



SMILE: A Semi-supervised Multiple Instance Learning Framework for Object Detection

Authors: Corey Snyder (ECE PhD Candidate), Molly Dasso (MS '21), Tian Ma (Sandia, Org), Prof. Minh N. Do (U of IL, Dept of Electrical and Computer

Background and Motivation

Deep learning has achieved state of the art performance on the task of object detection. However,

- Training a deep learning model requires large annotated datasets on the order of thousands to millions of labeled examples
- Collecting these annotations can be very time consuming and require expert knowledge while unlabeled data is often abundant

Semi-supervised learning (SSL) operates with a small labeled dataset and a much larger unlabeled dataset, e.g. 100:1 unlabeled:labeled



Figure 1: Increasing levels of supervision strength from left to right

Our Approach

A common framework for SSL is known as **pseudo-labeling** where confident model predictions are used as “pseudo” ground-truth for unlabeled data.

Pseudo-labels may be generated by a pre-trained “teacher” model (fixed) or the object detector as we train it (dynamic).

We compare the concepts of “weak” vs. “strong” semi-supervision. For object detection, we define:

- “strong” semi-supervision as full bounding boxes for “pseudo”-objects in one image
- “weak” semi-supervision as image-level class labels, i.e. which classes are present in one image

Main Hypothesis: there is a tradeoff between the reliability and quantity of information as semi-supervision becomes stronger.

Our SMILE algorithm trains an object detector f by computing a supervised loss over the labeled set and a novel MIL loss over the unlabeled set. For each batch of unlabeled images:

1. Construct MIL targets representing the classes present in the image
2. Apply data augmentation to the unlabeled images
3. Pass the augmented images through the network and compute unsupervised loss between the class probabilities and the MIL targets.

