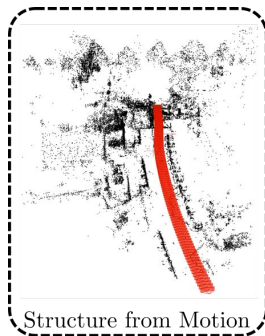# SfM-TTR: Structure from Motion for Test-Time Refinement of Single-View Depth Networks

Sergio Izquierdo, Javier Civera
University of Zaragoza

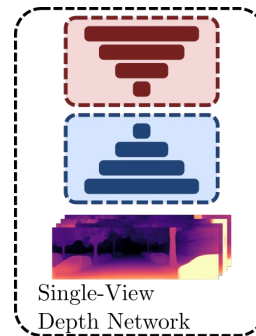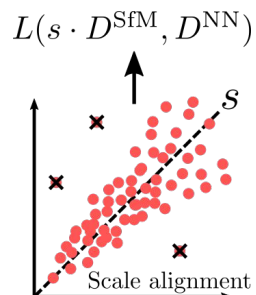**THU-PM-082**

# Scene Reconstruction

Multi-view
traditional methods



Structure from Motion

Single-view
deep learning



Single-View
Depth Network

# Method Overview

Multi-view
traditional methods

Single-view
deep learning

$$L(s \cdot D^{\mathrm{SfM}}, D^{\mathrm{NN}})$$

$s$
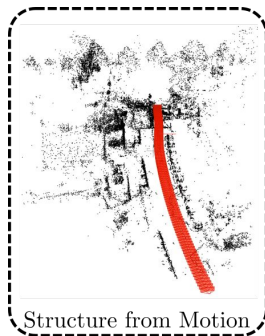
Scale alignment

Structure from Motion

Single-View
Depth Network

# Method Overview

Multi-view
traditional methods

Single-view
deep learning



$$L(s \cdot D^{\text{SfM}}, D^{\text{NN}})$$

$s$

Scale alignment

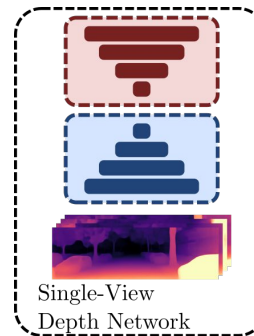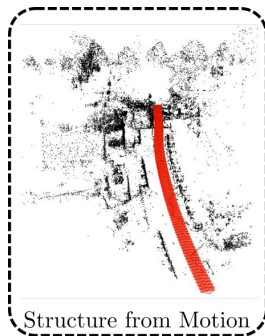Structure from Motion

Single-View
Depth Network

Refined Network

# Method Overview

Multi-view
traditional methods

Single-view
deep learning



$L(s \cdot D^{\mathrm{SfM}}, D^{\mathrm{NN}})$

$s$

Scale alignment

Structure from Motion

Single-View
Depth Network

Refined Network

- Improves Supervised and Self-supervised

# Method Overview

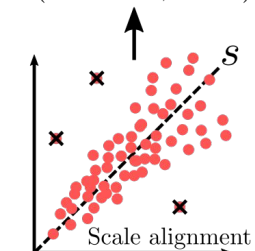Multi-view
traditional methods

Single-view
deep learning



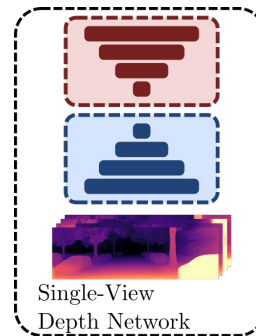$$L(s \cdot D^{\mathrm{SfM}}, D^{\mathrm{NN}})$$

$s$

Scale alignment

Structure from Motion

Single-View
Depth Network

Refined Network

- Improves Supervised and Self-supervised
- 27% RMSE reduction

# Method Overview



Multi-view
traditional methods

Single-view
deep learning

$L(s \cdot D^{\mathrm{SfM}}, D^{\mathrm{NN}})$

$s$

Scale alignment

Structure from Motion

Single-View
Depth Network

Refined Network

- Improves Supervised and Self-supervised
- 27% RMSE reduction
- Better estimates for further areas

