# Reducing Class Confusion in Semantic Segmentation

Sharat Agarwal, PhD17005
Divyam Anshumaan, 2017147

# Introduction: Semantic Segmentation

- We aim to assign every pixel in the image a class label.
- The class label corresponds to the object the pixel is representing.
- For example, in the segmentation given below, each pixel is assigned to one of three categories; 'Background', 'Cycle', 'Person'. ([Source](#))
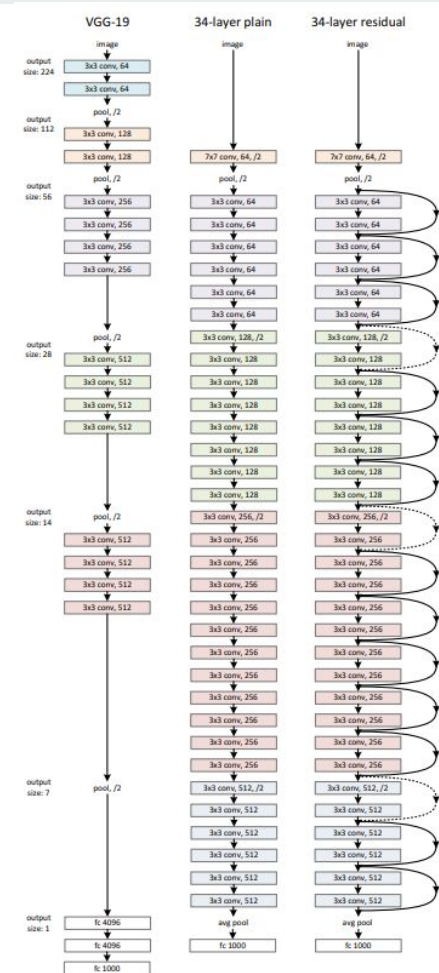


Person
Bicycle
Background

# DNN: ResNet-101

- Neural networks have certain problems that made them very difficult to train as they got deeper.
- A major problem was that of exploding or vanishing gradients, where partial derivatives calculated between the error and weights of a particular layer became very large or small as the intermediate values got multiplied.
- ResNets solved this problem by using skip connections between layers. Since the network could now be very deep, a novel bottleneck design with 1x1 convolution filters was used to reduce time complexity while preserving accuracy.
- ResNet-101 has 101 layers, and is used for feature extraction in this project.

# Dataset: Cityscapes

- Made for semantic understanding of urban street scenes. Contains pixel-level labeling for semantic segmentation.
- 30 classes, including; 'road', 'sidewalk', 'person', 'rider' etc.
- 5000 images with fine annotation, 20000 with coarse annotations.
- Diverse data. Images consist of a large number of dynamic objects. Collected from 50 European cities, over several months with varying backgrounds and scene layouts.



Example of a finely annotated segmentation from the dataset.

# Loss: Cross Entropy

- In information theory, cross-entropy between two distributions $p$ and $q$ over a set of common events measures the average number of bits required to identify an event drawn from the set if a coding scheme is used for the set is optimized for an estimated distribution $q$ rather than the true distribution $p$.
- In terms of classification, the true distribution $p(x)$ for an object $x$ would be $(y_1, y_2, \ldots, y_k, \ldots, y_n)$, where $y_k = 1$ for the true class k, and $y_i = 0$ for others. The distribution approximated by our model might predict $y_k$ to be some value less than one, while some $y_i$'s have a value greater than 0.
- Cross-entropy is minimized when the predicted distribution matches the true distribution.
- This is equivalent to minimizing the KL divergence between two distributions.
- Given by:

$$-\sum_{i=1}^{n} y_k log(p(x_k))$$

# Dice Coefficient

- For determining the goodness of the predicted segmentation, we use the Dice coefficient.
- It is given by:

$$\frac{2|X \cap Y|}{|X| + |Y|}$$

- Where $X$ and $Y$ are the set of points belonging the predicted segmentation and the true segmentation respectively.