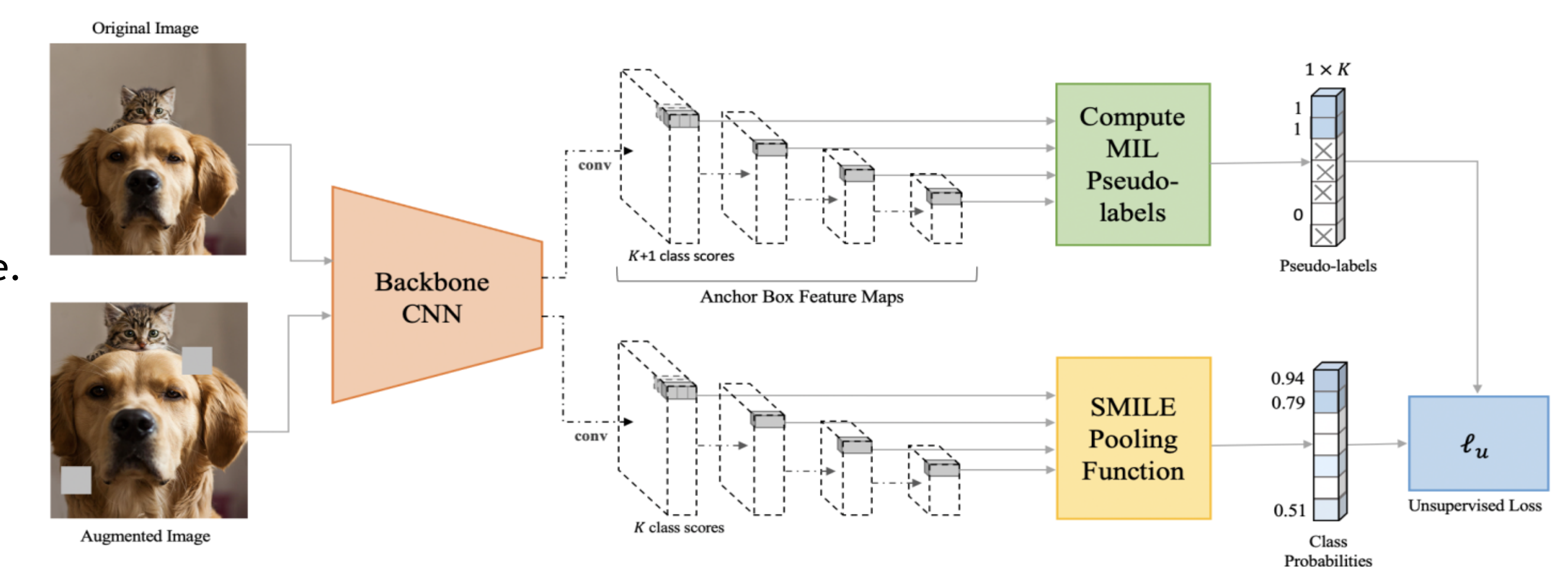


Figure 2: SMILE pipeline

ODD-MNIST

We synthesize our own dataset, **Object Detection with Distractors on MNIST (ODD-MNIST)** to act as an approachable benchmark for object detection algorithms.

To construct an image, we:

- randomly sample an image from the CIFAR-10 dataset to use as the background
 - sample digits from the MNIST dataset and paste them as **target** objects
 - sample letters from the E-MNIST dataset and paste them as **distractor** objects
- The targets are the bounding box coordinates and class labels for the MNIST digits.

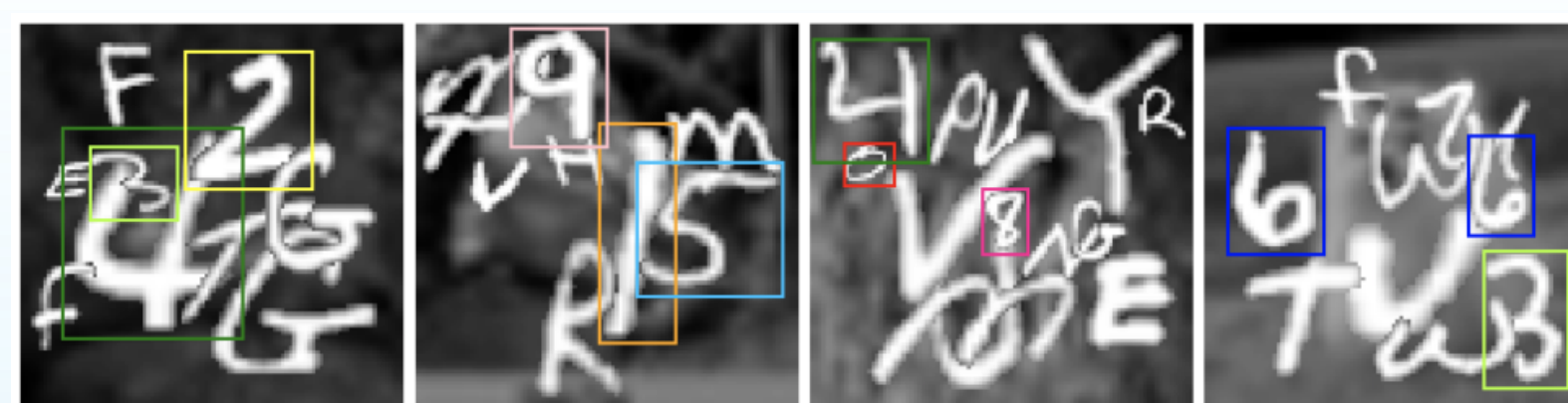


Figure 3: Example images from ODD-MNIST

Key Results

Table 1: Test mean absolute precision (mAP) results for supervised baseline, STAC (strong semi-supervision baseline), and our SMILE method on ODD-MNIST. Each column gives the percent of training data that is labeled.