# Method Overview

Multi-view
traditional methods

Single-view
deep learning



$L(s \cdot D^{\mathrm{SfM}}, D^{\mathrm{NN}})$

$s$

Scale alignment

Structure from Motion

Single-View
Depth Network

Refined Network

- Improves Supervised and Self-supervised
- 27% RMSE reduction
- Better estimates for further areas
- SOTA on KITTI

# Motivation

- Multi-view traditional methods
  - ✔ 3D geometry based
  - ✔ Accurate estimations
  - ✗ Sparse depth map
  - ✗ Not learned priors


- Single-view networks
  - ✔ Learned priors
  - ✔ Dense estimations
  - ✗ Vast collections of images
  - ✗ No geometry based



Structure from Motion

how to combine both?



Single-View Depth Network

# Method



Input sequence

# Method


Input sequence

SfM-TTR

# Method



Input sequence

SfM-TTR

Structure from Motion

Single-View
Depth Network

1. Compute SfM reconstruction
2. Compute Network predictions

# Method



Input sequence

SfM-TTR

Structure from Motion

Scale alignment

$S$

Single-View Depth Network

1. Compute SfM reconstruction
2. Compute Network predictions
3. Align both depth maps

JUNE 18-22, 2023
CVPR
VANCOUVER, CANADA

# Method



Input sequence

SfM-TTR

$L(s \cdot D^{\mathrm{SfM}}, D^{\mathrm{NN}})$

Scale alignment

Structure from Motion

Single-View
Depth Network

1. Compute SfM reconstruction
2. Compute Network predictions
3. Align both depth maps
4. Use sparse depth as pseudo-gt

$$\mathcal{L} = \frac{1}{|\mathcal{D}_j^{\mathrm{SfM}}|} \sum_l w_{l,j}^{\boldsymbol{\theta}} \| \hat{s} \cdot D_{l,j}^{\mathrm{SfM}} - D_{l,j}^{\mathrm{NN}} \|_1$$

# Method



Input sequence

SfM-TTR

$L(s \cdot D^{\mathrm{SfM}}, D^{\mathrm{NN}})$

$s$

Scale alignment

Structure from Motion

Single-View
Depth Network

1. Compute SfM reconstruction
2. Compute Network predictions
3. Align both depth maps
4. Use sparse depth as pseudo-gt

$$\mathcal{L} = \frac{1}{|\mathcal{D}_j^{\mathrm{SfM}}|} \sum_l w_{l,j}^{\boldsymbol{\theta}} \| \hat{s} \cdot D_{l,j}^{\mathrm{SfM}} - D_{l,j}^{\mathrm{NN}} \|_1$$

reprojection error

scale from alignment

$$w_{l,j}^{\boldsymbol{\theta}} = \exp(-\|\mathbf{r}_{l,j}\|_2^2)$$

# Method



Input sequence

Structure from Motion

SfM-TTR

$L(s \cdot D^{\mathrm{SfM}}, D^{\mathrm{NN}})$

Scale alignment

Single-View
Depth Network

1. Compute SfM reconstruction
2. Compute Network predictions
3. Align both depth maps
4. Use sparse depth as pseudo-gt

$$\mathcal{L} = \frac{1}{|\mathcal{D}_j^{\mathrm{SfM}}|} \sum_l w_{l,j}^{\boldsymbol{\theta}} \| \hat{s} \cdot D_{l,j}^{\mathrm{SfM}} - D_{l,j}^{\mathrm{NN}} \|_1$$