# Proposed Pairwise Loss

- Inspired by the Dice coefficient and the limitation of cross entropy loss (not being able to penalize misclassified samples) we propose a pairwise confusion reward that can be used with cross entropy.
- Considering the classes 'Road' (R) and 'Sidewalk' (S) whose pairwise confusion we wish to maximize, the loss is given by:

$$PW_{R,S} = \frac{2 \sum R_{gt} R_{pred} P_R}{\sum R_{gt} R_{pred} (1 + P_r) + \sum R_{gt} S_{pred} P_S}$$

- We wish to maximize this term as it must reduce confusion between the 'Road' and 'Sidewalk' pixels. We expect that this in turn improves accuracy.

# Smooth Cross Entropy

- The pairwise loss is weighted with cross entropy. This helps in reducing confusion between classes that are spatially nearby, while also improving performance.
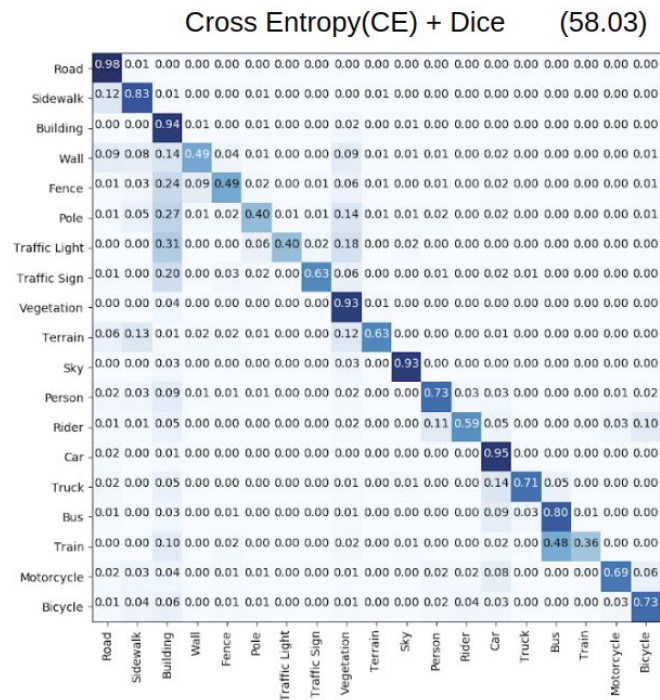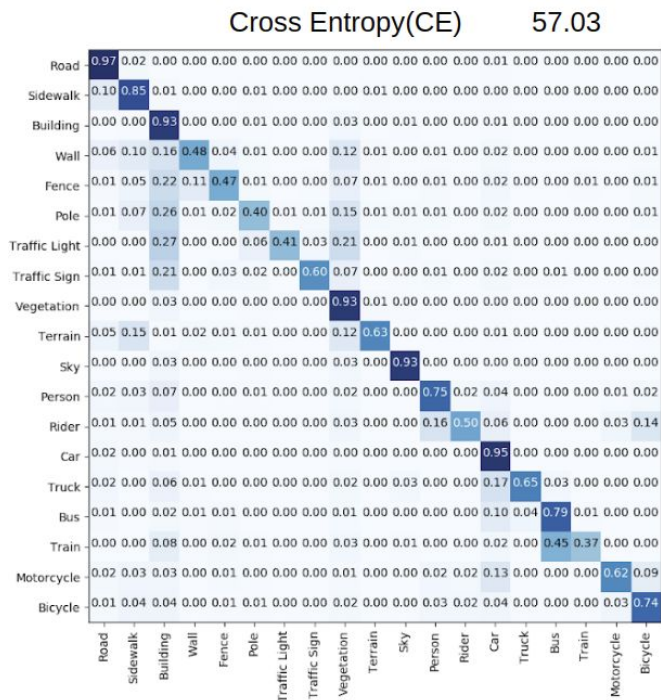- Given by:

