

Lab CudaVision

Learning Vision Systems on Graphics Cards (MA-INF 4308)

Lab Project

20.07.2021

PROF. SVEN BEHNKE, ANGEL VILLAR-CORRALES

Contact: villar@ais.uni-bonn.de

Stereo Depth Estimation

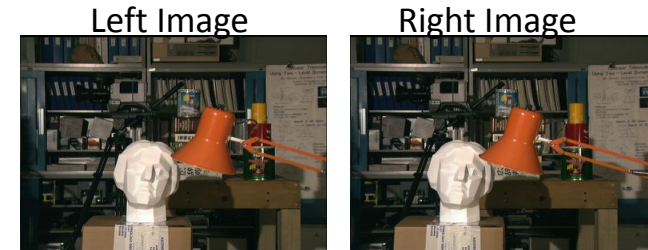
Depth Estimation

- Estimating distance and depth
 - Trivial for humans
 - Very challenging for machines
- Approaches
 - Radar/Lidar
 - Epipolar geometry
 - Deep Learning (mono or stereo)

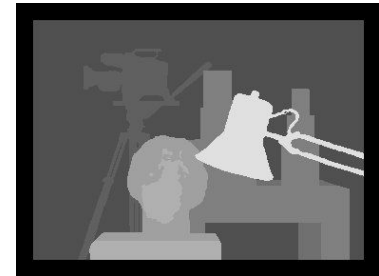


Stereo Depth Estimation

- Using information from left and right cameras
 - Imitates human physiology
- Disparity:** distance between two corresponding points in the left and right image of a stereo pair



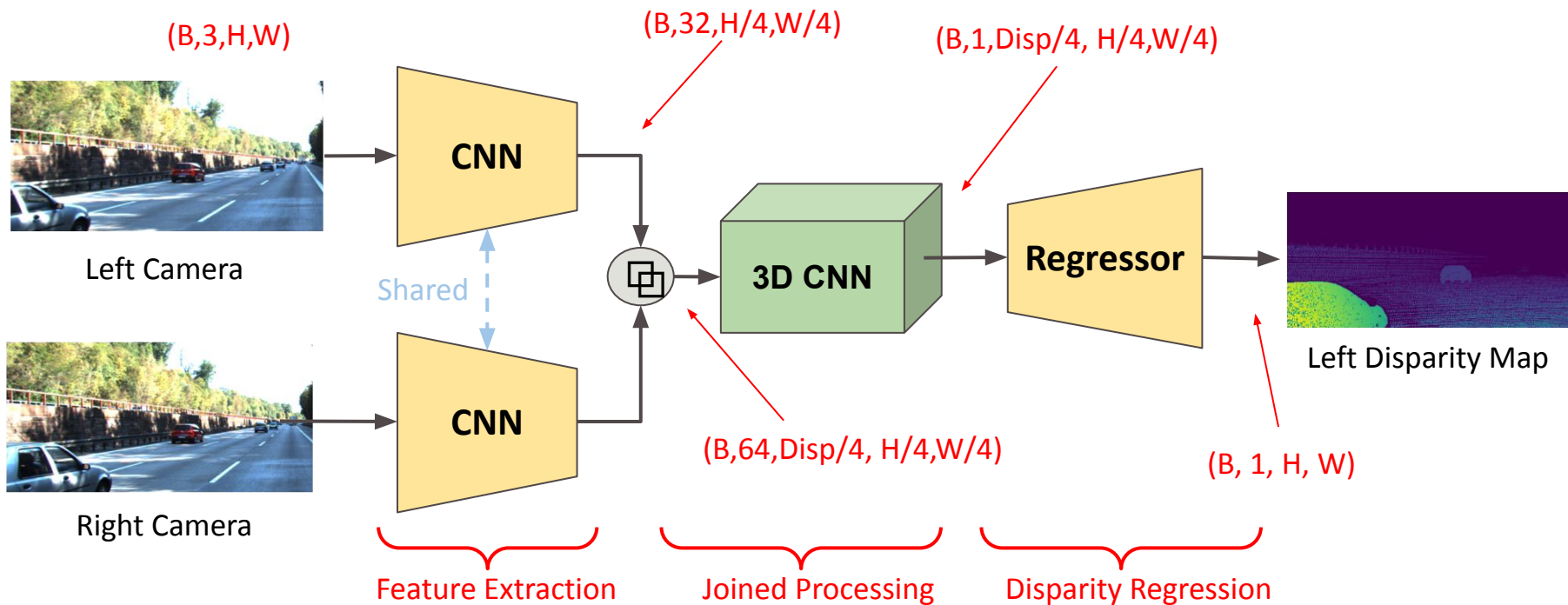
$$\text{Depth (mm)} = \frac{\text{focal length(pixels)} \times \text{Baseline(mm)}}{\text{Disparity (pixels)}}$$



