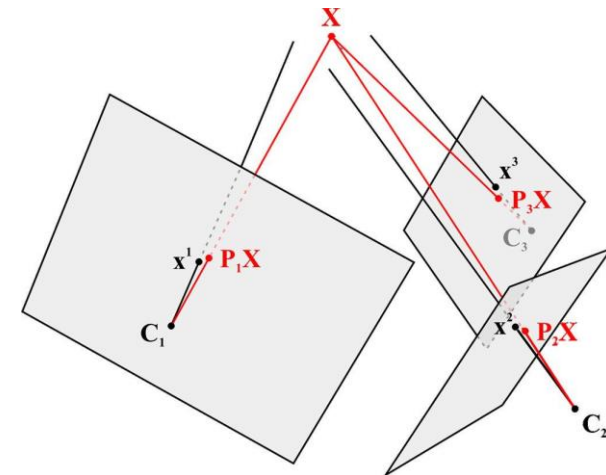
Mono camera image



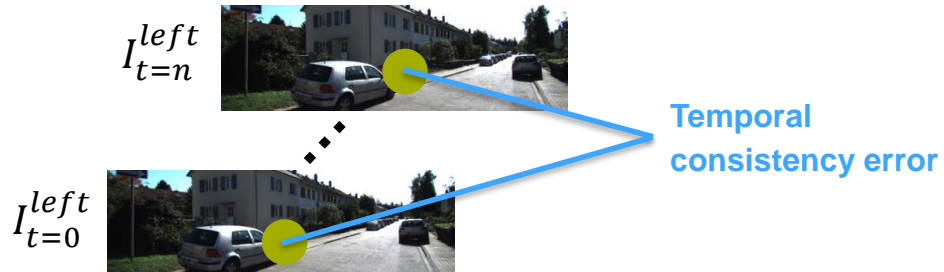
# Method – Definitions

- **Bundle adjustment (BA):** Optimize 3D points, camera motion & camera params based on 2D images
- **Direct bundle adjustment:** No reprojection error, optimize in image space
- **Sparse bundle adjustment:** Use only a subset of pixels
- **Windowed bundle adjustment:** BA with sliding-window (heuristic keyframes)