$$\lambda_1 H(p, q) + \lambda_2 P W_{R,S}$$

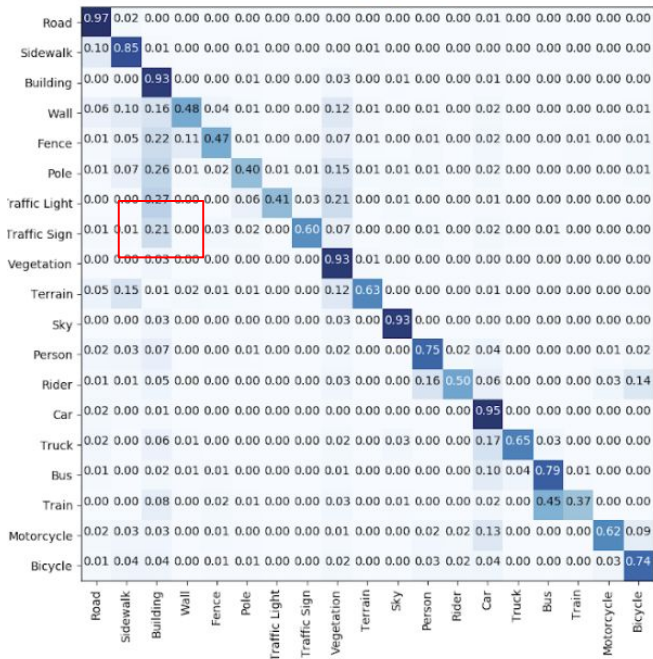- Where $\lambda_1$ and $\lambda_2$ are set to 0.8 and 0.2 respectively.

# Results

- The proposed loss is able to improve mIoU by at least 1% in all our experiments, as compared to traditional cross entropy loss.
- We reduce the confusion between classes selected for the pairwise confusion reward. For example, we were able to reduce confusion between the 'Road' and 'Sidewalk' classes from 0.21 to 0.09.
- We also managed to increase the separation between the feature vectors of classes. For example, the average distance between the feature vectors of the 'Bus' and 'Train' classes was 10.14 units while using cross entropy, but increased to 23.39 units with the proposed loss. This means that the confusion between the classes reduced.
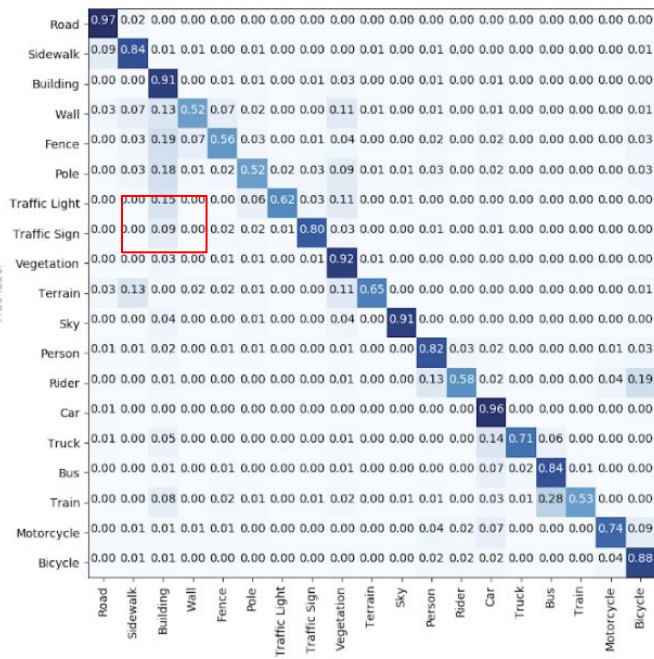
# Results



Cross Entropy(CE)     57.03

Cross Entropy(CE) + Dice     (58.03)