Model

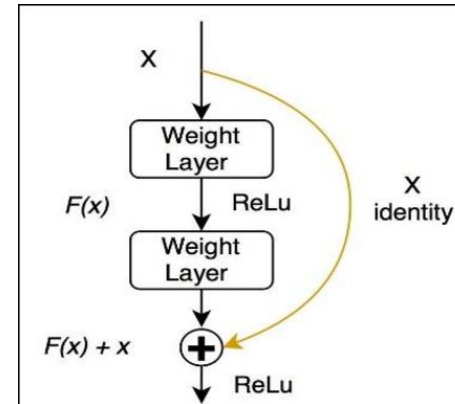
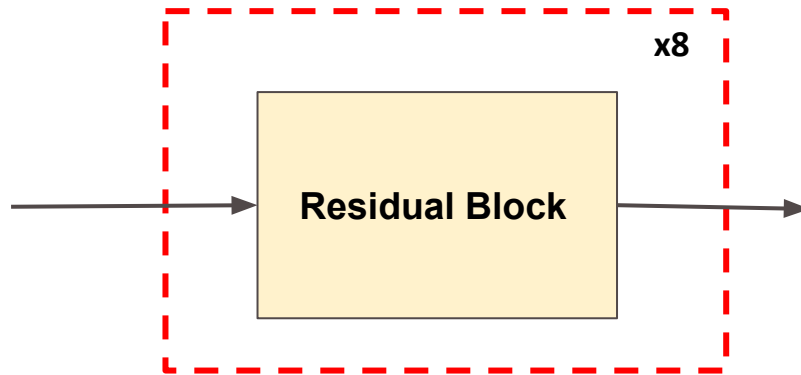
6

Baseline Model Pipeline



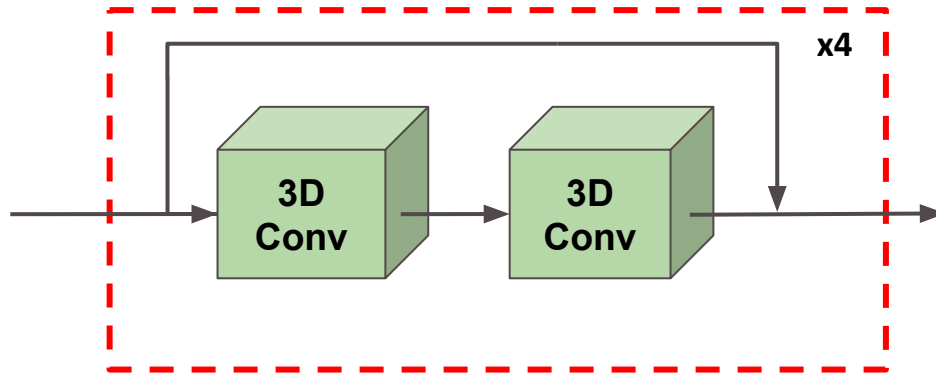
Feature Extraction

- Siamese residual CNN
 - Shared weights for left and right images
- Cascade of residual blocks (e.g., 8 blocks)