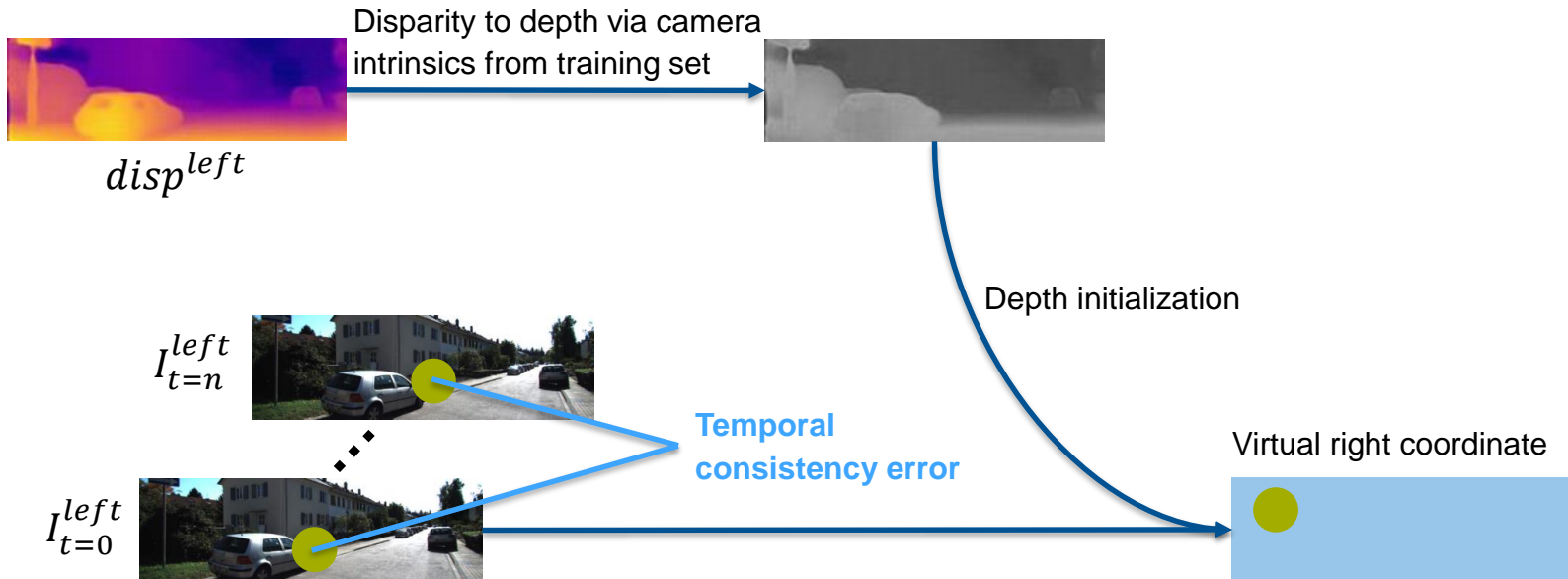
➤ **Here: Windowed direct sparse BA**



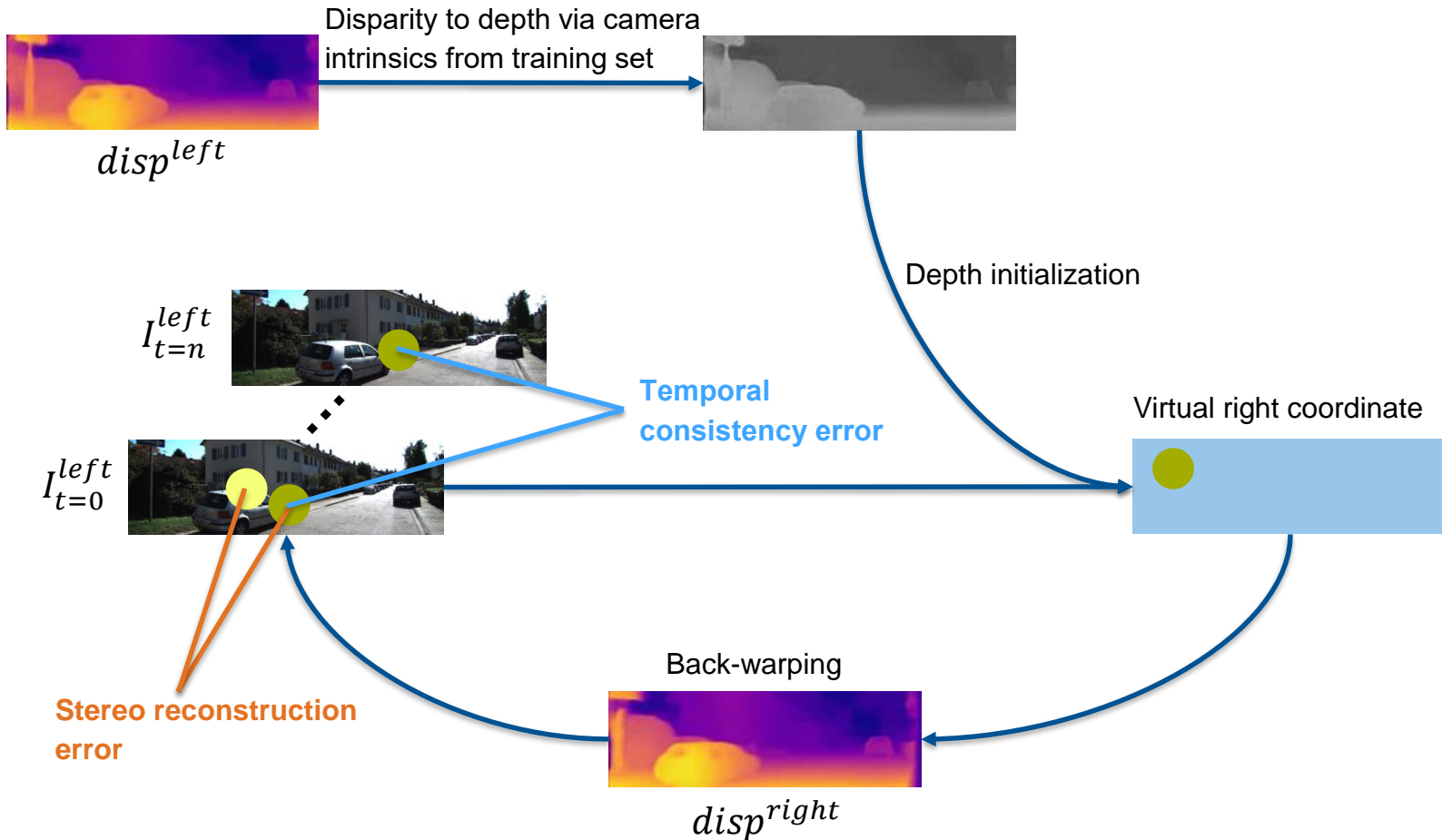
# Method – Virtual Stereo Photometric Error