# Results



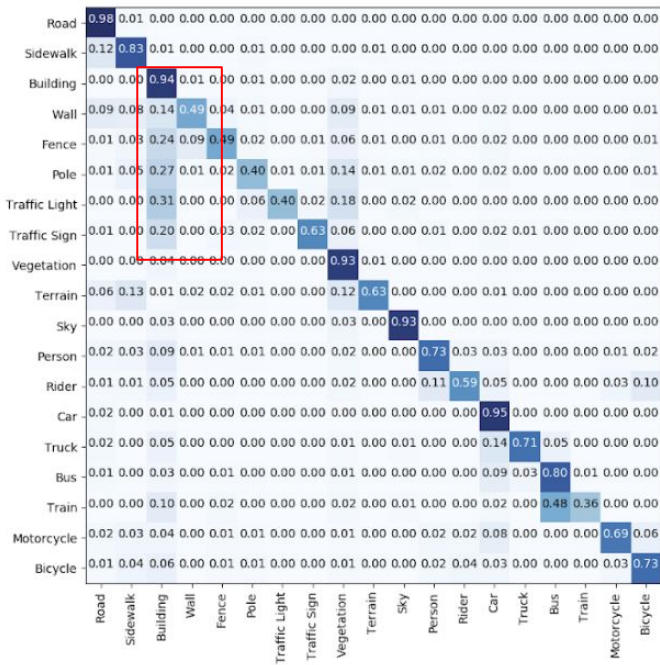Cross Entropy(CE)    57.03

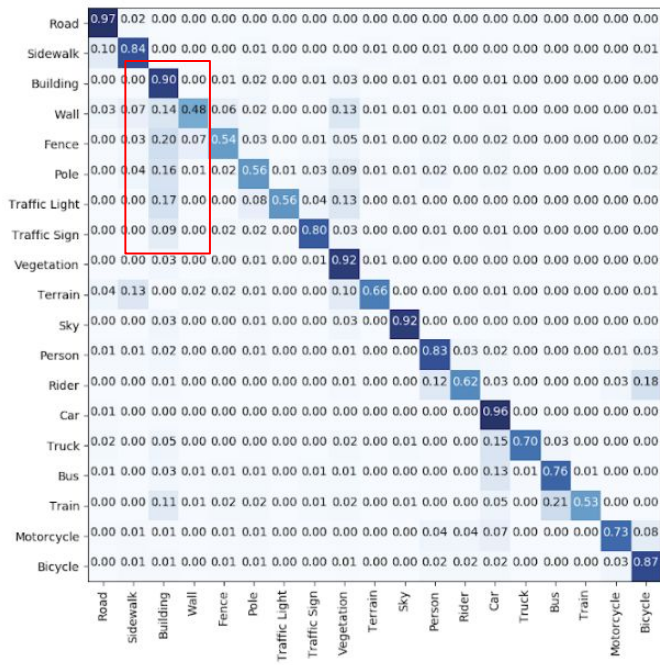(WCE) + PW_(Building , traffic sign)    58.35

# Results



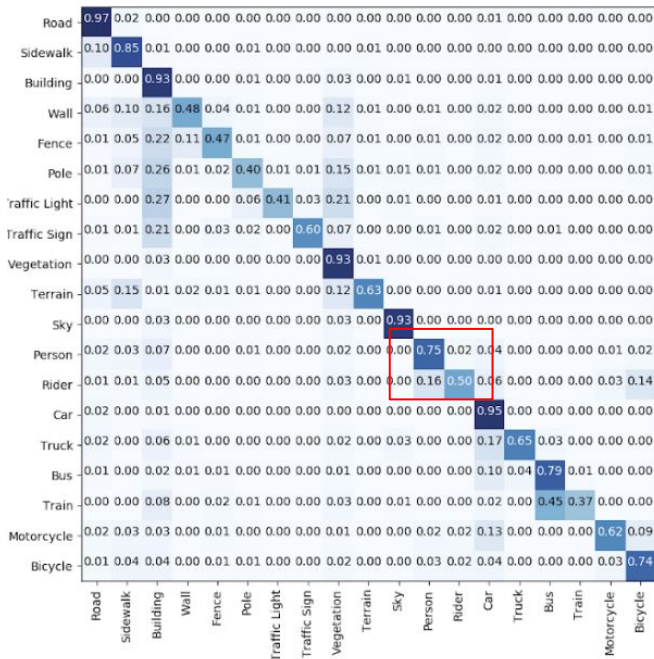Cross Entropy(CE) + Dice    (58.03)