Joined Processing

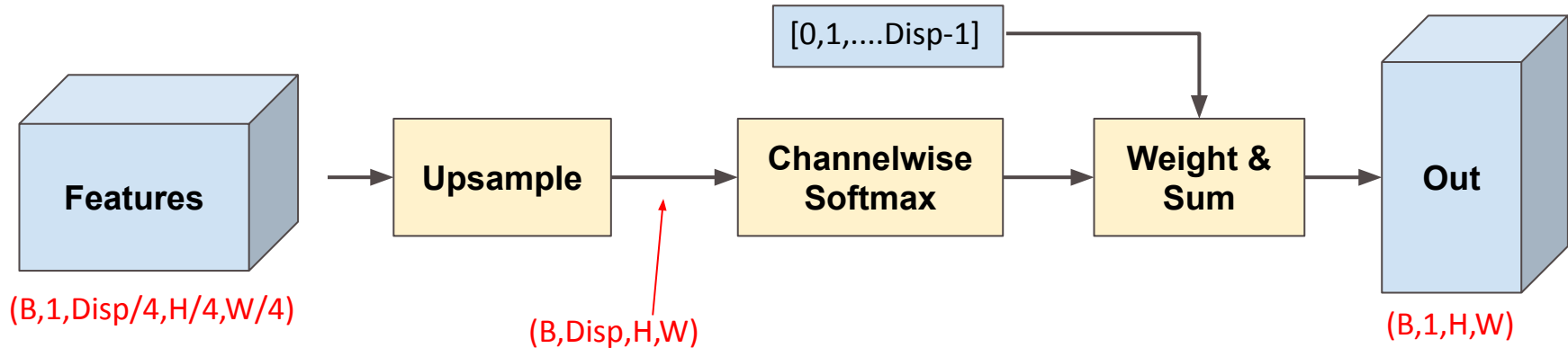
- Concatenate features from both views into a cost volume
 - $(B, 2 \cdot C, \text{Disp}/4, H/4, W/4)$
- Process volume with 3D-Convolutions



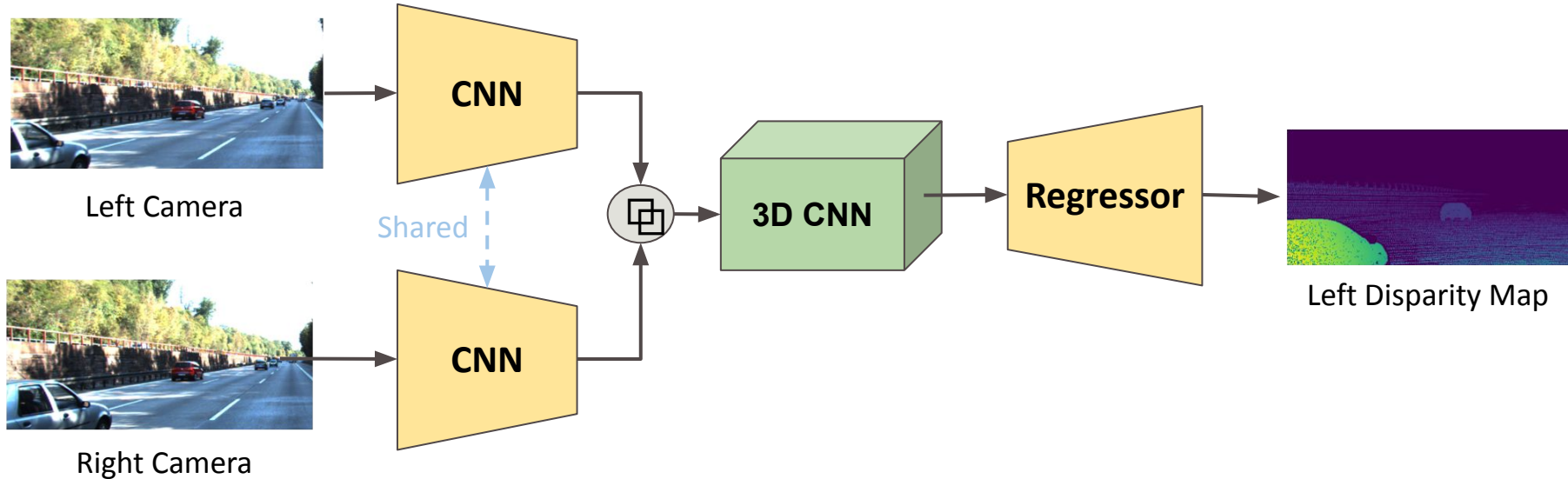
```
"""
making cost volume: concatenating features across channel dimension
for each disparity level
"""
cost = torch.Tensor(B, C*2, self.max_disp//4, H//4, W//4).to(device)
for i in range(self.max_disp // 4):
    if(i == 0):
        cost[:, :C, i, :, :] = left_feats
        cost[:, C:, i, :, :] = right_feats
    else:
        cost[:, :C, i, :, i:] = left_feats[:, :, :, i:]
        cost[:, C:, i, :, i:] = right_feats[:, :, :, i:]
```

Disparity Regression

- Upscales the volumetric features
- Performs a **soft-regression** of the disparity values