5. Optimize the encoder

# Method



Input sequence

**SfM-TTR**

$L(s \cdot D^{\mathrm{SfM}}, D^{\mathrm{NN}})$

Structure from Motion

Scale alignment

Single-View
Depth Network

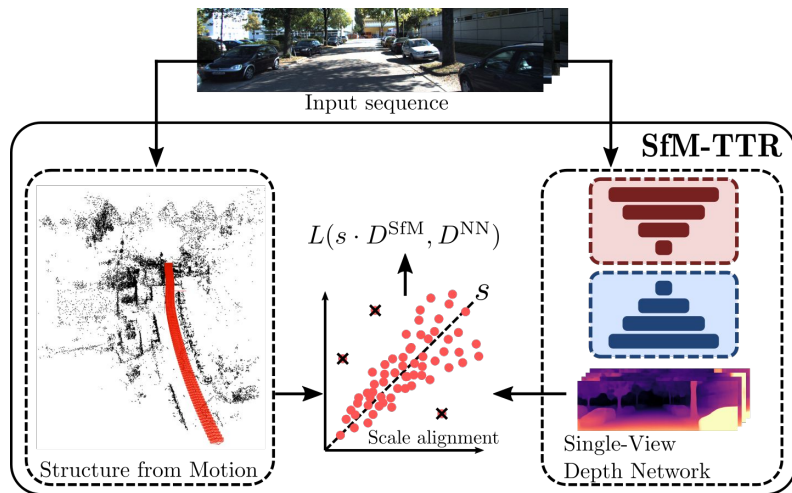Refined depth predictions

1. Compute SfM reconstruction
2. Compute Network predictions
3. Align both depth maps
4. Use sparse depth as pseudo-gt

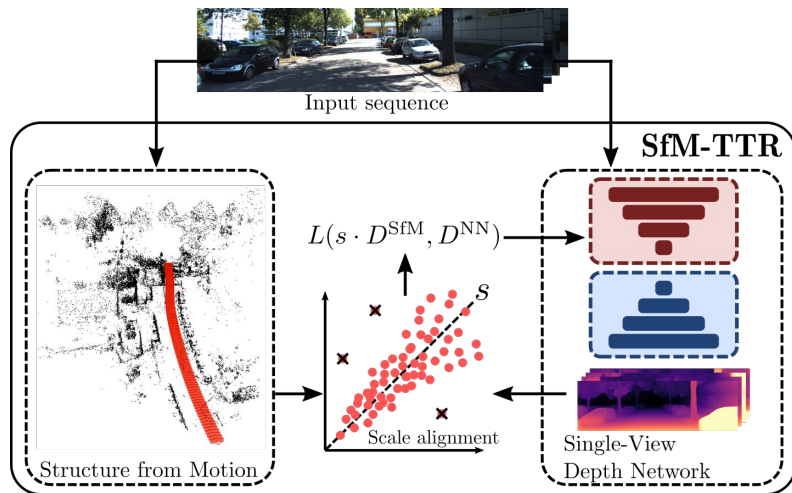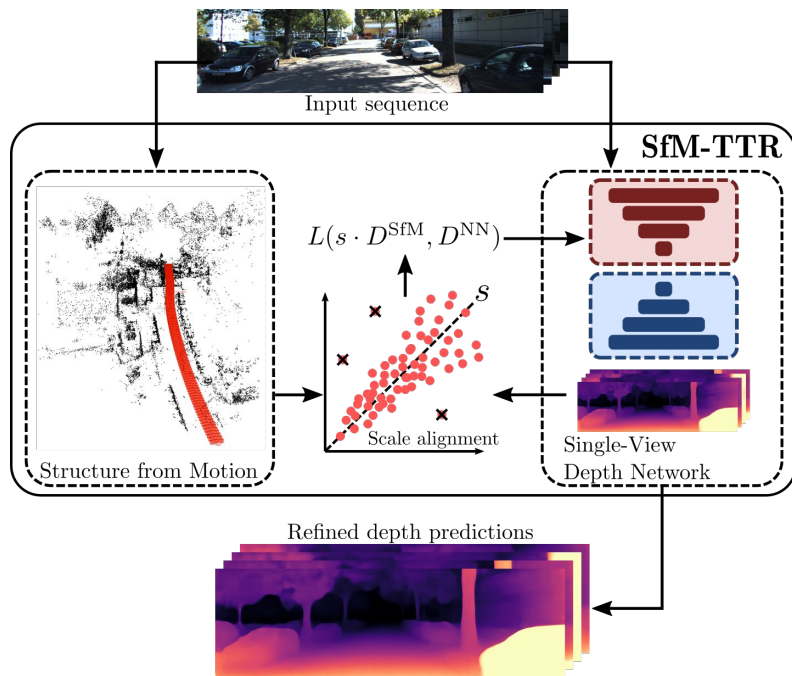$$\mathcal{L} = \frac{1}{|\mathcal{D}_j^{\mathrm{SfM}}|} \sum_l w_{l,j}^{\boldsymbol{\theta}} \| \hat{s} \cdot D_{l,j}^{\mathrm{SfM}} - D_{l,j}^{\mathrm{NN}} \|_1$$