(WCE) + PW_(2-7)    59.29
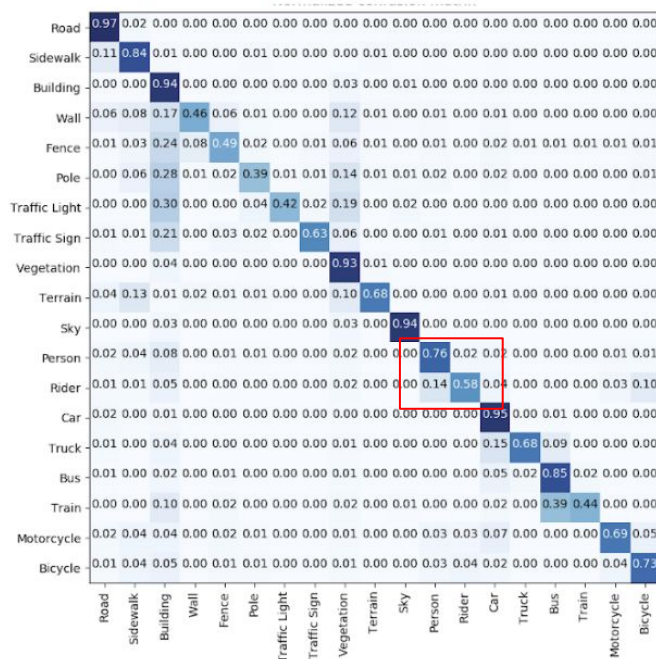
# Results



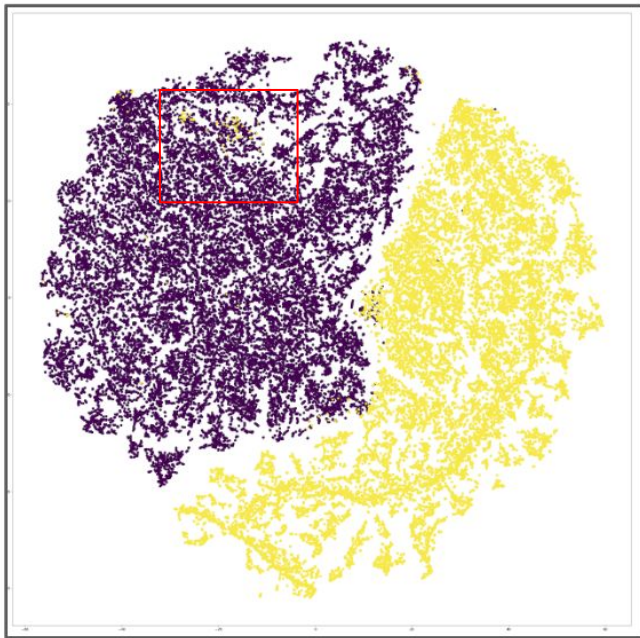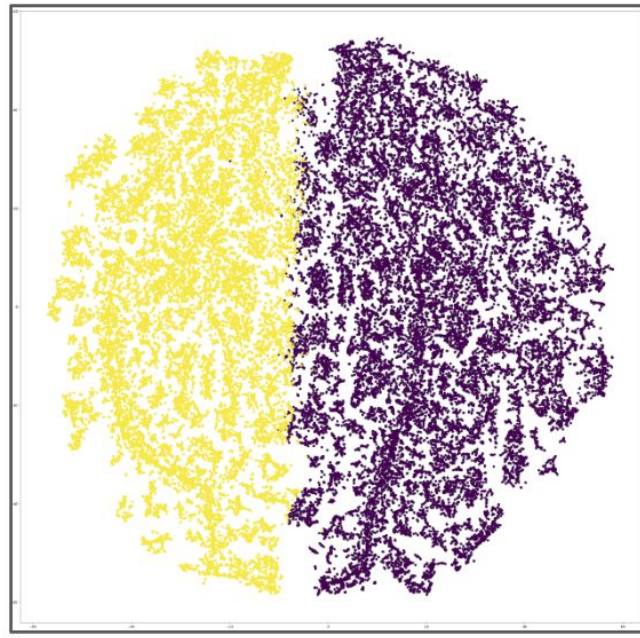Cross Entropy(CE)    57.03          (WCE) + PW_(person , rider)    58.45

# Results



CE Bus_Train    10.14

CE PW Bus_Train    23.39

# Conclusion

- We analyse the shortcomings of traditional cross entropy wherein it does not penalize misclassification made by the network.
- We also note that the loss does not consider the confusion caused by neighboring pixels during classification.
- The proposed loss attempts to incorporate these and shows that it is capable of reducing inter-class confusion and improving performance.
- However it was seen that reducing the confusion between two classes caused confusion to increase elsewhere. More research can be done regarding this in the future.