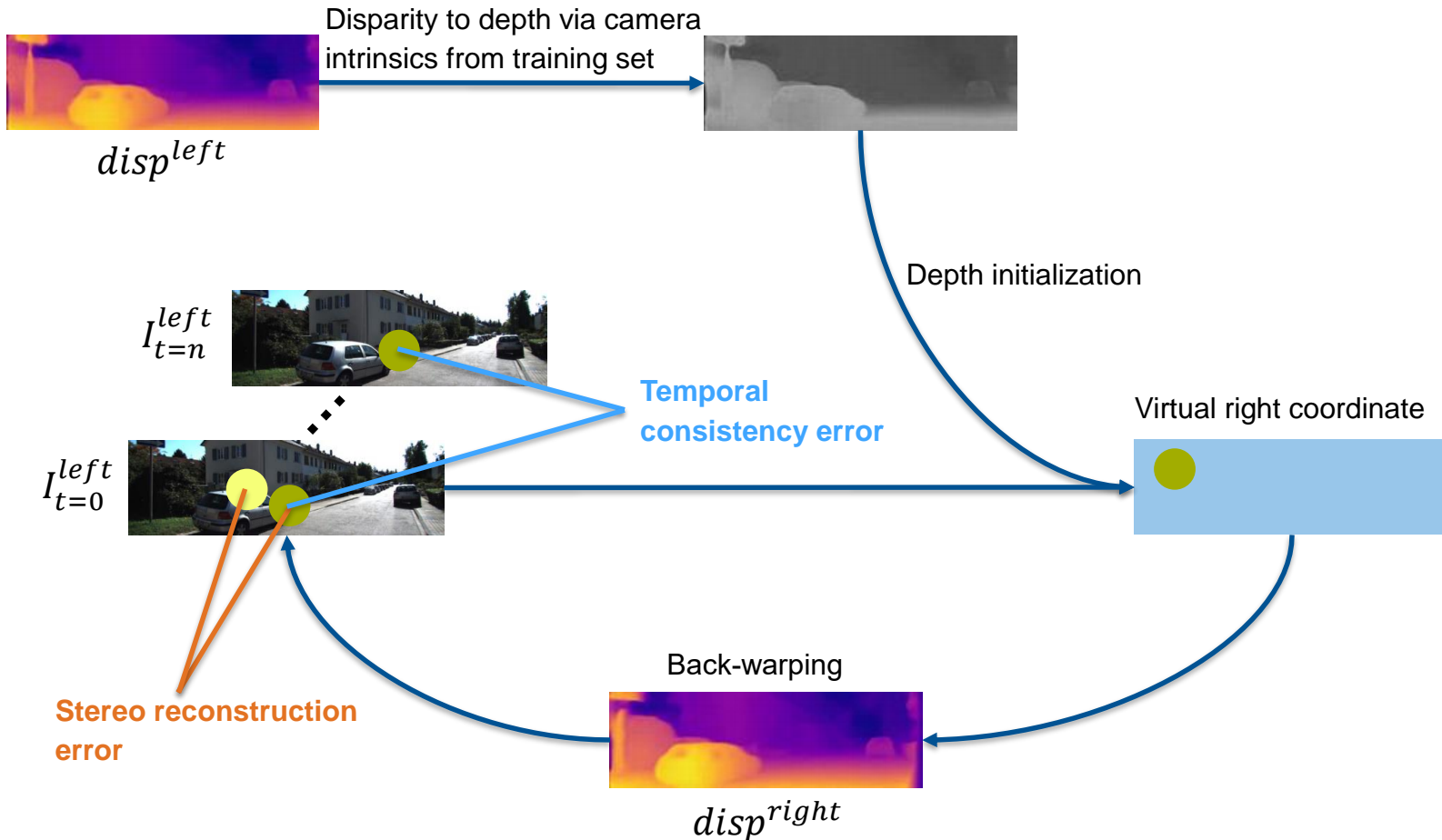
# Method – Virtual Stereo Photometric Error



