



YOLOnet

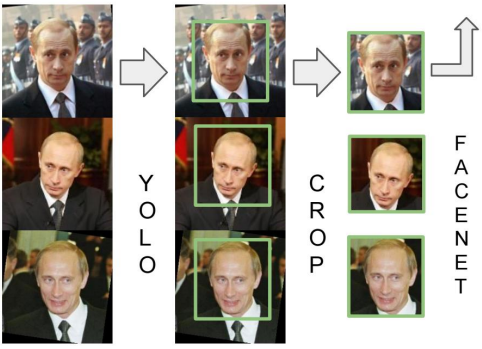
Will Lauer and Hannah DeBalsi
(wlaue, hdebalsi)@stanford.edu



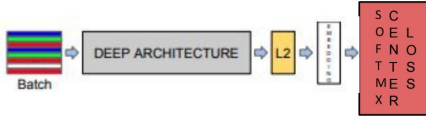
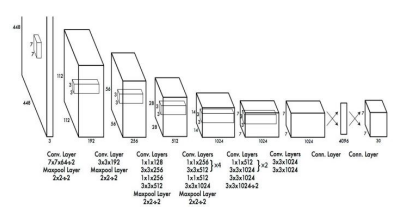
Idea

We applied two well-known and successful algorithms, YOLO and FaceNet, in an attempt to combine the speed of YOLO with the accuracy of FaceNet on facial recognition.

In our project, we relied on Allan Zelener’s Yad2k and David Sandberg’s FaceNet.



Architectures



Analysis

	Accuracy	Relative Error	Image Size
Baseline	99.2	-	250x250 pix
Crop pre-training	98.2	+1.0	166x166 pix
Crop post-training	98.7	+0.5 (50% reduction in error disparity)	166x166 pix

Data

We utilized WIDER Face^[3] images and Labeled Faces in the Wild.^[4]

For each LFW image, we created 5 randomized croppings, with each dimension 2/3 that of the original. These were split into 50%/50% train/dev sets, since we were using a pretrained FaceNet model^[2] and thus needed less new data.

Data: WIDER Face



Data: Labeled Faces in the Wild



Loss Functions

$$\lambda_{\text{coord}} \sum_{i=1}^S \sum_{j=0}^{H_j} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right]$$

$$+ \lambda_{\text{smooth}} \sum_{i=0}^{S-2} \sum_{j=0}^{H_j} \left[(\sqrt{w_i} - \sqrt{w_{i+1}})^2 + (\sqrt{h_i} - \sqrt{h_{i+1}})^2 \right]$$

$$+ \sum_{i=0}^{S-1} \sum_{j=0}^{H_j} (C_i - \hat{C}_i)^2$$

$$+ \lambda_{\text{conf}} \sum_{i=0}^{S-1} \sum_{j=0}^{H_j} (C_i - \hat{C}_i)^2$$

$$+ \sum_{i=0}^{S-1} \sum_{c \in \{S, \text{background}\}} (p_i(c) - \hat{p}_i(c))^2$$

$$\mathcal{L} = \mathcal{L}_S + \lambda \mathcal{L}_C$$

$$= - \sum_{i=1}^n \log \frac{e^{W_i^T x_i + b_i}}{\sum_{j=1}^n e^{W_j^T x_i + b_j}} + \frac{\lambda}{2} \sum_{i=1}^n \|x_i - c_{y_i}\|^2$$

Conclusion and Future Work

By simulating YOLO output by using random croppings on LFW images, we were able to reduce base image size by 54%. Further training on FaceNet allowed for comparable error rates.

In its original state, FaceNet relies on MTCNN, an image-pyramid based bounding box locator; far slower than YOLO. Cutting out the alignment step in favor of YOLO should allow faster pre-processing, without loss of accuracy in classification.

Before the final write-up, we will run two more additional tests with a modified FaceNet, to see how speed improves with full-sized inputs of size 250 and 166, respectively. Eliminating image size normalization will give us a definitive idea of how much YOLO croppings can speed up FaceNet.

Given an additional six months, we would optimize the entire pipeline by merging the two models into a single framework, speeding it up by eliminating excess file reading.

References: ^[1] YOLO <https://github.com/allanzelener/YAD2K>, ^[2] FaceNet <https://github.com/davidsandberg/faceNet>, ^[3] WIDER Face <http://mmlab.ie.cuhk.edu.hk/projects/WIDERFace/>, ^[4] LFW <http://vis-www.cs.umass.edu/lfw/>