# 2nd workshop on Unsupervised Learning for Automated Driving
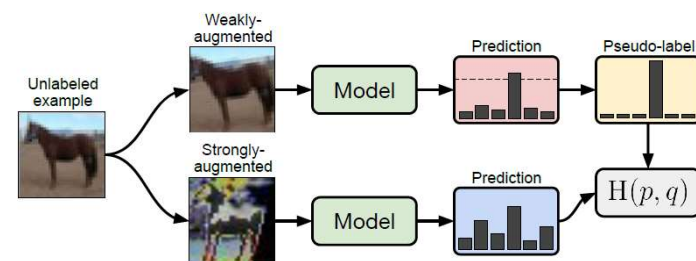
## Leveraging single and cross-modal unlabeled data for learning with limited labels

**Zsolt Kira**
**Assistant Professor**
**School of Interactive Computing**
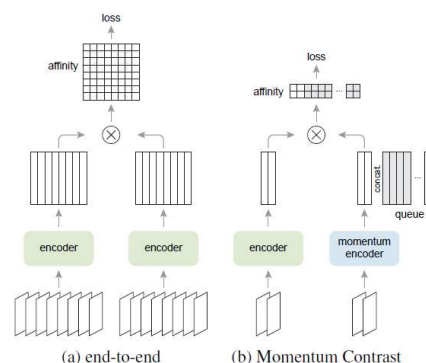**Georgia Tech**
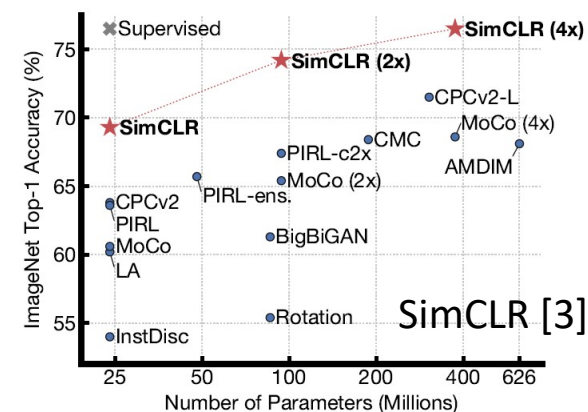
# The Fast-Paced Landscape of Reducing Labels

- The past two years have seen tremendous progress in:
  - Zero-shot learning
  - Few-shot learning
  - Semi-supervise learning
  - Self-supervised learning
  - Domain adaptation/generalization
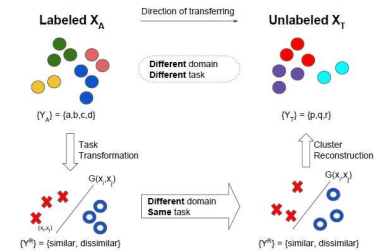  - Weakly supervised learning
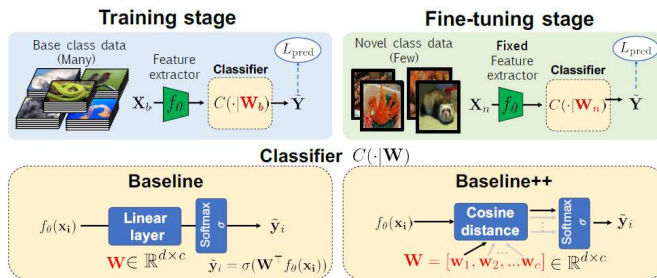  - Long-tailed datasets

FixMatch [1]

MoCo-v2 [2]

SimCLR [3]

[1] FixMatch: Simplifying Semi-Supervised Learning with Consistency and Confidence, Sohn et al.
[2] Improved Baselines with Momentum Contrastive Learning, Chen et al.
[3] A Simple Framework for Contrastive Learning of Visual Representations, Chen et al.
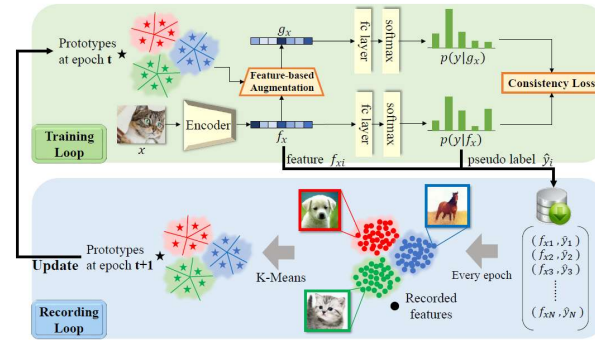
# Our Contributions



**Pairwise Similarity for Cross-Task Object Discovery**
[ICLR 2018, 2019]

**Complex Data Augmentation Domain Generalization/SSL**
[ECCV 2020]



**Closer Look @ Few-Shot (w/ VT)**
[ICLR 2019]

**Video Domain adaptation**
[CVPR 2019, ICCV 2019]

**Student-Teacher for Semi-Supervised Object Detection**
[in submission, with FB]

Teacher

Weakly augmented

Original

strongly augmented

RPN

ROIHead

Prediction

Pseudo-labeling

Shared Weights

Supervision

RPN

ROIHead

Student

**2D to 3D Inflation for semi-supervised learning**
[https://arxiv.org/abs/2008.10592, with Argo AI]

# The Methods are Surprisingly Simple

- A handful of common techniques:
  - **Data augmentation**
  - **Pseudo-labeling** / distillation
  - Surrogate tasks / contrastive losses
  - Temperature scaling / Entropy maximization
  - Cosine/metric learning
  - **Prototypes**
  - **Graph neural networks**
  - Meta-learning

# Setting: Semi-Supervised Learning
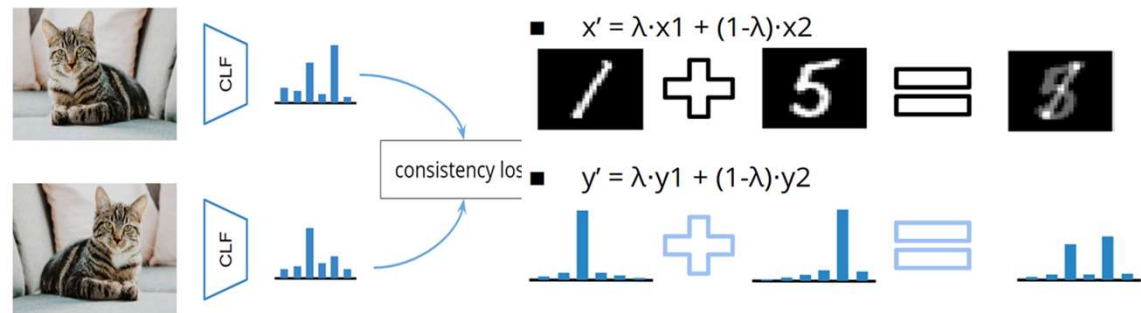
- Setting
  - Small amount of labeled data
  - Large amount of unlabeled data