# Method – Virtual Stereo Photometric Error



# Method – Virtual Stereo Photometric Error



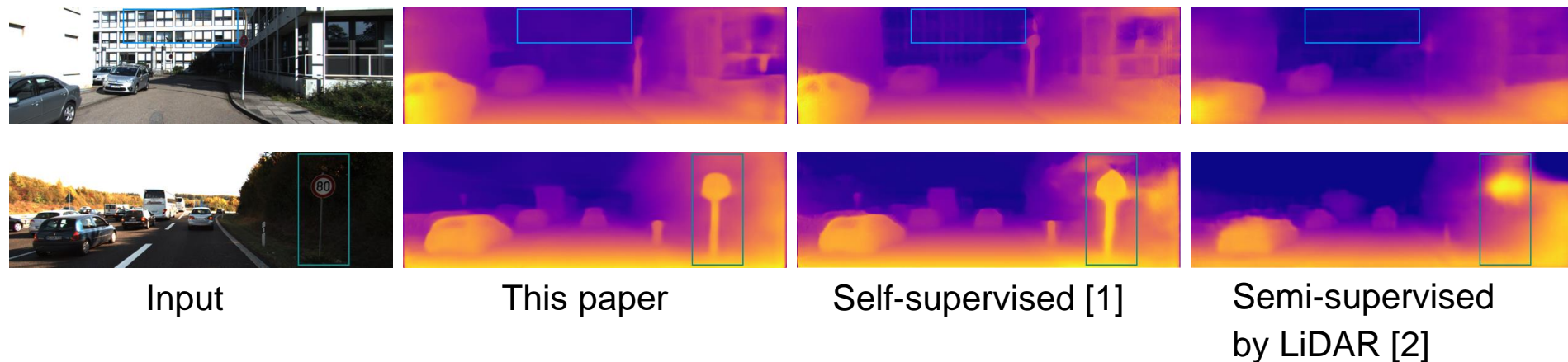
3D-points and camera motion are recovered by jointly optimizing **temporal consistency** and **stereo reconstruction error** using Gauss-Newton

# Experiments & Results – Depth Estimation

**Training schedule** („easy task first, then harder task, then easy again to smooth outliers“):

1. Train SimpleNet: semi-supervised → self-supervised → semi-supervised
2. Freeze SimpleNet, train ResidualNet: semi-supervised → self-supervised → semi-supervised