5. Optimize the encoder
6. Get refined predictions

# Alignment

# Alignment

# Alignment



- Heteroscedasticity

# Alignment

- Heteroscedasticity
- Uneven distribution of points

# Alignment



- Heteroscedasticity
- Uneven distribution of points



- Many outliers (both distributions)

# Alignment

Two Steps: Strict & Relaxed Model

# Alignment
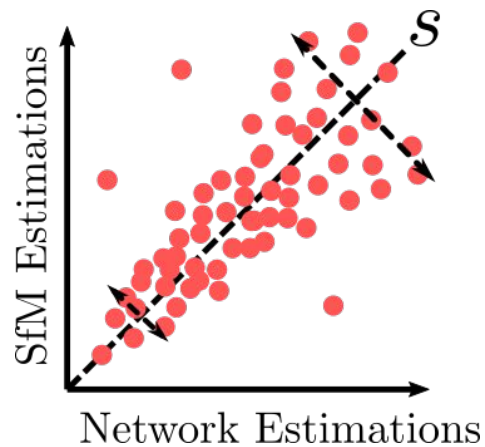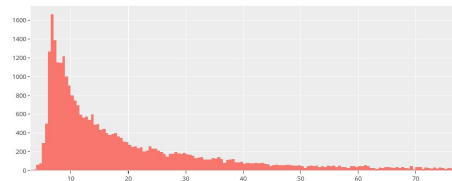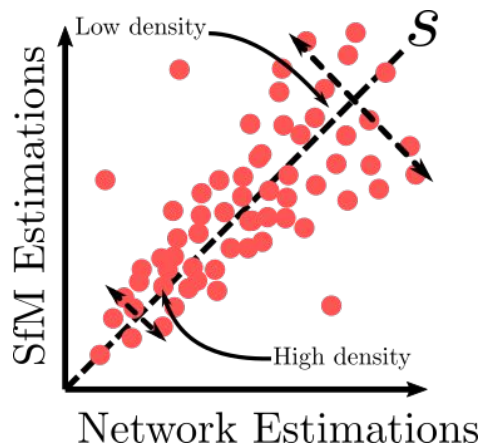


Two Steps: Strict & Relaxed Model

1. Strict model

   • RANSAC to remove outliers

   $$\frac{\left(s_{l,j} \cdot D_{l',j'}^{\text{SfM}} - D_{l',j'}^{\text{NN}}\right)^2}{s_{l,j} \cdot D_{l',j'}^{\text{SfM}}} \leq \tau$$

   • Weighted Least Squares

   $$\hat{s} = \underset{s}{\arg\min} \sum_{j} \sum_{l} w_{l,j}^{s} \left(s \cdot D_{l,j}^{\text{SfM}\checkmark} - D_{l,j}^{\text{NN}\checkmark}\right)^2$$

# Alignment



Two Steps: Strict & Relaxed Model

1. Strict model

   - RANSAC to remove outliers

$$\frac{\left(s_{l,j} \cdot D_{l',j'}^{\text{SfM}} - D_{l',j'}^{\text{NN}}\right)^2}{s_{l,j} \cdot D_{l',j'}^{\text{SfM}}} \leq \tau$$

   - Weighted Least Squares

$$\hat{s} = \underset{s}{\arg\min} \sum_j \sum_l w_{l,j}^s \left(s \cdot D_{l,j}^{\text{SfM}\checkmark} - D_{l,j}^{\text{NN}\checkmark}\right)^2$$

2. Relaxed model

   - Use $\hat{s}$ to include outliers and correct them

# Experiments

- SfM-TTR improves depth estimations

# Experiments

- SfM-TTR improves depth estimations

# Experiments

- SfM-TTR improves depth estimations

- Bigger improvement than photometric refinement

- Specially for further areas

# Experiments

# Experiments

- Improvement in different network architectures
- Improvement in self and supervised models