- Example Datasets
  - SVHN
  - CIFAR-10
  - CIFAR-100
  - mini-ImageNet

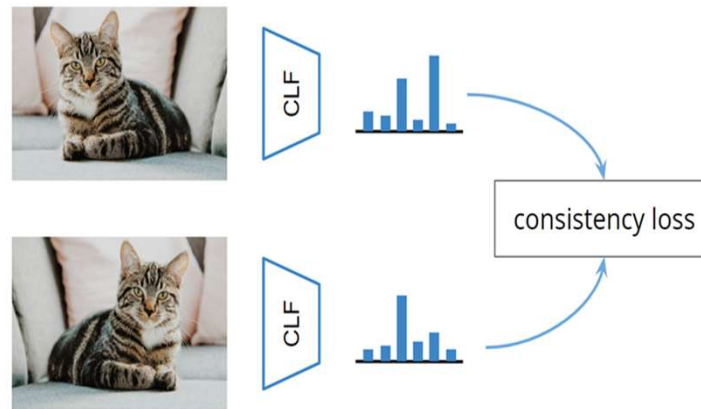- Previous SoA method: MixMatch [1], with key contributions:
  - Consistency
  - Mixup
  - Mean teacher (for more reliable pseudo-labeling)



$$x' = \lambda \cdot x1 + (1-\lambda) \cdot x2$$

consistency los

$$y' = \lambda \cdot y1 + (1-\lambda) \cdot y2$$

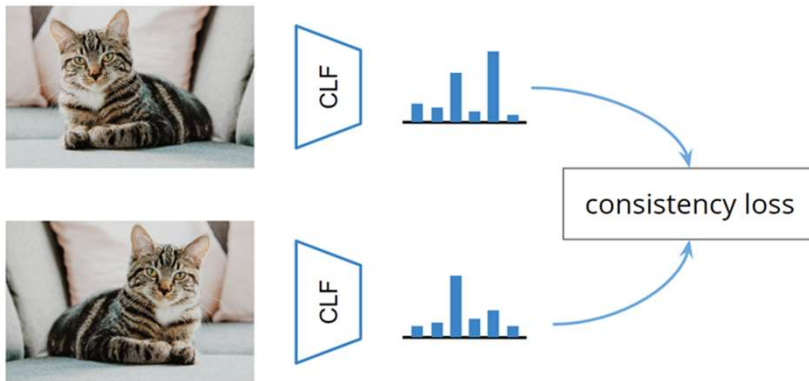[1] MixMatch: A Holistic Approach to Semi-Supervised Learning

# Data Augmentation

- Data augmentation key to many different areas, including:
  - Semi-supervised learning
  - Self-supervised learning
  - Showing up in few-shot learning, etc.

# Limitations of Consistency-Based Method

- Data augmentation only operates in **image space**
  - limits the possible transformations to textural or geometric within images.

- Data augmentation operates within a single instance
  - fails to transform data with the knowledge of other instances **(manifold structure)**

- Mixup method (sort of) addressed these issues.
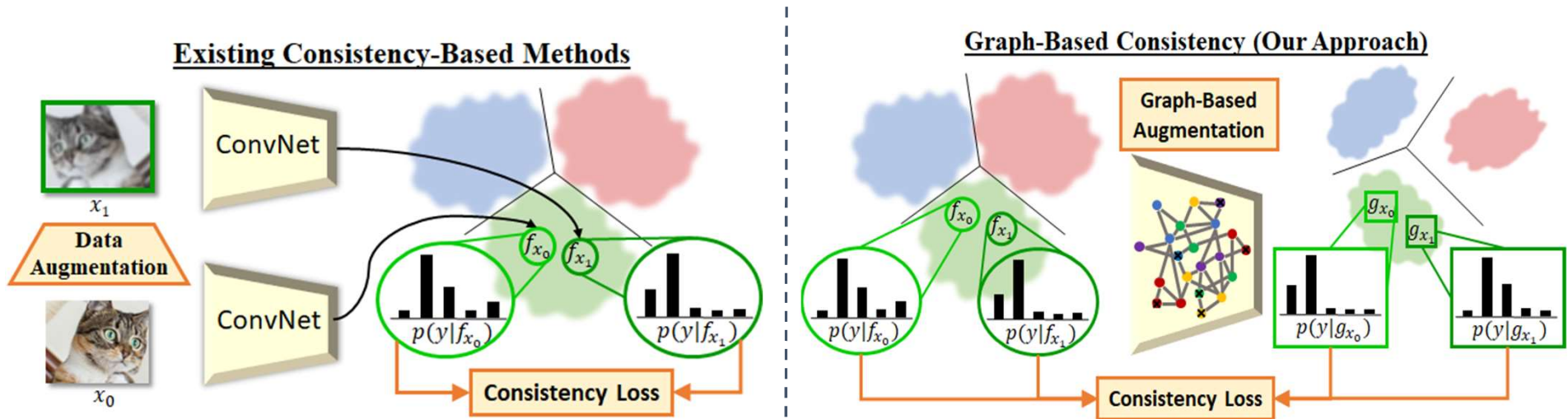  - However, the transformations are still in image space
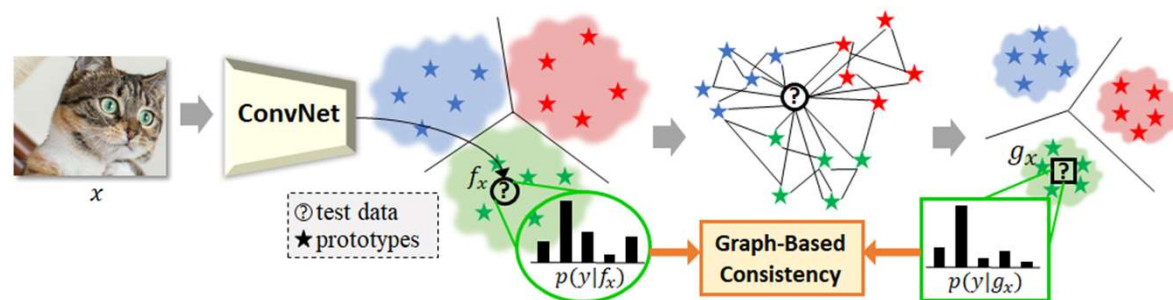
# FeatMatch: Proposed Graph-Based Consistency Method

- Construct a graph neural network in the feature space for feature augmentation.

- Method is orthogonal to consistency-based methods (here combined with MixMatch, can also use FixMatch)



Kuo et al., FeatMatch: Feature-Based Augmentation for Semi-Supervised Learning, ECCV 2020

# Graph Module

- The graph is constructed between image features and a set of prototypes
  - Computationally intractable to construct graph over the entire dataset

- Prototype extraction
  - K-mean centroids for each class every K iterations
  - Small set of labeled data: extract features on the fly
  - Large set of unlabeled data: record previous features ([2] has similar idea)



[2] Momentum Contrast for Unsupervised Visual Representation Learning (https://arxiv.org/abs/1911.05722)

# Updating Prototypes

# Comparison to Other Methods

| | ReMixMatch [3] | MixMatch [4] | Mean Teacher [26] | ICT [30] | PLCB [1] | FeatMatch (Ours) |
|---|---|---|---|---|---|---|
| **Feature-Based Augmentation** | - | - | - | - | - | ✔ |
| Image-Based Augmentation | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ |
| Temporal Ensembling | ✔ | ✔ | ✔ | - | - | - |
| Self-Supervised Loss | ✔ | - | - | - | - | - |
| Alignment of Class Distribution | ✔ | - | - | - | ✔ | - |