## Results:



Source:

[1] Godard, C., Mac Aodha, O., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. arXiv preprint arXiv:1609.03677 (2016)

[2] Kuznetsov, Y., Stückler, J., Leibe, B.: Semi-supervised deep learning for monocular depth map prediction. arXiv preprint arXiv:1702.02706 (2017)

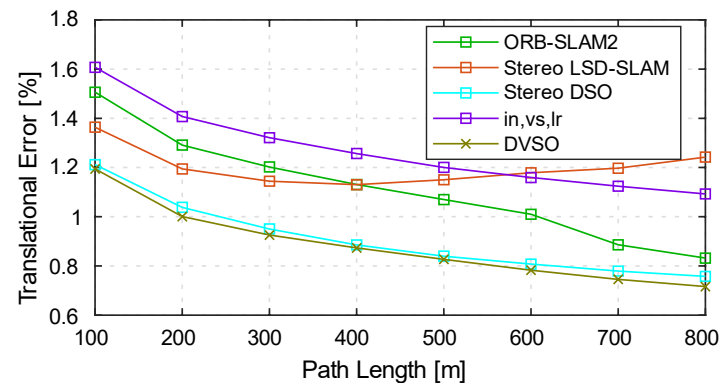
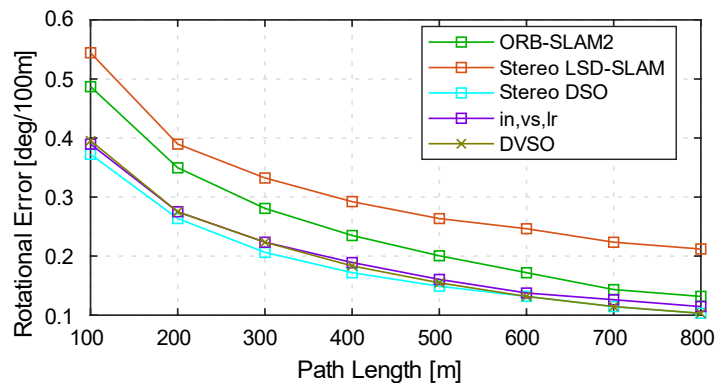


# Experiments & Results – Monocular Visual Odometry

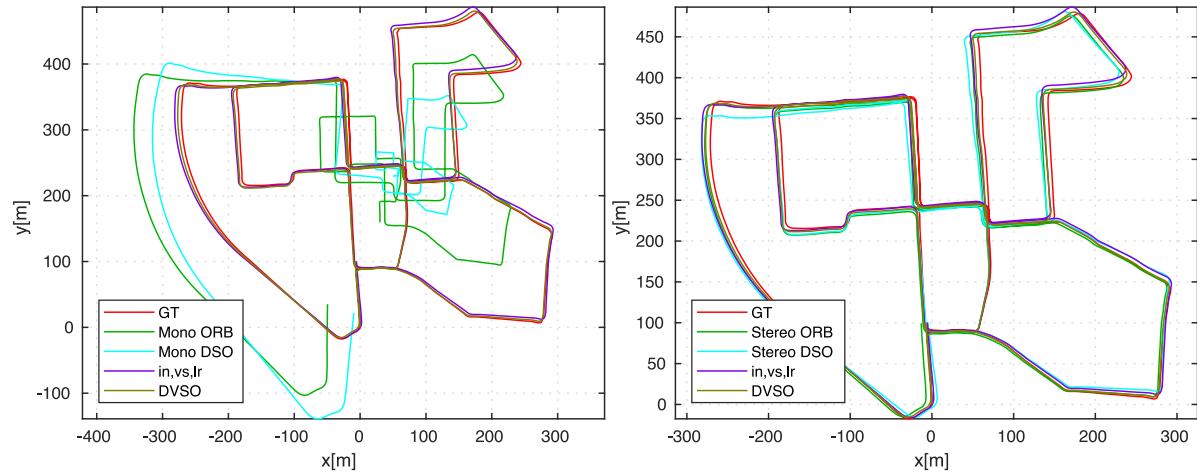
## Ablation study:

