

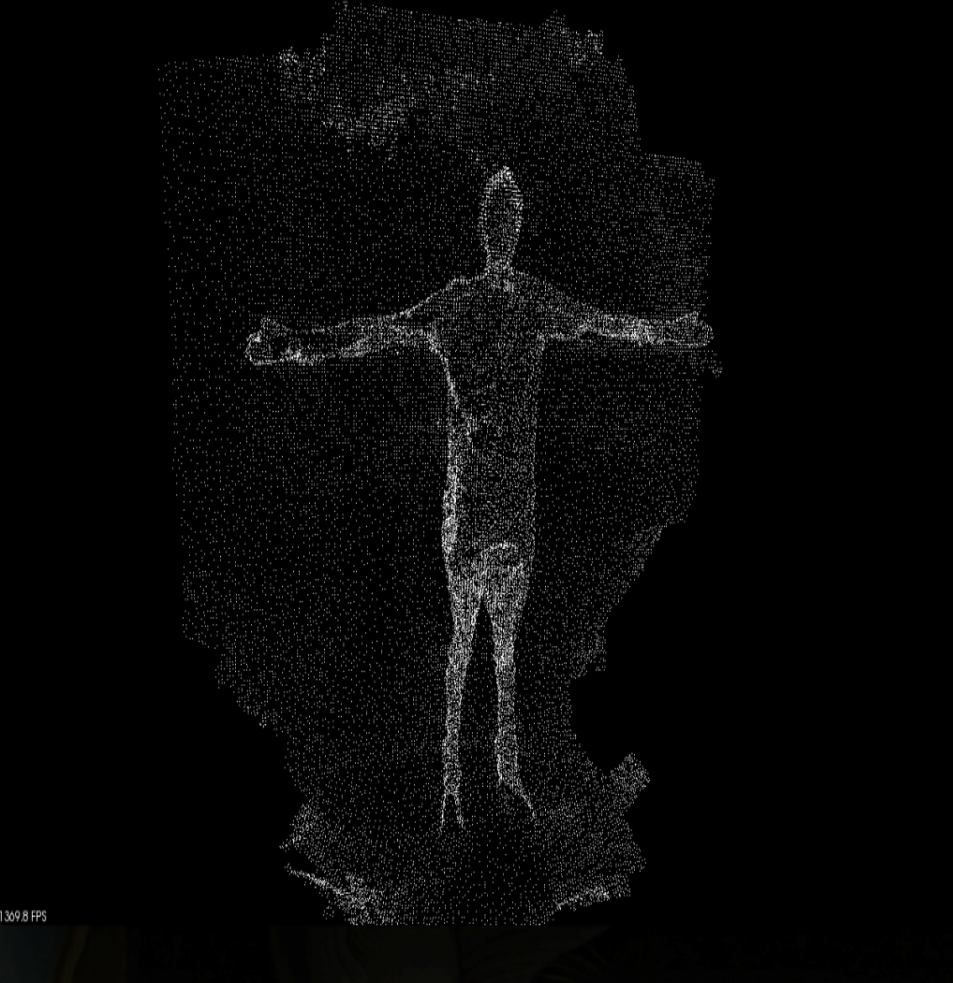


RGB-D Fusion for Real Time Object Detection

Chayan Patodi, Nikhil Mehra and Raghav Nandwani



UNIVERSITY OF
MARYLAND



Motivation

Exploring RGB+Depth Fusion for Real-Time Object Detection

By: Tanguy Ophoff, Kristof Van Beeck and Toon Goedemé

Objective: Whether fusion of depth data with the RGB data can help increase the performance of current state of the art single shot networks.