# Quantitative Results (Standard Datasets)

| Method | Model (param.) | CIFAR-10 # Labeled samples | | | SVHN # Labeled samples | | |
|---|---|---|---|---|---|---|---|
| | | 250 | 1,000 | 4,000 | 250 | 1,000 | 4,000 |
| SSL with Memory [5] | | - | - | $11.91 \pm 0.22$ | 8.83 | 4.21 | - |
| Deep Co-Training [22] | | - | - | $8.35 \pm 0.06$ | - | $3.29 \pm 0.03$ | - |
| Weight Averaging [2] | | - | $15.58 \pm 0.12$ | $9.05 \pm 0.21$ | - | - | - |
| ICT [30] | CNN-13 (3M) | - | $15.48 \pm 0.78$ | $7.29 \pm 0.02$ | $4.78 \pm 0.68$ | $3.89 \pm 0.04$ | - |
| Label Propagation [14] | | - | $16.93 \pm 0.70$ | $10.61 \pm 0.28$ | - | - | - |
| SNTG [18] | | - | $18.41 \pm 0.52$ | $9.89 \pm 0.34$ | $4.29 \pm 0.23$ | $3.86 \pm 0.27$ | - |
| PLCB [1] | | - | $6.85 \pm 0.15$ | $5.97 \pm 0.15$ | - | - | - |
| $\Pi$-model [25] | | $53.02 \pm 2.05$ | $31.53 \pm 0.98$ | $17.41 \pm 0.37$ | $17.65 \pm 0.27$ | $8.60 \pm 0.18$ | $5.57 \pm 0.14$ |
| PseudoLabel [17] | | $49.98 \pm 1.17$ | $30.91 \pm 1.73$ | $16.21 \pm 0.11$ | $21.16 \pm 0.88$ | $10.19 \pm 0.41$ | $5.71 \pm 0.07$ |
| Mixup [13] | | $47.43 \pm 0.92$ | $25.72 \pm 0.66$ | $13.15 \pm 0.20$ | $39.97 \pm 1.89$ | $16.79 \pm 0.63$ | $7.96 \pm 0.14$ |
| VAT [19] | WRN (1.5M) | $36.03 \pm 2.82$ | $18.68 \pm 0.40$ | $11.05 \pm 0.31$ | $8.41 \pm 1.01$ | $5.98 \pm 0.21$ | $4.20 \pm 0.15$ |
| Mean Teacher [26] | | $47.32 \pm 4.71$ | $17.32 \pm 4.00$ | $10.36 \pm 0.25$ | $6.45 \pm 2.43$ | $3.75 \pm 0.10$ | $3.39 \pm 0.11$ |
| MixMatch [4] | | $11.08 \pm 0.87$ | $7.75 \pm 0.32$ | $6.24 \pm 0.06$ | $3.78 \pm 0.26$ | $3.27 \pm 0.31$ | $2.89 \pm 0.06$ |
| ReMixMatch [3] | | $\mathbf{6.27 \pm 0.34}$ | $\mathbf{5.73 \pm 0.16}$ | $5.14 \pm 0.04$ | $\mathbf{3.10 \pm 0.50}$ | $\mathbf{2.83 \pm 0.30}$ | $\mathbf{2.42 \pm 0.09}$ |
| FeatMatch (Ours) | | $7.50 \pm 0.64$ | $\mathbf{5.76 \pm 0.07}$ | $4.91 \pm 0.18$ | $\mathbf{3.34 \pm 0.19}$ | $\mathbf{3.10 \pm 0.06}$ | $2.62 \pm 0.08$ |

# Quantitative Results (Larger Datasets)

- Our method is scalable to larger datasets and categories

| | CIFAR-100 | | mini-ImageNet | |
| | # Labeled samples | | # Labeled samples | |
| Method | 4,000 | 10,000 | 4,000 | 10,000 |
|---|---|---|---|---|
| $\Pi$-model [25] | - | $39.19 \pm 0.36$ | - | - |
| SNTG [18] | - | $37.97 \pm 0.29$ | - | - |
| SSL with Memory [5] | - | $34.51 \pm 0.61$ | - | - |
| Deep Co-Training [22] | - | $34.63 \pm 0.14$ | - | - |
| Weight Averaging [2] | - | $33.62 \pm 0.54$ | - | - |
| Mean Teacher [26] | $45.36 \pm 0.49$ | $36.08 \pm 0.51$ | $72.51 \pm 0.22$ | $57.55 \pm 1.11$ |
| Label Propagation [14] | $43.73 \pm 0.20$ | $35.92 \pm 0.47$ | $70.29 \pm 0.81$ | $57.58 \pm 1.47$ |
| PLCB [1] | $37.55 \pm 1.09$ | $32.15 \pm 0.50$ | $56.49 \pm 0.51$ | $46.08 \pm 0.11$ |
| FeatMatch (Ours) | $\mathbf{31.06 \pm 0.41}$ | $\mathbf{26.83 \pm 0.04}$ | $\mathbf{39.05 \pm 0.06}$ | $\mathbf{34.79 \pm 0.22}$ |

# Quantitative Results (DomainNet)

- Our method is more robust to out-of-distribution unlabeled data
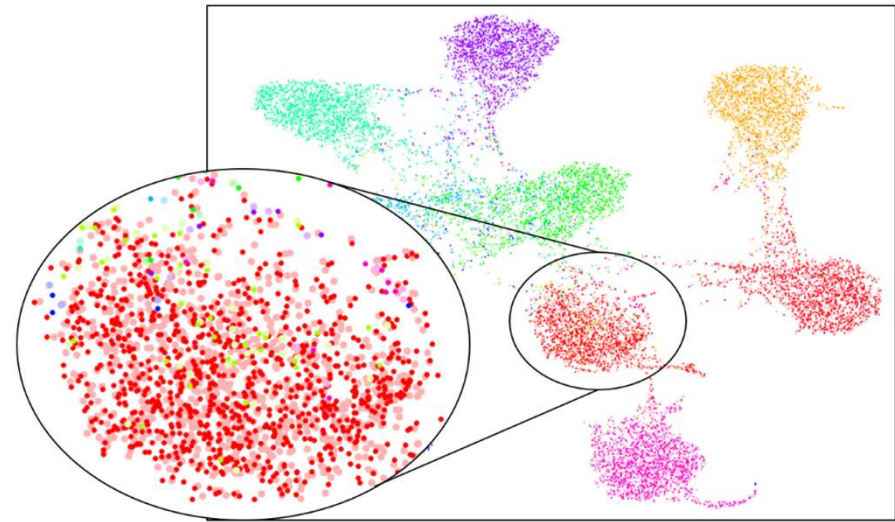  - Half of data from shifted domain for r = 50%

| Method (5% labeled samples) | $r_u = 0\%$ | $r_u = 25\%$ | $r_u = 50\%$ | $r_u = 75\%$ |
|---|---|---|---|---|
| (Semi-supervised) Baseline | $56.63 \pm 0.17$ | $62.44 \pm 0.67\ \%$ | $65.82 \pm 0.07$ | $70.50 \pm 0.51$ |
| FeatMatch (Ours) | $\mathbf{40.66 \pm 0.60}$ | $\mathbf{46.11 \pm 1.15}$ | $\mathbf{54.01 \pm 0.66}$ | $\mathbf{58.30 \pm 0.93}$ |
| Supervised baseline (5% labeled samples, lower bound) | | $77.25 \pm 0.52$ | | |
| Supervised baseline (100% labeled samples, upper bound) | | $31.91 \pm 0.15$ | | |

# Qualitative Results



Data Augmentation

Graph-Based
Feature Augmentation

15

# Prototypes and Similarity



(a) t-SNE of selected prototypes.

(b) Leaned attention weights.