Baseline Model Pipeline



Datasets

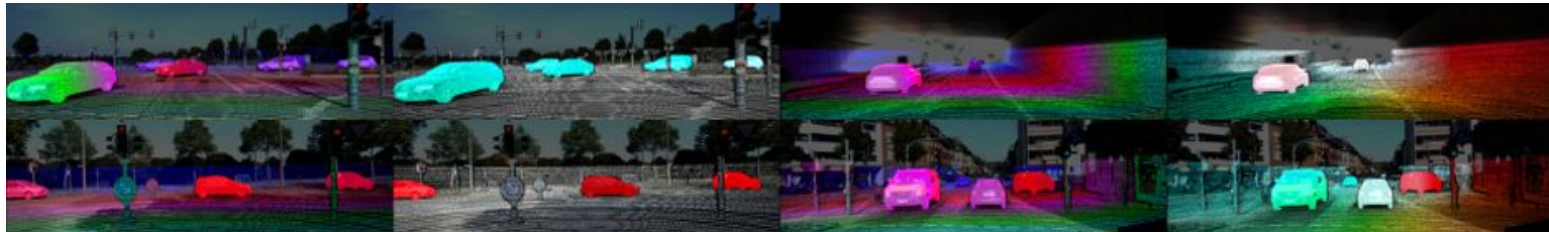
SceneFlow Dataset

- Synthetic dataset
 - 35,454 training images
 - 4,370 test images
 - Dense disparity maps
- Use for pretraining model
- <https://lmb.informatik.uni-freiburg.de/resources/datasets/SceneFlowDatasets.en.html>



KITTI-2015 Dataset

- Real-world autonomous driving dataset
- Small size for stereo depth:
 - 200 stereo pairs with sparse disparity maps
 - 150/50 training and evaluation split
- http://www.cvlibs.net/datasets/kitti/eval_scene_flow.php?benchmark=stereo



Training & Evaluation

Train/Eval on KITTI

- Train with image crops of size: (3, 256, 512)
- Evaluation on original size: \approx (3, 376, 1240) with batch size of 1
- Predict left disparity map with occluded pixels (*dips_occ*)
- Dataset splits:
 - First 150 image pairs for training
 - Final 50 images for evaluation
- Maximum disparity: Disp=192
- *SmoothL1* loss function for training
 - Considering only pixels with non-zero disparity

Evaluation Metrics

- Only evaluate pixels where true disparity is non-zero
- **SmoothL1** regression loss on evaluation set ($\beta=1$)
- **3-pixel error (3PE)**: Percentage of pixels where
 - disparity error is less than 3 pixels
 - error is less than 5% of the true disparity

$$\text{Smooth } \ell_1 = \frac{1}{N} \sum_i^N l_i$$

$$l_i = \begin{cases} \frac{1}{2\beta} (\mathbf{X}_i - \hat{\mathbf{X}}_i)^2 & |\mathbf{X}_i - \hat{\mathbf{X}}_i| \leq \beta \\ |\mathbf{X}_i - \hat{\mathbf{X}}_i| - 0.5 \cdot \beta & \text{otherwise} \end{cases}$$

$$3PE(\mathbf{X}, \hat{\mathbf{X}}) = 1 - \frac{1}{N} \sum_{i=1}^N \text{CorrectDisp}(\hat{x}_i, x_i)$$

$$\text{CorrectDisp}(\hat{x}_i, x_i) = \begin{cases} 1 & ; \text{if } |x_i - \hat{x}_i| < 3 \\ 1 & ; \text{if } |x_i - \hat{x}_i| < 0.05 \cdot x_i \\ 0 & ; \text{otherwise} \end{cases}$$

Project Goals and Deliverables

Passing Requirements

1. Implement model, pipelines and utils
2. Beat a weak and simple baseline
 - a. No pretraining
 - b. Little parameter optimization
3. Create overview notebook
4. Write project report