	0.1%	0.2%	0.4%	1%	2%
Supervised	13.82 ± 1.70	23.20 ± 1.45	34.53 ± 0.52	48.50 ± 0.03	56.46 ± 0.16
STAC [1]	19.38 ± 2.44	31.78 ± 1.64	42.98 ± 0.39	55.02 ± 0.18	60.59 ± 0.06
SMILE	15.68 ± 4.04	33.91 ± 2.49	44.47 ± 0.60	53.08 ± 0.23	57.67 ± 0.07

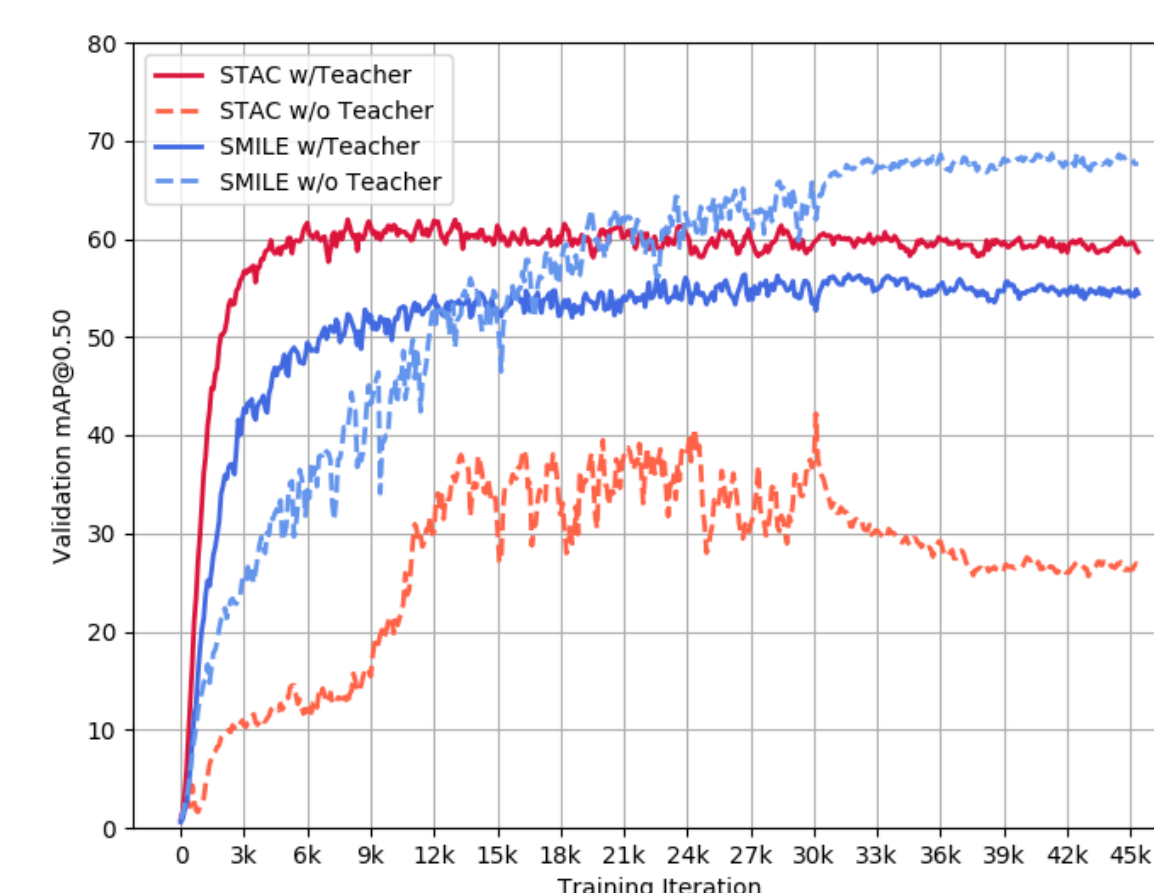


Figure 4: Validation mAP during training for SMILE and STAC with and without a teacher model.