(c) Nearest image neighbors of prototypes

- We can see variability in prototypes, and similarity function largely focuses on same-class prototypes

# Augmentation Visualization



- Graph-based augmentation produces more variable, uniformly distributed augmentations
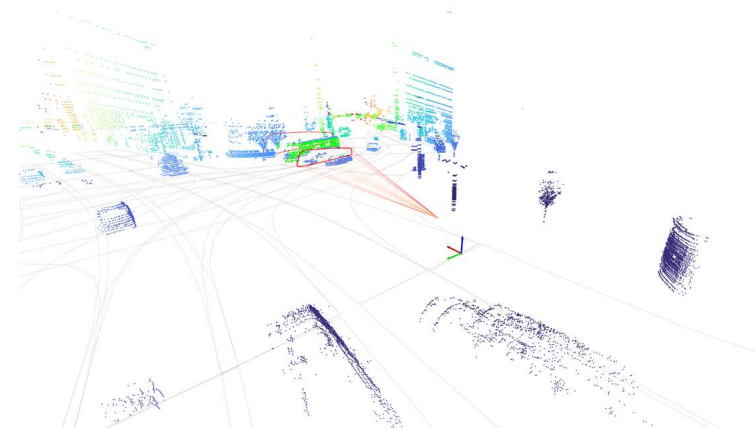
# The Methods are Surprisingly Simple

- A handful of common techniques:
  - **Data augmentation**
  - **Pseudo-labeling** / distillation
  - Surrogate tasks / contrastive losses
  - Temperature scaling / Entropy maximization
  - Cosine/metric learning
  - **Prototypes**
  - **Graph neural networks**
  - Meta-learning

- Autonomous vehicles:
  - **Object Detection?**
  - **Multi-modal?**

What about Object Detection?



What about 3D?

# Motivation: Object Detection

- Few papers address 2-stage detector under semi-supervised setting.
  - ISD [1] and CSD [2] mainly focus on 1-stage detector
- 2-stage detectors generally have more accurate predictions.
  - Previous works focus on ROIheads training, and RPN net training is seldom explored.

Image Classification
(Single Object)



"Cat"

Object Detection
(Multiple objects + Bounding boxes)



"Bike": (120, 201, 356, 347)
"Car": (345, 318, 945, 847)
"Truck": (420, 512, 601, 782)
"Traffic Light": (430, 60, 467, 123)

[1] "Consistency-based Semi-supervised Learning for Object detection", Jeong et al., NeurIPS 2019
[2] "Interpolation-based semi-supervised learning for object detection", Jeong et al., arXiv June 2020

19

1. Two-stage detector
   - 3 modules – Feature Backbone, Region Proposal Network, and ROIHead

# 1. Two-stage detector

- 4 predictions/losses
  - RPN - Foreground/Background Detection, Bounding Box Regression
  - ROI head – Patch Classification, Bounding Box Regression

# Observation 1: A good RPN is necessary!

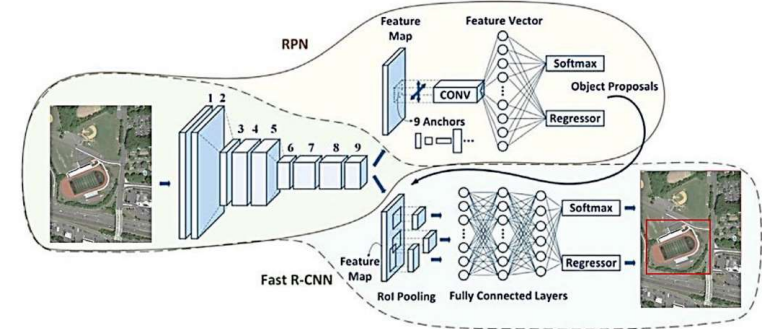Goal: Verify a good RPN is important; A good ROIhead requires a good RPN

Train a model with
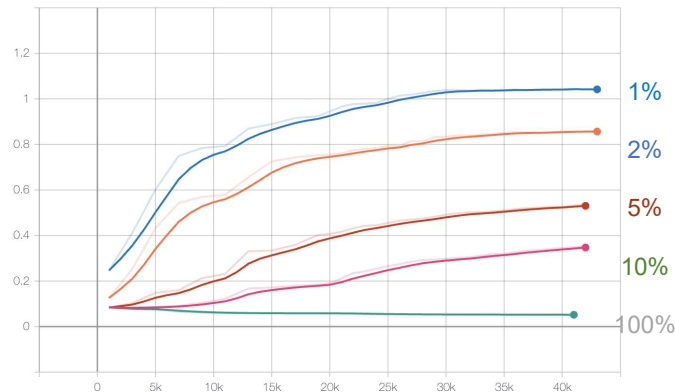- RPN: 1% supervised data
- ROIHead: 100% supervised data

# Observation 2: Overfitting!

2. When the labels are insufficient …

- Overfitting!
  - Foreground/Background Classification



**Validation** Fg-Bg Classification Loss
(RPN)



1%
2%
5%
10%
100%

**Validation** Box Regression Loss
(RPN)



1%
2%
5%
10%
100%

23

# Observations

2. When the labels are insufficient …

- Overfitting!
  - Foreground/Background Classification
  - Patch Classification
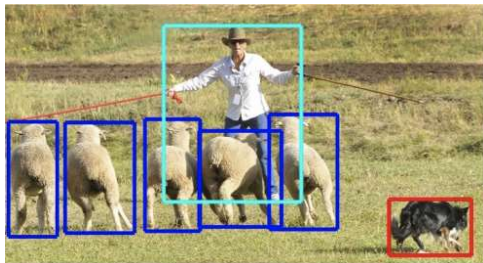


**Validation** Patch Classification Loss
(ROIHead)



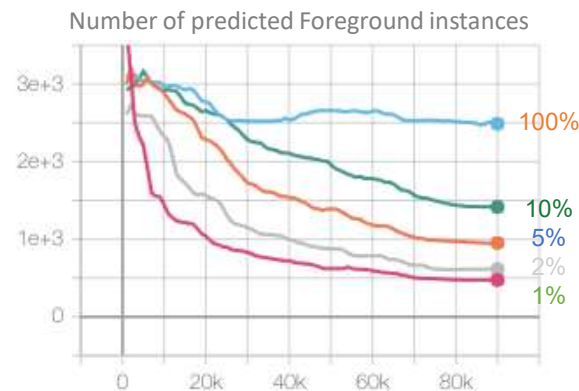**Validation** Box Regression Loss
(ROIHead)