Model Architecture

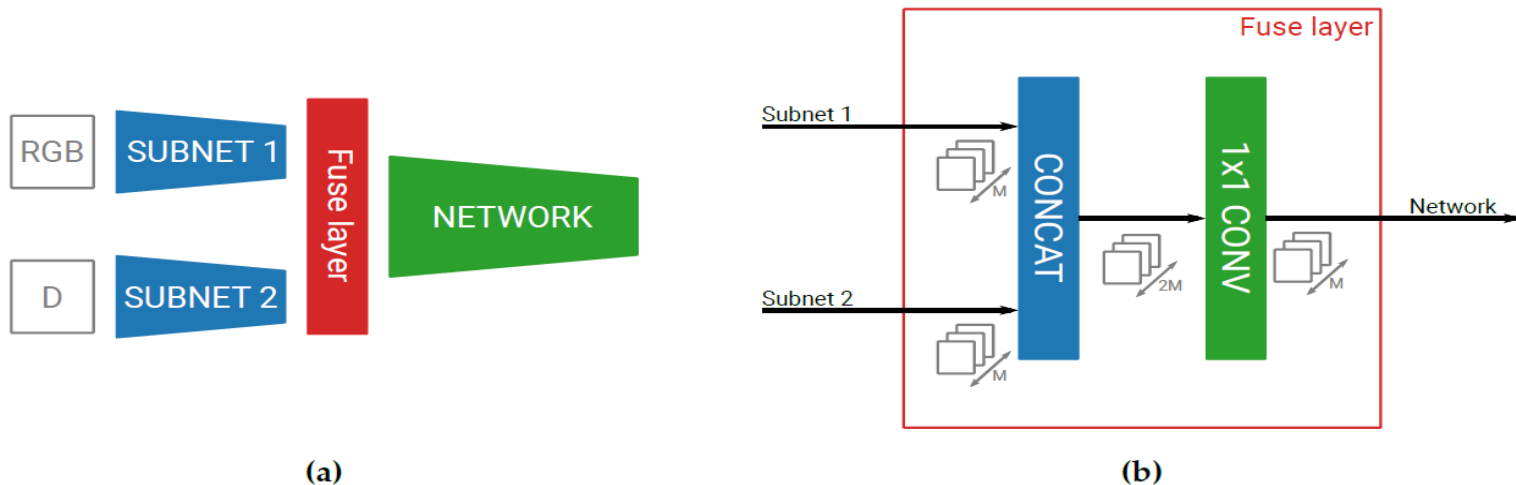


Figure 1. The main building blocks of our parameterizable fusion network. (a) The fuse layer can be transparently implemented after any arbitrary layer, allowing for a parameterizable fusion level. (b) The fuse layer combines both information streams and divides the number of output channels by two.

Network training

- Used ImageNet pretrained weights instead of random initialization
- Same weights for depth network, just removed the weights of first layer

Network training

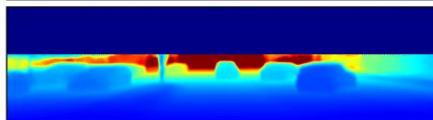
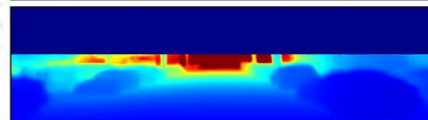
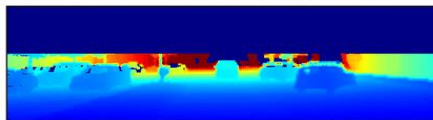
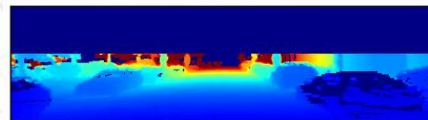
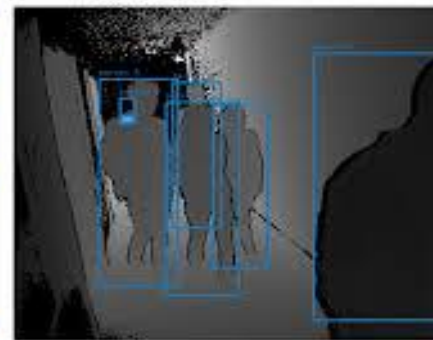
- Used ImageNet pretrained weights instead of random initialization
- Same weights for depth network, just removed the weights of first layer

Why ??

- The networks looks for similar features in both RGB and Depth images
- If the depth subnetwork does not provide any substantial information compared to the RGB network, the fusion layer could possibly ignore those feature maps.