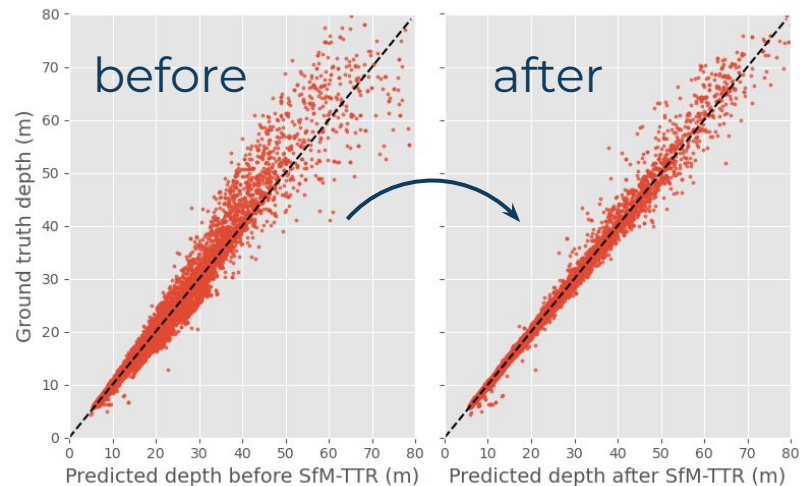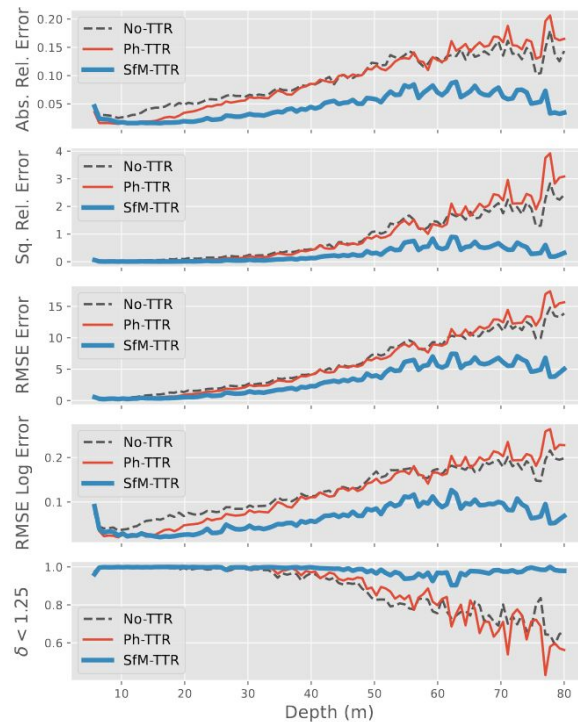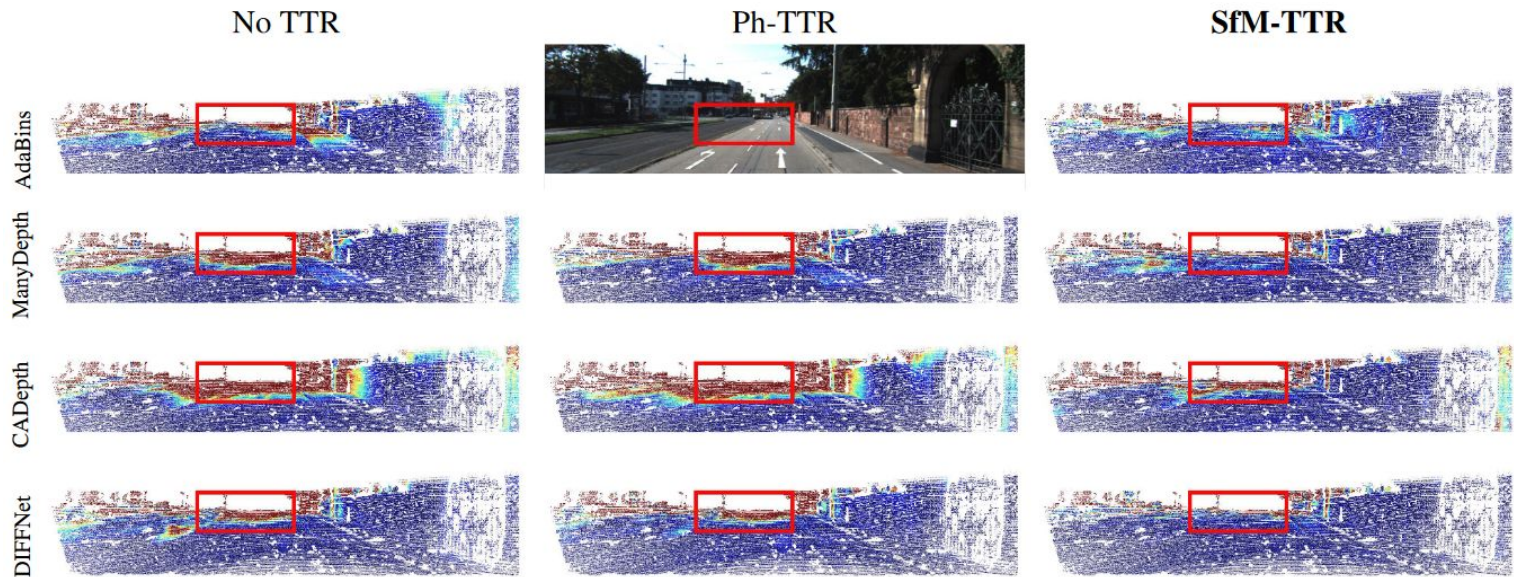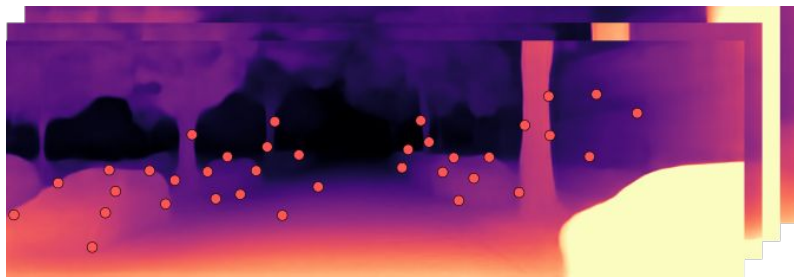| TTR | Method | Abs Rel ↓ | Sq Rel ↓ | RMSE ↓ | RMSE log ↓ | $\delta < 1.25$ ↑ | $\delta < 1.25^2$ ↑ | $\delta < 1.25^3$ ↑ |
|---|---|---|---|---|---|---|---|---|
| ✗ | AdaBins [6] ◇ † | 0.058 | 0.190 | 2.360 | 0.088 | 0.964 | 0.995 | **0.999** |
| ✓ | **AdaBins [6] + SfM-TTR †** | **0.054** | **0.138** | **1.885** | **0.078** | **0.978** | **0.996** | **0.999** |
| ✗ | ManyDepth [48] * | 0.059 | 0.297 | 2.960 | 0.097 | 0.954 | 0.991 | <u>**0.998**</u> |
| ✓ | ManyDepth [48] + Ph-TTR * | **0.053** | **0.252** | 2.774 | **0.089** | 0.962 | **0.993** | <u>**0.998**</u> |
| ✓ | **ManyDepth [48] + SfM-TTR** | 0.054 | **0.252** | **2.510** | **0.089** | **0.966** | 0.992 | <u>**0.998**</u> |