# Why does overfitting occur?

1. Foreground-background imbalance
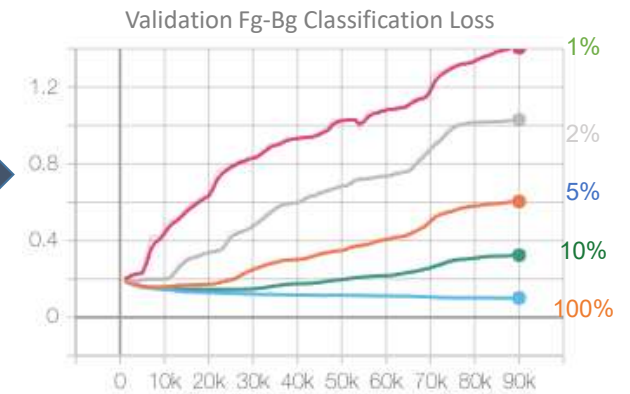2. Foreground Classes imbalance

Foreground : Background = 1 : 3
(Ground-truth data)
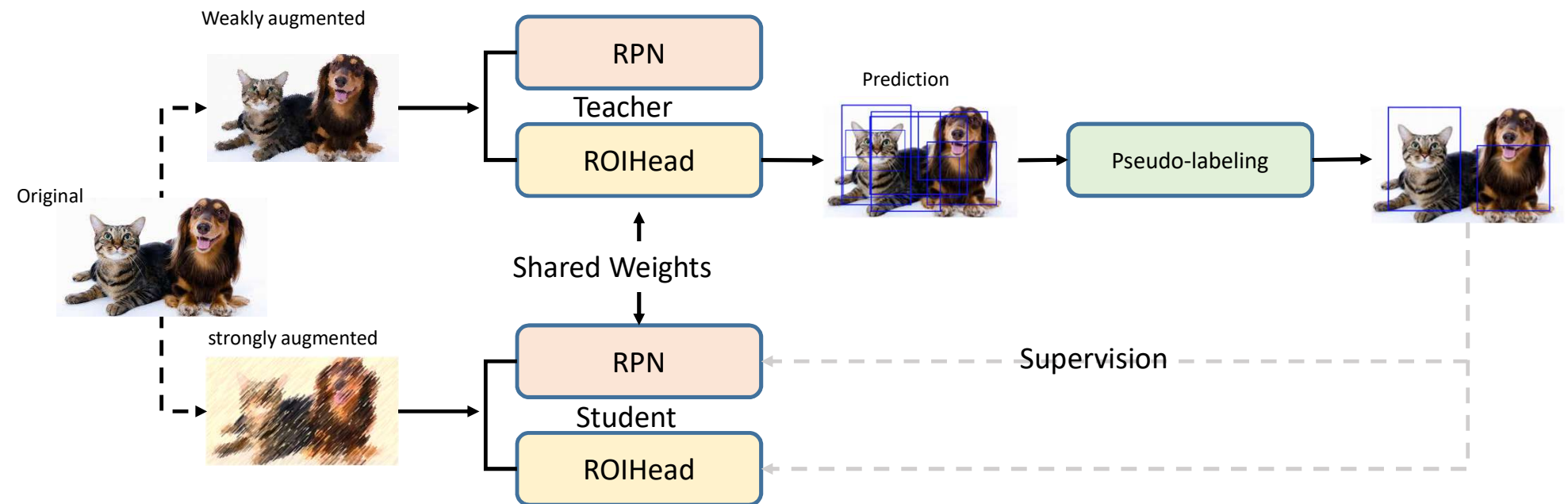
Model biases to Background

Overfitting for fg-bg prediction



Number of predicted Foreground instances

Validation Fg-Bg Classification Loss

# Semi-supervised Object Detection
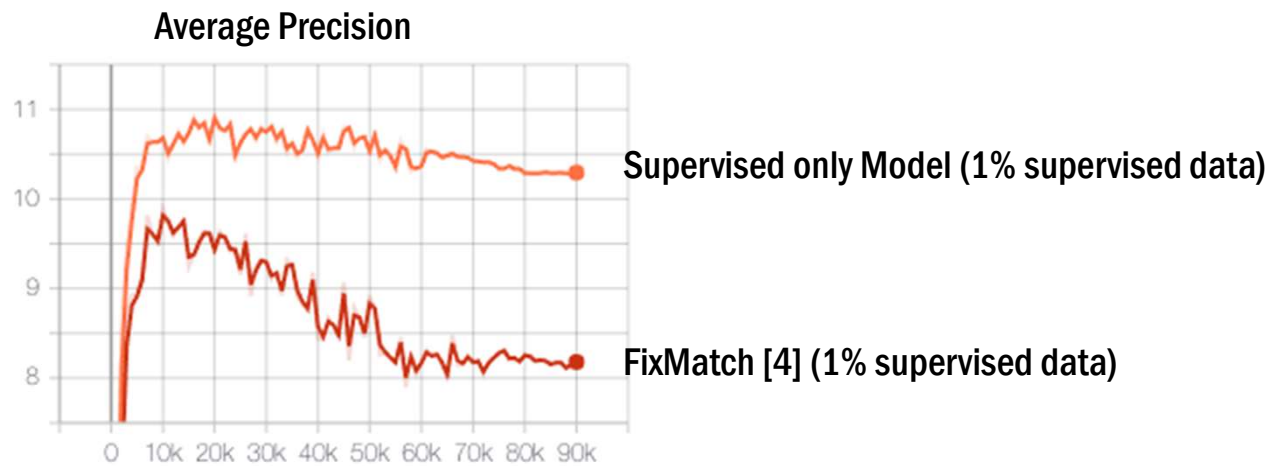
## Problem 1: Labeled data is insufficient

Trial: Using the state-of-the-art semi-supervised classification method [4]?

# Semi-supervised Object Detection

## Problem 1: Labeled data is insufficient

Trial: Using the state-of-the-art semi-supervised classification method [4]?
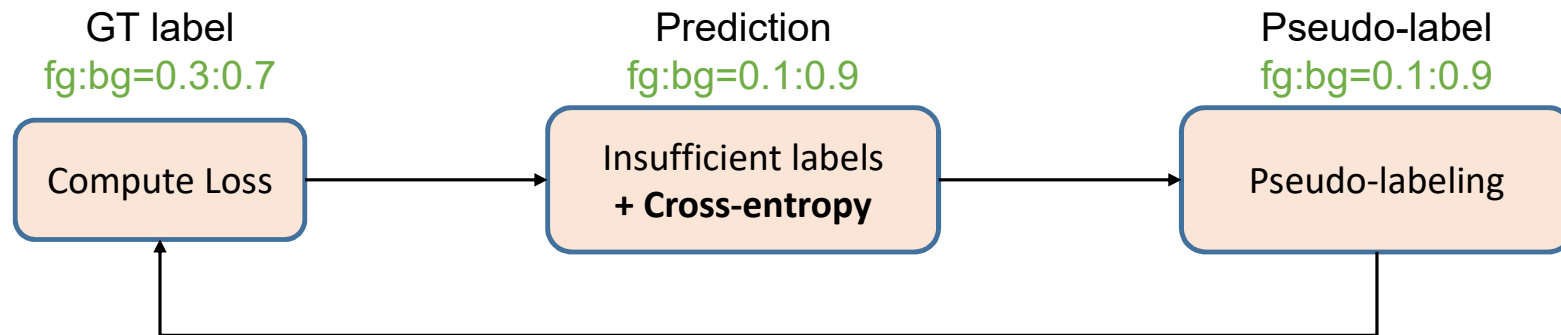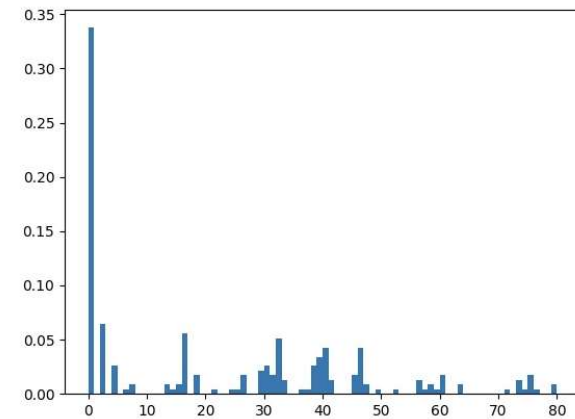
**Average Precision**



Supervised only Model (1% supervised data)

FixMatch [4] (1% supervised data)

[4] "FixMatch: Simplifying Semi-Supervised Learning with Consistency and Confidence", Sohn et al. arXiv, Jan 2020