Dataset

- EPFL pedestrian depth dataset
- KITTI depth map



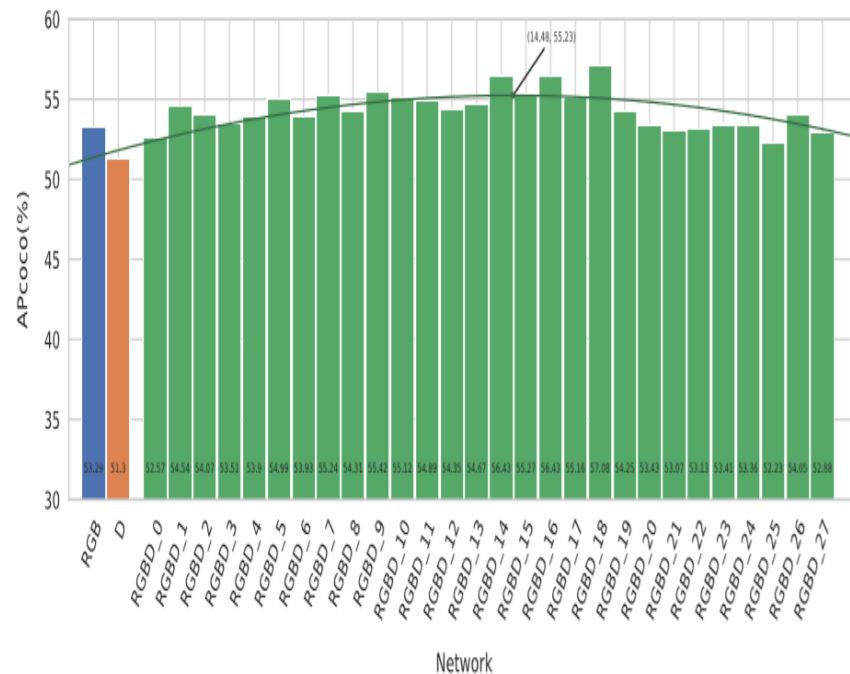
Evaluation

The main advantages of fusing depth data is in that the clearly distinguishable **silhouettes** in the **depth maps** allow for more accurate bounding boxes

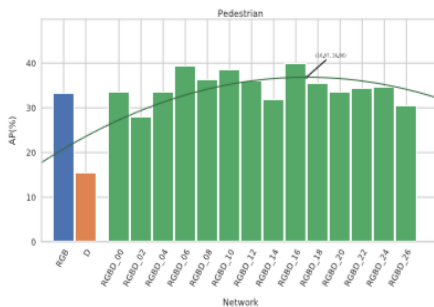
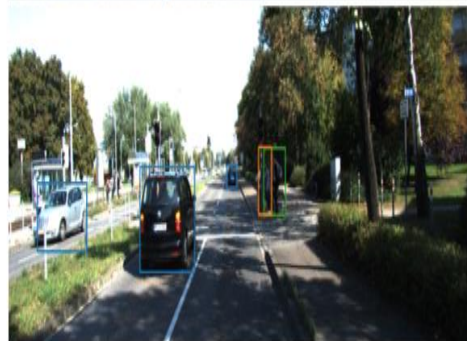
AP of the networks was measured using COCO IoU thresholding scheme, which is defined as follows.

$$AP = \frac{\sum_{IoU \in I} AP_{IoU}(Annotations, Detections)}{I} ; I = \{0.50, 0.55, 0.60, \dots, 0.95\}$$

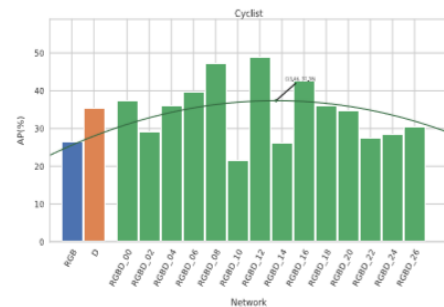
Results



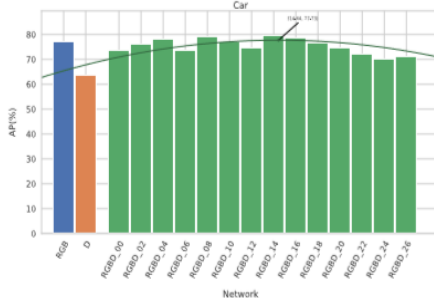
Results Continued



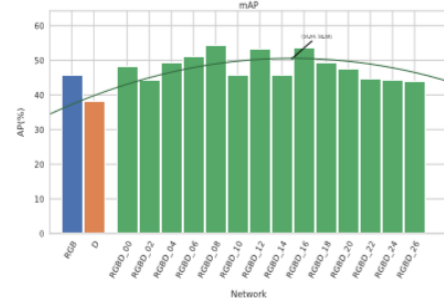
(a)



(b)