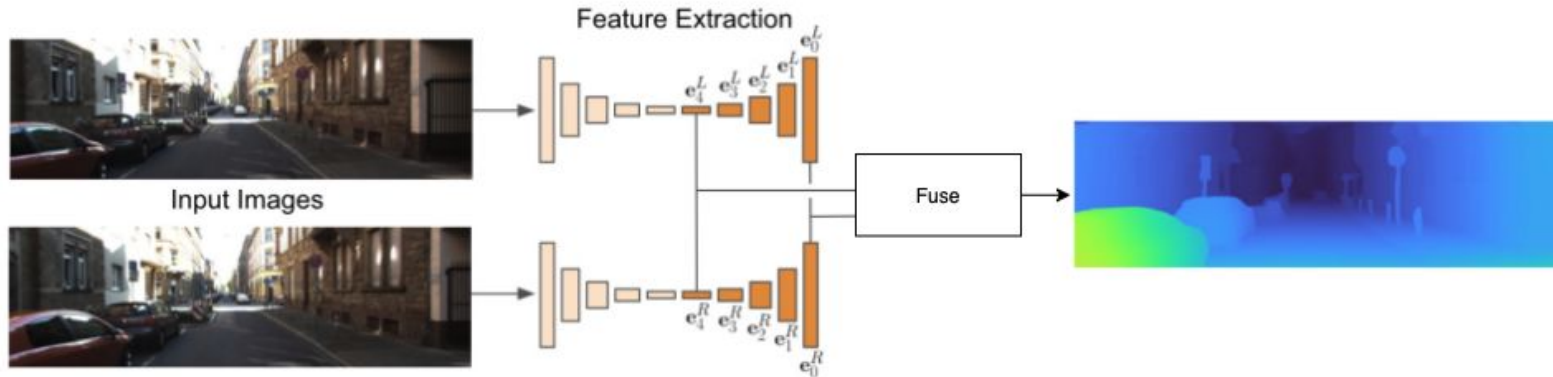
	Mean Loss	Mean 3PE
Baseline	1.1	9.0%

Improvement Ideas

- Test and debug your code
- KITTI dataset is very small. I strongly recommend:
 - Pretrain on SceneFlow
 - Fine-tune on KITTI with data augmentation
- Tweak the model
 - Use a pretrained feature extractor (e.g., ResNet18)
 - Change modules (num. layers, num. kernels, ...)
 - Slightly change the building blocks
- Hyper-parameter and training optimization
 - Optimize for LR, optimizer, batch size, ...
 - Scheduling, learning rate warmup

If Motivated :)

- Implement a model with the following architecture
 - Shared Hourglass/UNet backbone
 - Feature fusion
- Try to obtain the best results possible



Deliverables

- Complete codebase
 - Clean and structured
 - Not just a notebook!
- Trained model checkpoint and tensorboard logs
- Overview notebook (.ipynb & .html) showing main functionalities:
 - Load data
 - Load pretrained model
 - Display some results
- Project report

Project Report

- Document your work in the project report
- Try to be brief, but readable and informative
- Include figures and tables
- Use *BibTex* for the references
- I expect 6-10 pages
- Use the following template
 - <https://www.overleaf.com/read/tmnvhrsdmjrp>

Important Dates

- **20.07:** Starting date
- **03.09:** Draft submission due
- **30.08-15.09:** Revision session
- **30.09:** Final submission:

Questions?



References

1. http://www.cvlibs.net/datasets/kitti/eval_scene_flow.php?benchmark=stereo
2. <https://lmb.informatik.uni-freiburg.de/resources/datasets/SceneFlowDatasets.en.html>
3. Geiger, Andreas, et al. "Vision meets robotics: The kitti dataset." The International Journal of Robotics Research 32.11 (2013): 1231-1237.
4. Mayer, Nikolaus, et al. "A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.
5. Luo, Wenjie, Alexander G. Schwing, and Raquel Urtasun. "Efficient deep learning for stereo matching." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.
6. Kendall, Alex, et al. "End-to-end learning of geometry and context for deep stereo regression." Proceedings of the IEEE International Conference on Computer Vision. 2017.
7. Chang, Jia-Ren, and Yong-Sheng Chen. "Pyramid stereo matching network." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2018.

References

8. Smolyanskiy, Nikolai, Alexey Kamenev, and Stan Birchfield. "On the importance of stereo for accurate depth estimation: An efficient semi-supervised deep neural network approach." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops. 2018.
9. Garg, Divyansh, et al. "Wasserstein distances for stereo disparity estimation." arXiv preprint arXiv:2007.03085 (2020).
10. Menze, Moritz, and Andreas Geiger. "Object scene flow for autonomous vehicles." Proceedings of the IEEE conference on computer vision and pattern recognition. 2015.