1. Initializing depth with left disparity prediction
2. Using right disparity for virtual stereo term in windowed bundle adjustment
3. Checking left-right disparity consistency for sparse point selection
4. Tuning virtual stereo baseline

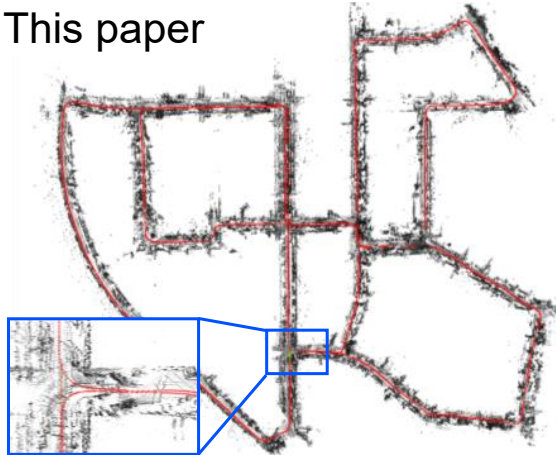
## Results:



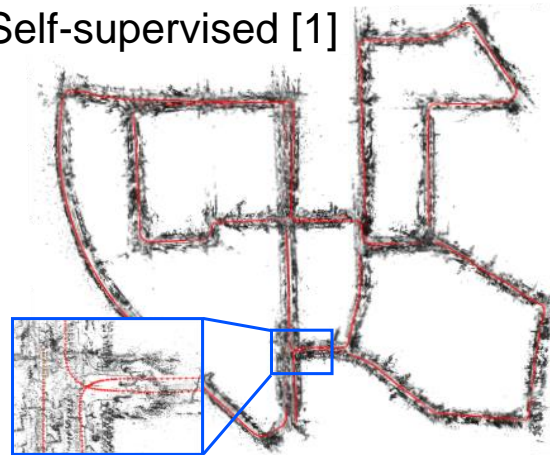
# Experiments & Results – Monocular Visual Odometry



This paper



Self-supervised [1]



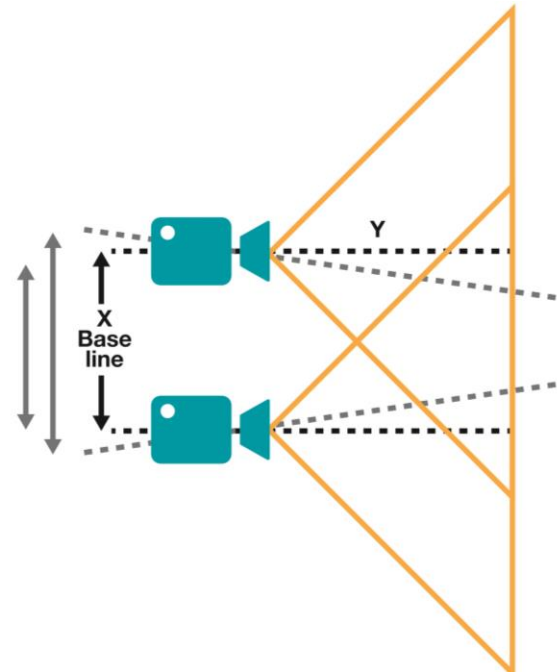
Source:

[1] Godard, C., Mac Aodha, O., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. arXiv preprint arXiv:1609.03677 (2016)