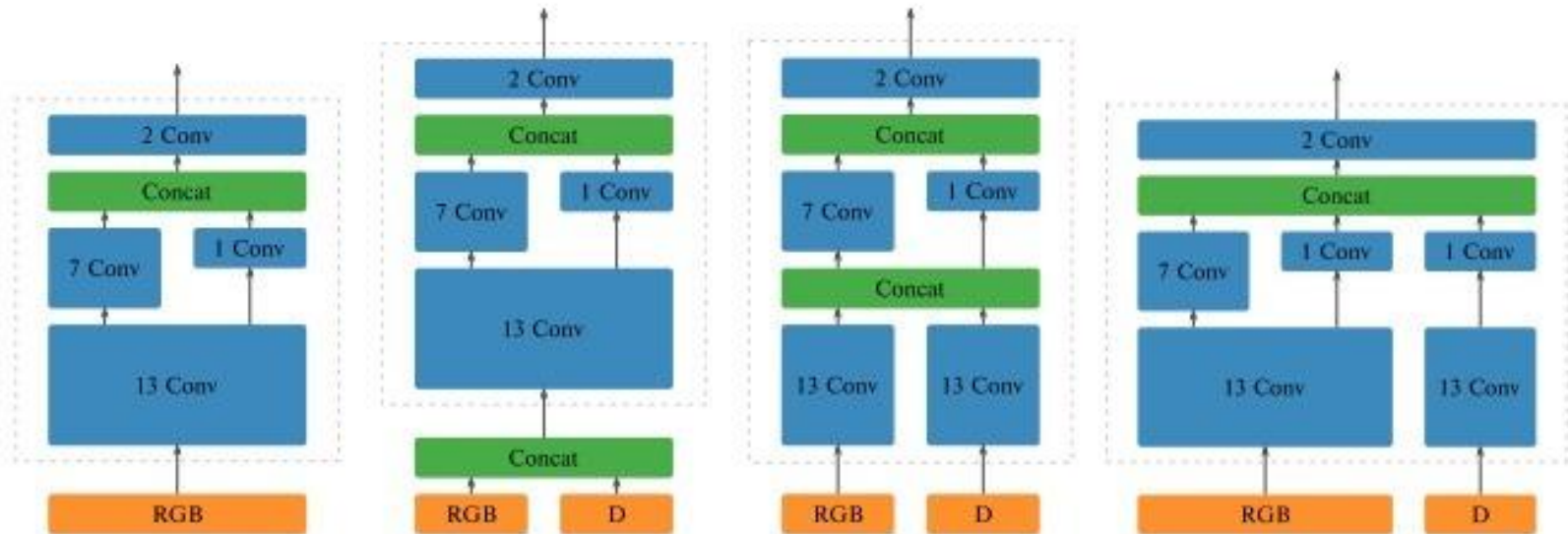


(c)



(d)