27

# Semi-supervised Object Detection

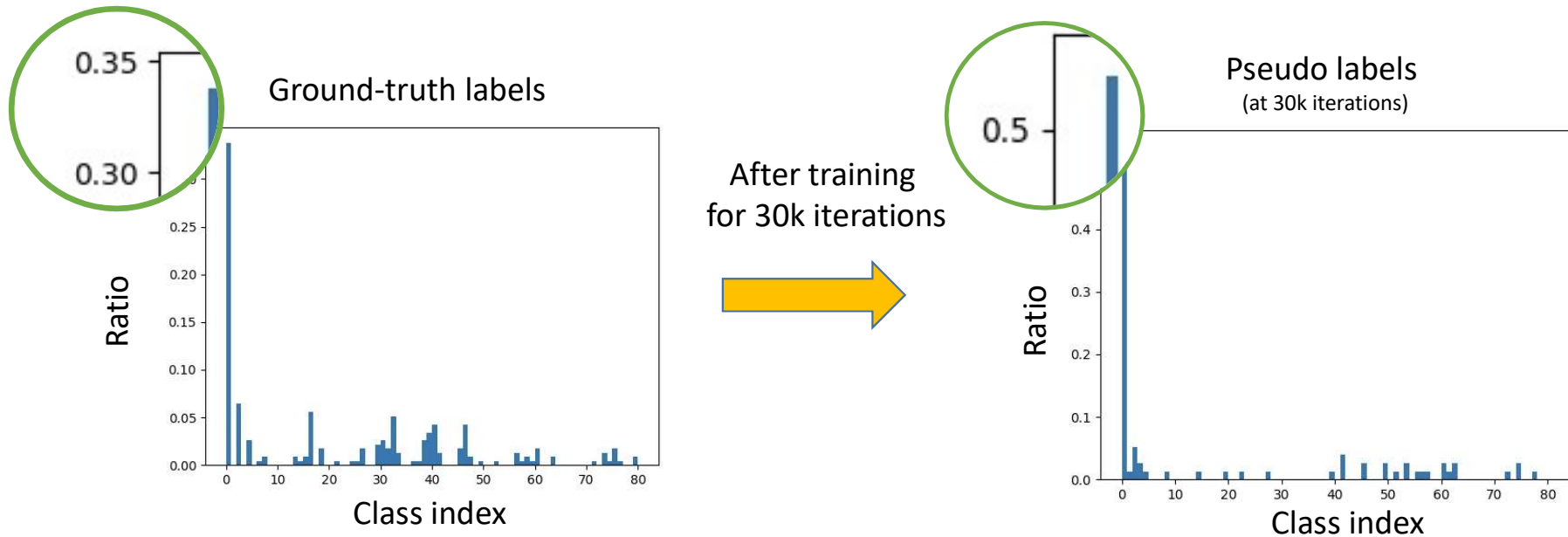Why SOTA semi-supervised classification cannot work?
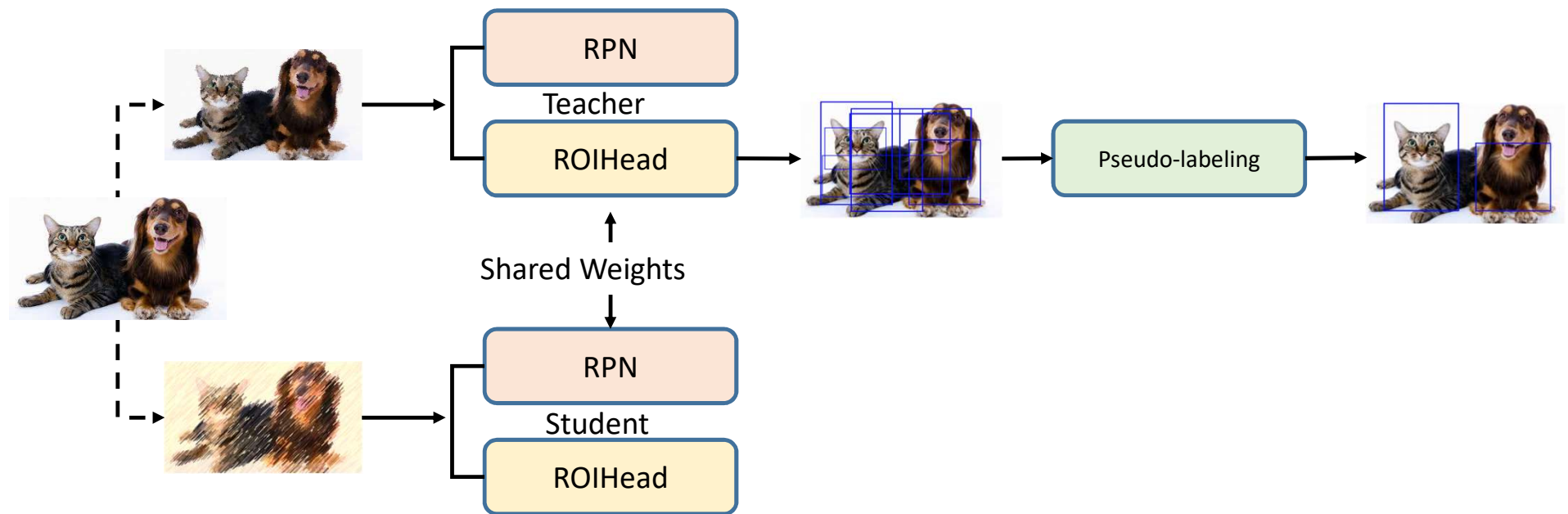Reason: Class imbalance!
Pseudo-labeling is a closed loop



| GT label | Prediction | Pseudo-label |
|----------|------------|--------------|
| fg:bg=0.3:0.7 | fg:bg=0.1:0.9 | fg:bg=0.1:0.9 |
| Compute Loss | Insufficient labels **+ Cross-entropy** | Pseudo-labeling |

