

Project 9d: Monocular Depth Estimation via Transfer Learning

Supervisor: Yongbo Chen

Who are we?

Students from Robotics at TU Delft and Autonomous Systems at DTU

Weihaio Xuan, Ruijie Ren, Siyuan Wu, Liangchen Sui, Shaohang Han

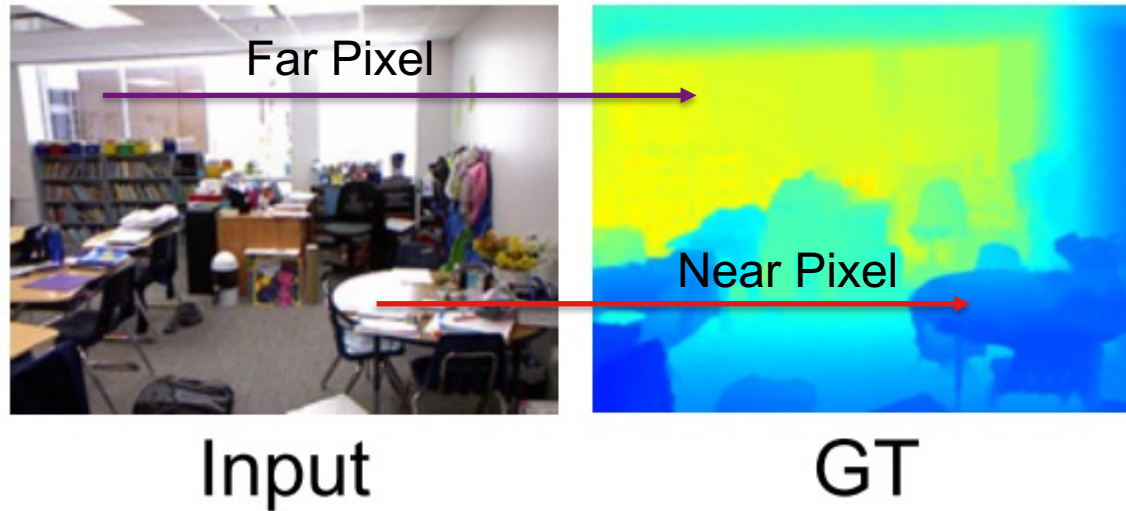


Presentation by Liangchen Sui

Project 9: Monocular Depth Estimation via Transfer Learning

Task Description:

- Pixel to pixel regression¹



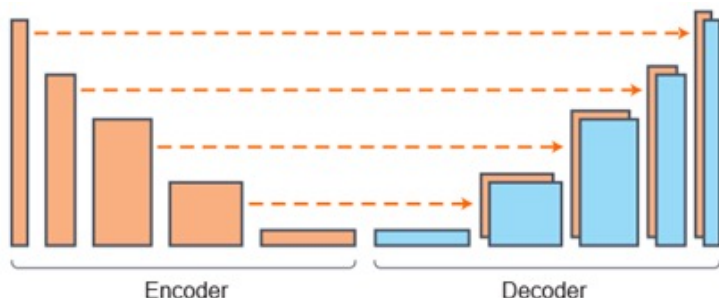
¹Alhashim, Ibraheem, and Peter Wonka. "High quality monocular depth estimation via transfer learning." *arXiv preprint arXiv:1812.11941* (2018).

Encoder & Decoder Design:

- Encoder
 - Backbone of DenseNet-169 transferred from ImageNet classification
- Decoder
 - Successive bilinear upsampling blocks followed by shallow convolutional layers to realize pixel-pixel regression (high complexity proved to backfire)
 - 1x1 Convs to establish inter-channel connections
 - Skip connections to retrieve information and add high-resolution details

Project 9: Monocular Depth Estimation via Transfer Learning

Encoder & Decoder Design:



LAYER	OUTPUT	FUNCTION
INPUT	$480 \times 640 \times 3$	
CONV1	$240 \times 320 \times 64$	DenseNet CONV1
POOL1	$120 \times 160 \times 64$	DenseNet POOL1
POOL2	$60 \times 80 \times 128$	DenseNet POOL2
POOL3	$30 \times 40 \times 256$	DenseNet POOL3
...
CONV2	$15 \times 20 \times 1664$	Convolution 1×1 of DenseNet BLOCK4

Encoder

UP1	$30 \times 40 \times 1664$	Upsample 2×2
CONCAT1	$30 \times 40 \times 1920$	Concatenate POOL3
UP1-CONVA	$30 \times 40 \times 832$	Convolution 3×3
UP1-CONVB	$30 \times 40 \times 832$	Convolution 3×3
UP2	$60 \times 80 \times 832$	Upsample 2×2
CONCAT2	$60 \times 80 \times 960$	Concatenate POOL2
UP2-CONVA	$60 \times 80 \times 416$	Convolution 3×3
UP2-CONVB	$60 \times 80 \times 416$	Convolution 3×3
UP3	$120 \times 160 \times 416$	Upsample 2×2
CONCAT3	$120 \times 160 \times 480$	Concatenate POOL1
UP3-CONVA	$120 \times 160 \times 208$	Convolution 3×3
UP3-CONVB	$120 \times 160 \times 208$	Convolution 3×3
UP4	$240 \times 320 \times 208$	Upsample 2×2
CONCAT3	$240 \times 320 \times 272$	Concatenate CONV1
UP2-CONVA	$240 \times 320 \times 104$	Convolution 3×3
UP2-CONVB	$240 \times 320 \times 104$	Convolution 3×3
CONV3	$240 \times 320 \times 1$	Convolution 3×3