| ✗ | CADepth [51] * | 0.073 | 0.359 | 3.287 | 0.112 | 0.941 | 0.990 | **0.997** |
| ✓ | CADepth [51] + Ph-TTR * | 0.082 | 0.426 | 3.565 | 0.124 | 0.923 | 0.986 | **0.997** |
| ✓ | **CADepth [51] + SfM-TTR** | **0.060** | **0.263** | **2.620** | **0.096** | **0.962** | **0.992** | **0.997** |
| ✗ | DIFFNet [58] * | 0.066 | 0.318 | 3.078 | 0.103 | 0.953 | 0.992 | <u>**0.998**</u> |
| ✓ | DIFFNet [58] + Ph-TTR * | 0.053 | 0.252 | 2.778 | 0.090 | 0.965 | 0.993 | <u>**0.998**</u> |
| ✓ | **DIFFNet [58] + SfM-TTR** | <u>**0.052**</u> | <u>**0.229**</u> | <u>**2.444**</u> | <u>**0.085**</u> | <u>**0.973**</u> | <u>**0.994**</u> | <u>**0.998**</u> |

KITTI

# Conclusions

- SfM serves as good pseudo ground truth
- Significant improvement of the estimations
  - Specially in distant points
- Better results than other refinements



Paper+code

JUNE 18-22, 2023
CVPR
VANCOUVER, CANADA