# Semi-supervised Object Detection

## Problem 2: Pseudo-label biases

- Cross-entropy → Model biases majority classes



Ground-truth labels

After training
for 30k iterations
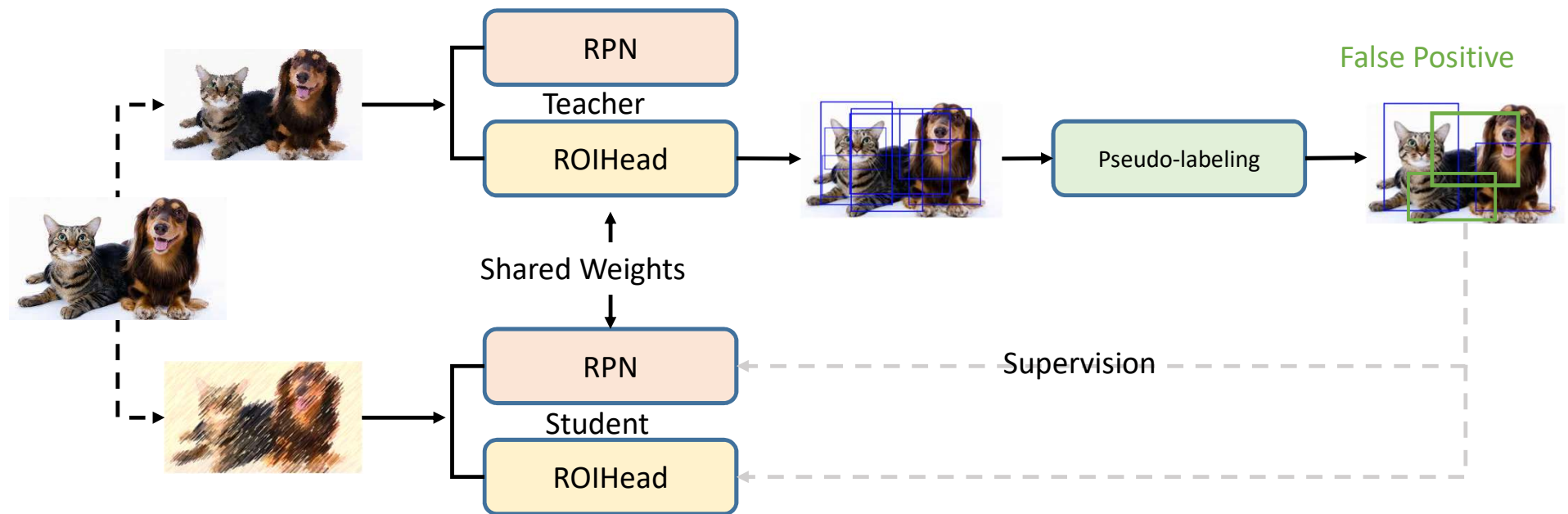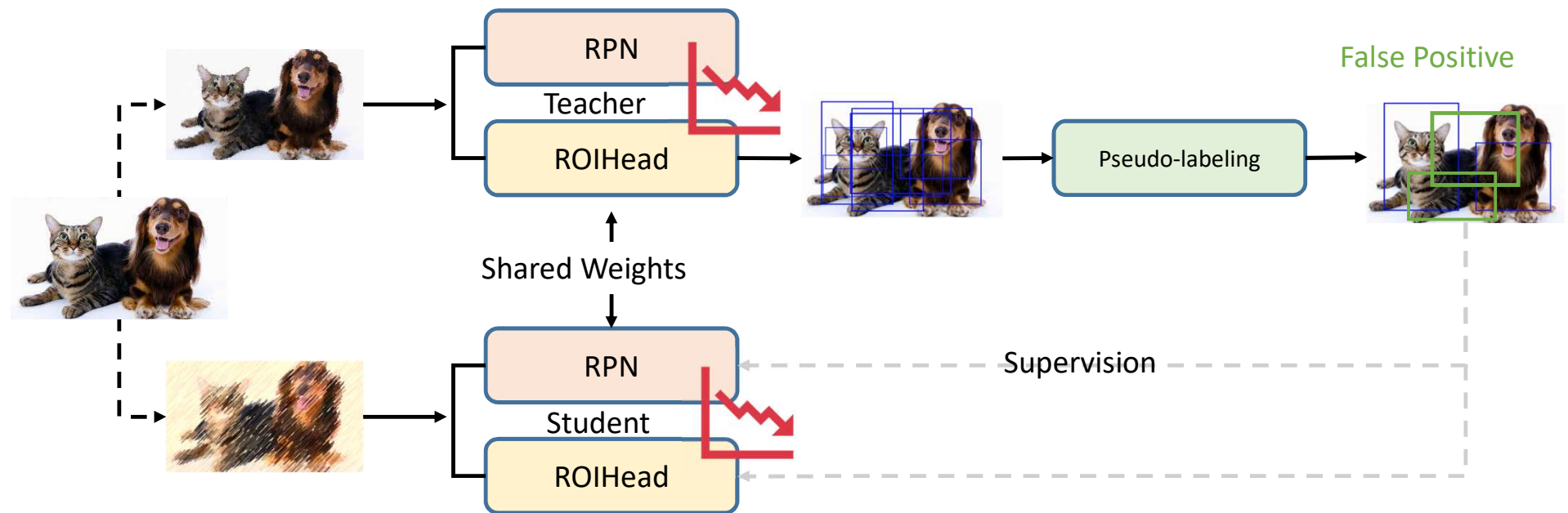
Pseudo labels
(at 30k iterations)

# Semi-supervised Object Detection

Problem 3: Pseudo-label generation is not stable; False positive samples are detrimental

# Semi-supervised Object Detection

Problem 3: Pseudo-label generation is not stable; False positive samples are detrimental
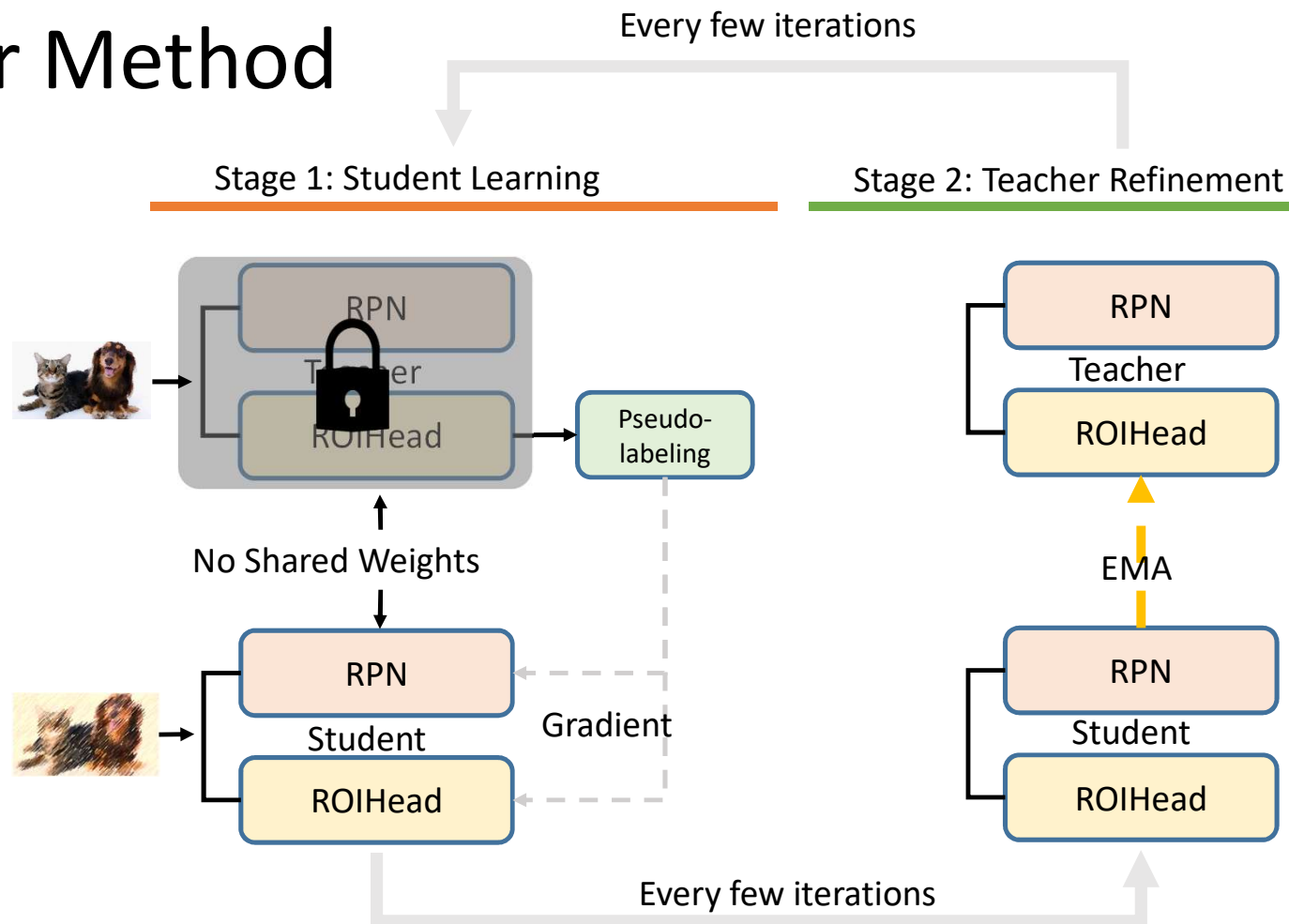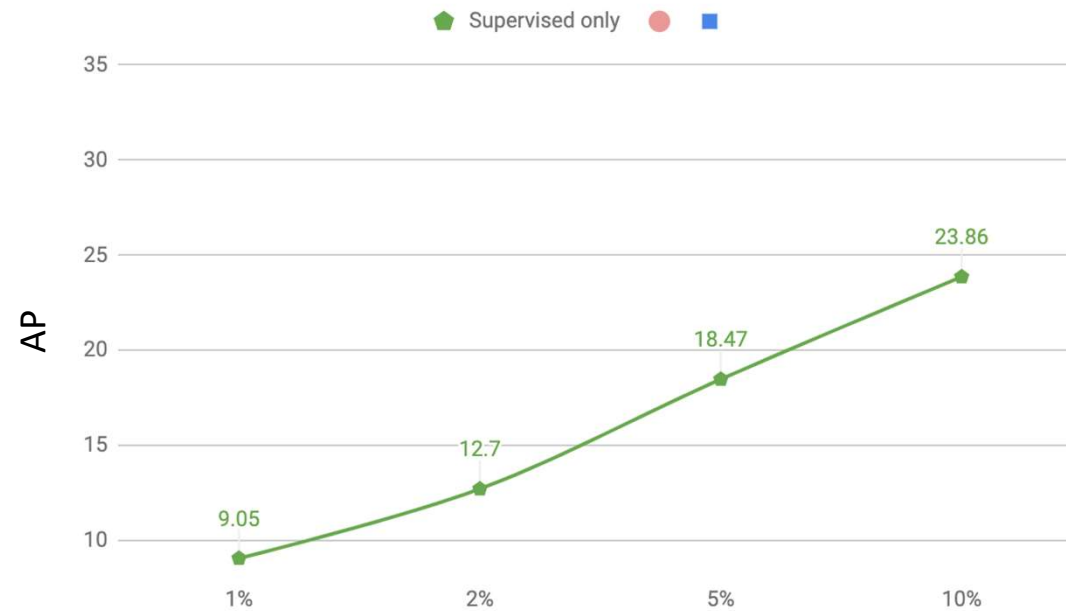
# Semi-supervised Object Detection