Decoder

Loss Function Design

- L1 Loss(Depth + Gradient) + SSIM → High Granularity
- L2 Loss(Depth + Gradient) + SSIM → Fast Convergence
- Berhu Loss(Depth + Gradient) + SSIM → Make a Compromise

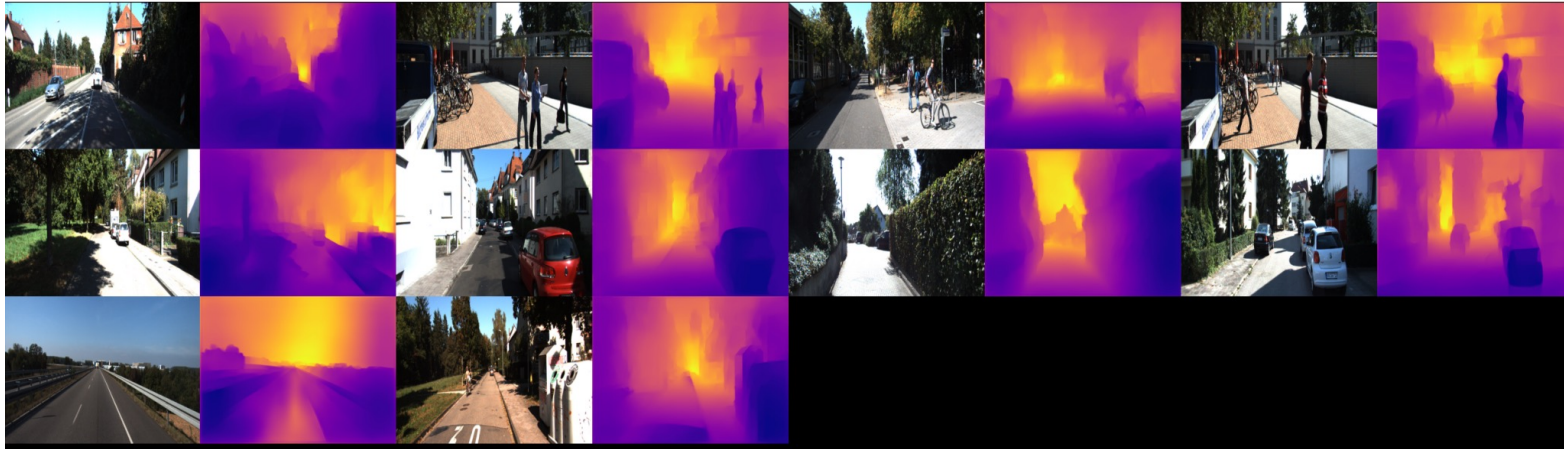
Data Augmentation

- To achieve higher generalizations, series of operations simulating natural deformations are implemented on the dataset
 - Random crops
 - Horizontal flips
 - Vertical flips
 - Rotations
 - Affine transformations

Project 9: Monocular Depth Estimation via Transfer Learning

Results

- RGB Inputs & Depth Predictions(fine-tuned on KITTI)



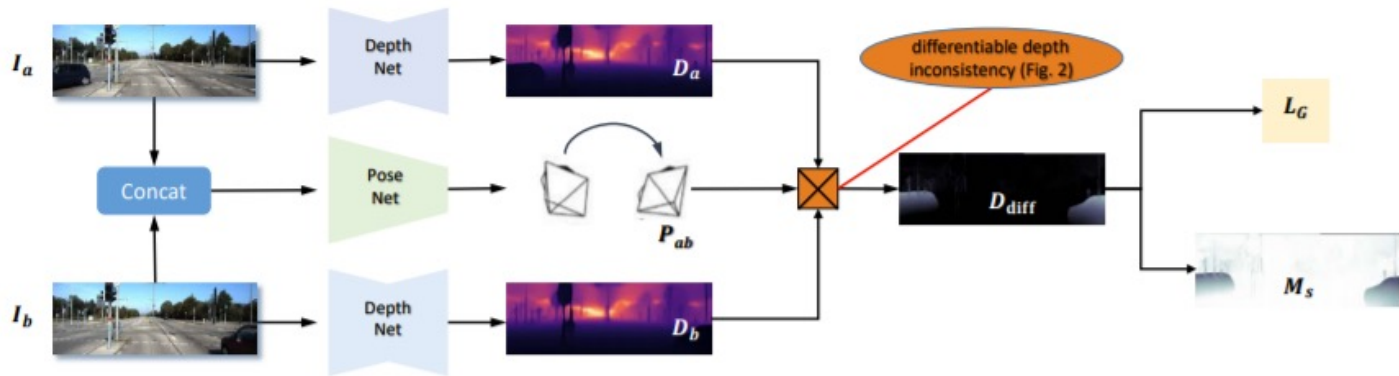
- Evaluation metrics are regrettably not available due to the key files missing in the benchmark

Future Work

- An interesting discovery: specific scenarios such as **outdoors** and **over exposures** are not considerably collected in NYU Depth v2 dataset, which could jeopardize the **robustness** and **generalization** of the model under some extreme circumstances, especially when the **depth values are large and light intensities are high**.

Future Work

- The structure of ideal model based on supervised learning¹ could be used as Depth Net component in unsupervised scale-consistent depth learning from video, which tends to serve as an important part in visual SLAM to provide depth information.



¹Bian, Jia-Wang, et al. "Unsupervised Scale-consistent Depth Learning from Video." *International Journal of Computer Vision* (2021): 1-17.

Thank you for listening!