Our Approach



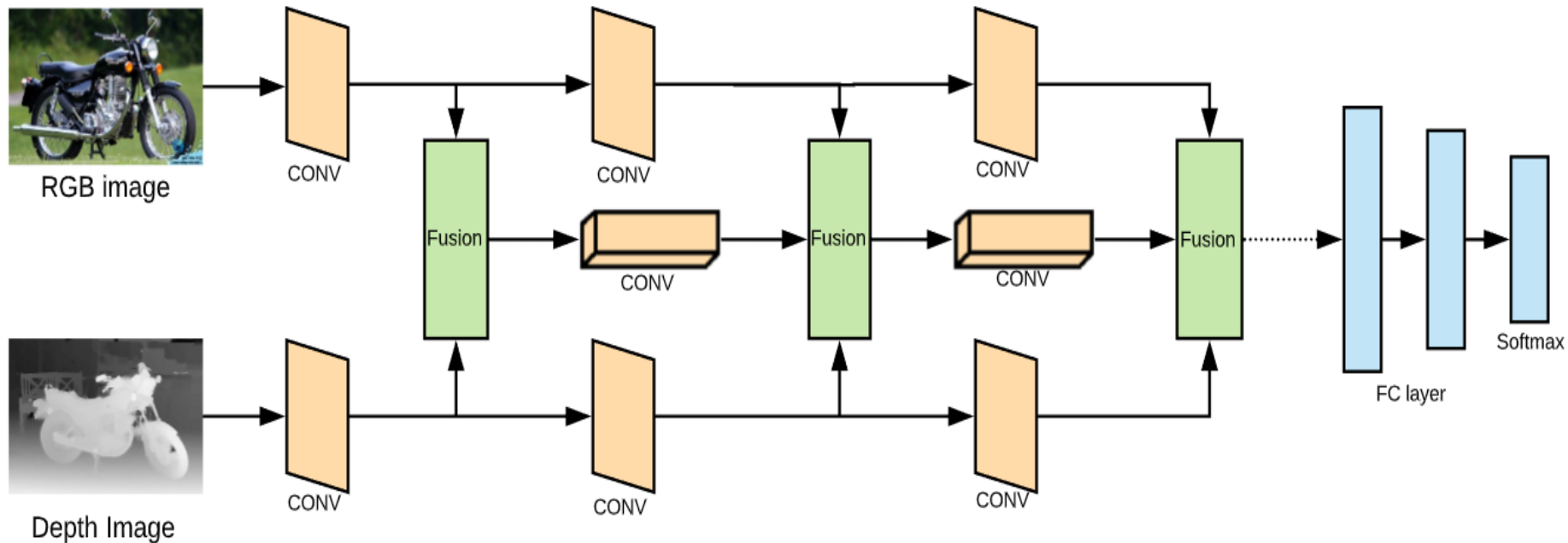
(a) YOLOv2

(b) YOLOv2 with Early Fusion

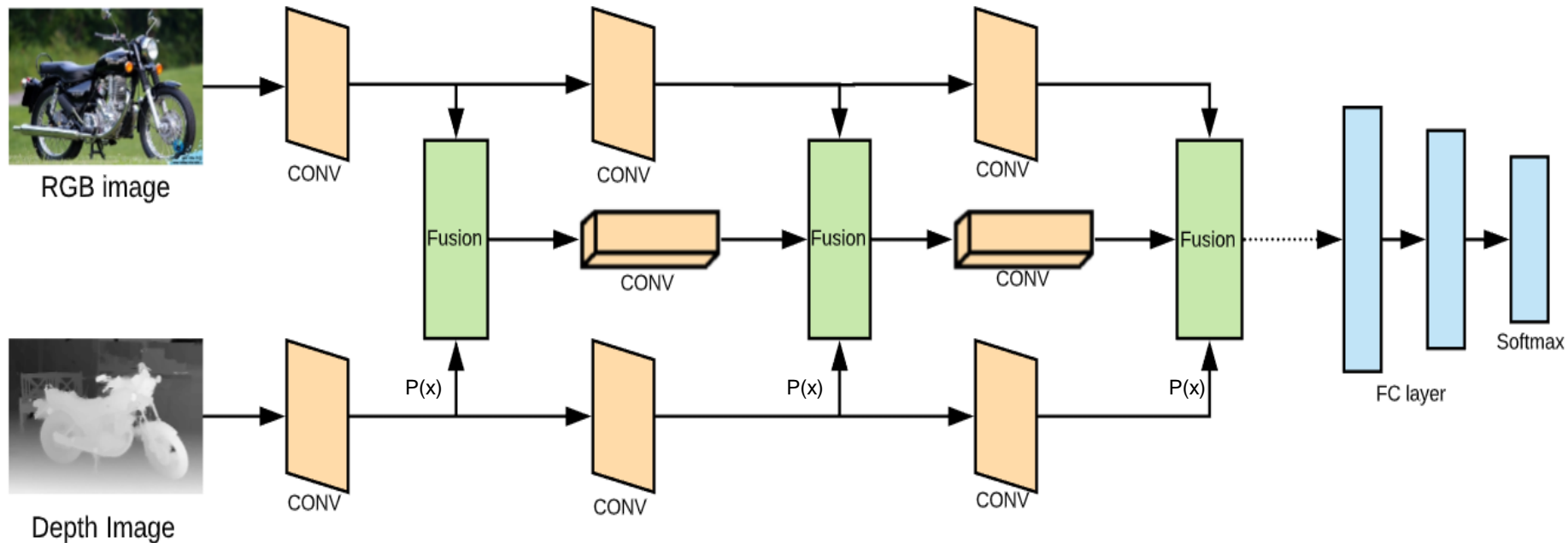
(c) YOLOv2 with Mid Fusion

(d) YOLOv2 with Late Fusion

Proposed Architecture



Proposed Arch. Continued



Expectation

- Time taken for a detection will increase, as the number of fusion layer increases.
- Difficulty in training as the number of parameters increases.
- Might get increase in the accuracy/AP metric.

[illegible]



UNIVERSITY OF
MARYLAND

Thank You