# Personal Comments

## Baseline tuning is cheating

- Stereo information at test time required
- Future work: Online fine tuning



# Personal Comments

## Necessity for Deep Learning

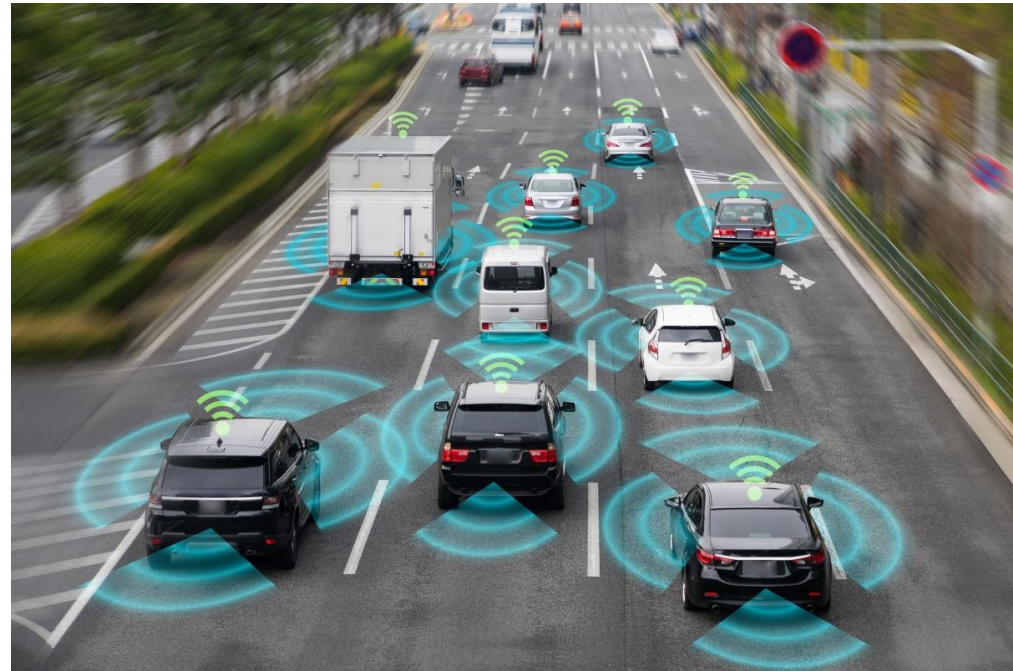
- Stereo cameras cheap
- Deep Learning needs good training data
- Safety & error accountability



# Personal Comments

## Enables fleet learning

- Few vehicles with stereo camera to provide training data
- Remaining vehicles use DL and operate in a similar enough domain



# Summary: Deep Virtual Stereo Odometry

- **Monocular Input + Deep Depth Predictions → Stereo Performance**
  - Scale drift eliminated
- **Semi-Supervised Disparity Predictions**
  - Stereo reconstruction for self-supervision
  - Sparse DSO for supervision
- **Incorporation into Visual Odometry Pipeline**
  - Metric scale initialization
  - Virtual stereo error

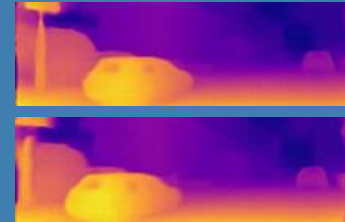


# Discussion

## Deep Learning Pipeline



Mono camera image



Left disparity map

Right disparity map

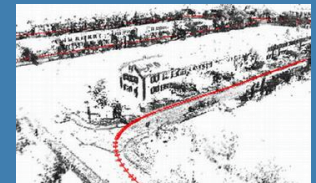
Consistent metric  
scale initialization

Virtual stereo  
photometric error

## Classic Visual Odometry Pipeline



Mono camera image



# Backup: From Disparity to Depth

$$d(\vec{p}) = \frac{b \cdot f_x}{\text{disp}(\vec{p})}$$

- $d(\vec{p})$ : Depth at point  $\mathbf{p}$
- $b$ : Baseline
- $f_x$ : Focal length (camera intrinsic)
- $\text{disp}(\vec{p})$ : Disparity at point  $\mathbf{p}$