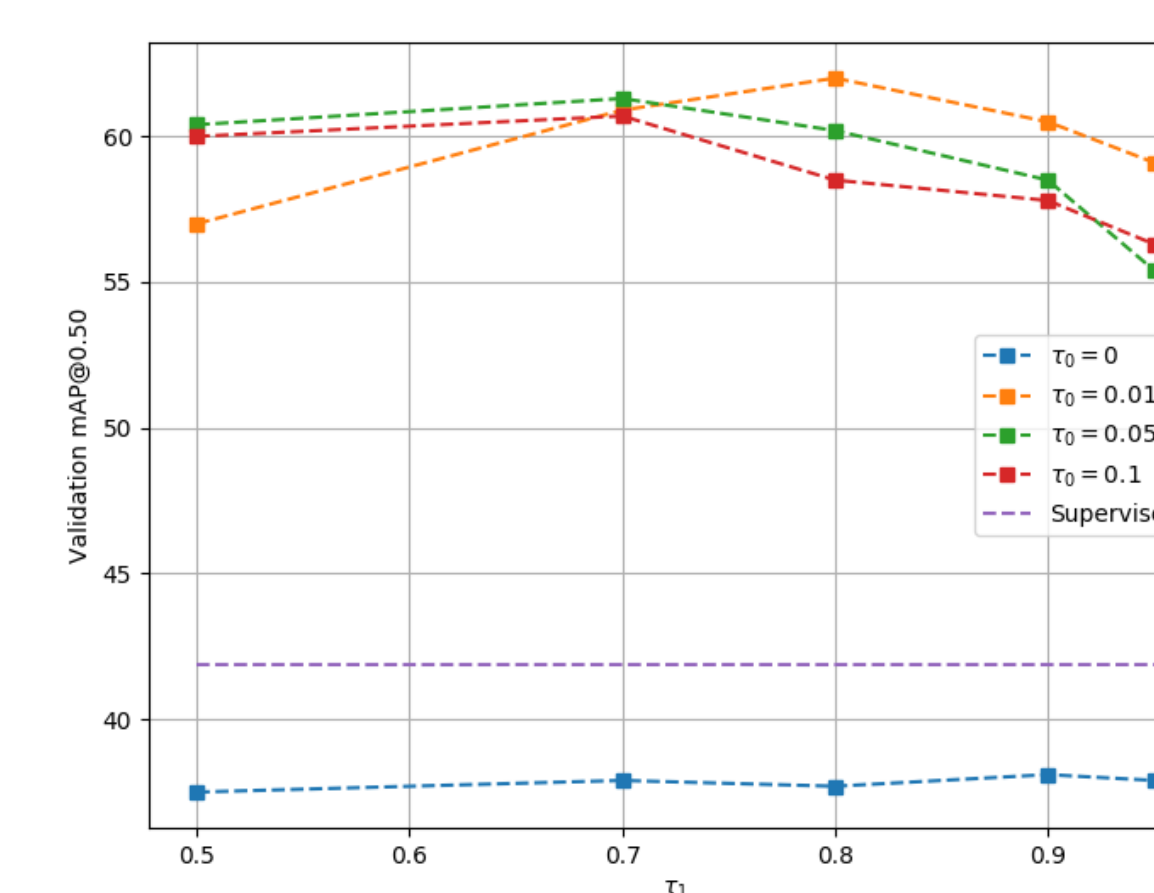


Figure 5: Validation mAP of SMILE against different choices of pseudo-labeling thresholds τ_0 and τ_1 .

Next Steps

To strengthen a conference submission, we plan to:

- Conduct experiments with a popular two-stage object detector, e.g. Faster R-CNN
- Benchmark SMILE on popular natural image datasets like PASCAL and COCO
- Compile further ablation studies regarding quality and quantity of weak vs. strong semi-supervision

[1] Kihyuk Sohn, Zizhao Zhang, Chun-Liang Li, Han Zhang, Chen-Yu Lee, and Tomas Pfister. A simple329semi-supervised learning framework for object detection. InarXiv:2005.04757, 2020