Problem 3: Pseudo-label generation is not stable; False positive samples are detrimental
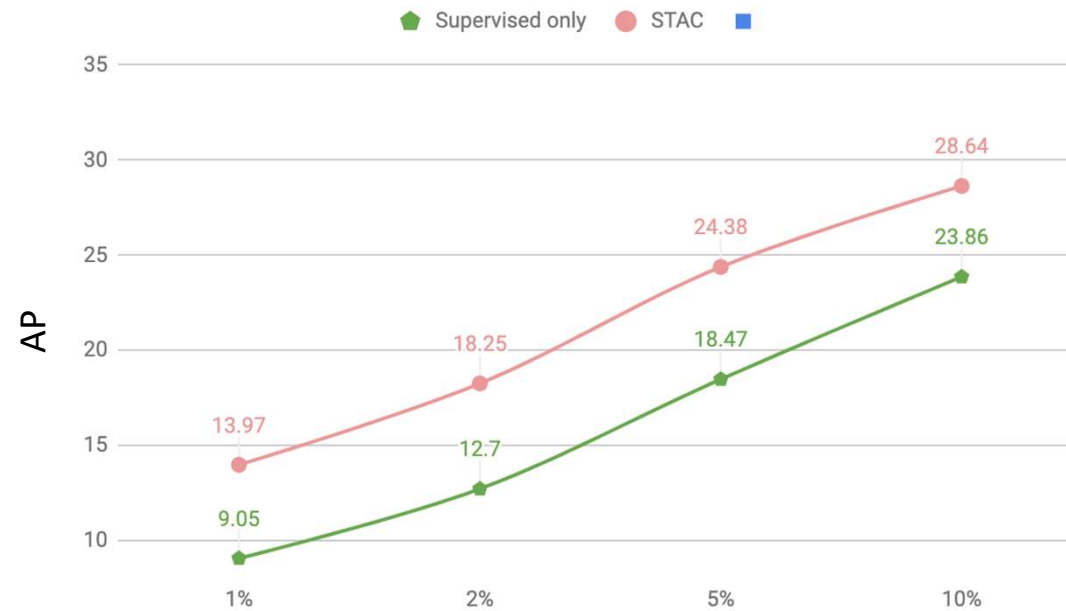
# Our Method

Every few iterations



Stage 1: Student Learning

Stage 2: Teacher Refinement

RPN

Teacher

ROIHead

Pseudo-labeling

No Shared Weights

RPN

Student

ROIHead

Gradient

RPN

Teacher

ROIHead

EMA

RPN

Student

ROIHead

Every few iterations

33

# Experiments



| | 1% | 2% | 5% | 10% |
|---|---|---|---|---|
| Supervised only | 9.05 | 12.70 | 18.47 | 23.86 |

# Experiments



| | 1% | 2% | 5% | 10% |
|---|---|---|---|---|
| Supervised only | 9.05 | 12.70 | 18.47 | 23.86 |
| STAC [2]<br>(SOTA from Google) | 13.97 (+4.92) | 18.25 (+5.55) | 24.38 (+5.91) | 28.64 (+4.78) |

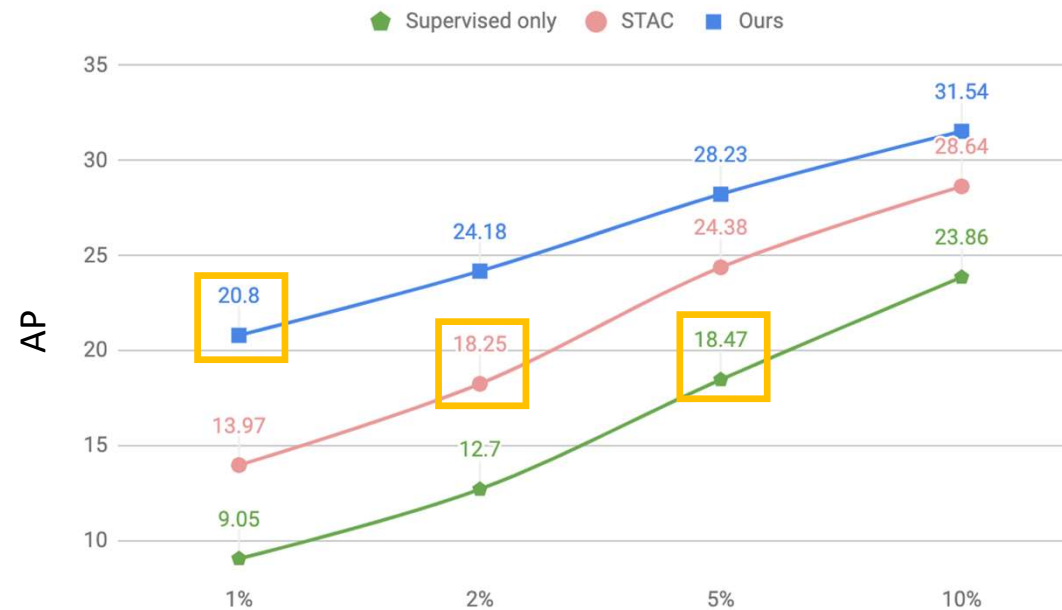[2]"A Simple Semi-Supervised Learning Framework for Object Detection", Sohn et al., arXiv May 2020

# Experiments



| | 1% | 2% | 5% | 10% |
|---|---|---|---|---|
| Supervised only | 9.05 | 12.70 | 18.47 | 23.86 |
| STAC [2]<br>(SOTA from Google) | 13.97 (+4.92) | 18.25 (+5.55) | 24.38 (+5.91) | 28.64 (+4.78) |
| Ours | 20.80 (+11.75) | 24.18 (+11.48) | 28.23 (+9.76) | 31.54 (+7.68) |

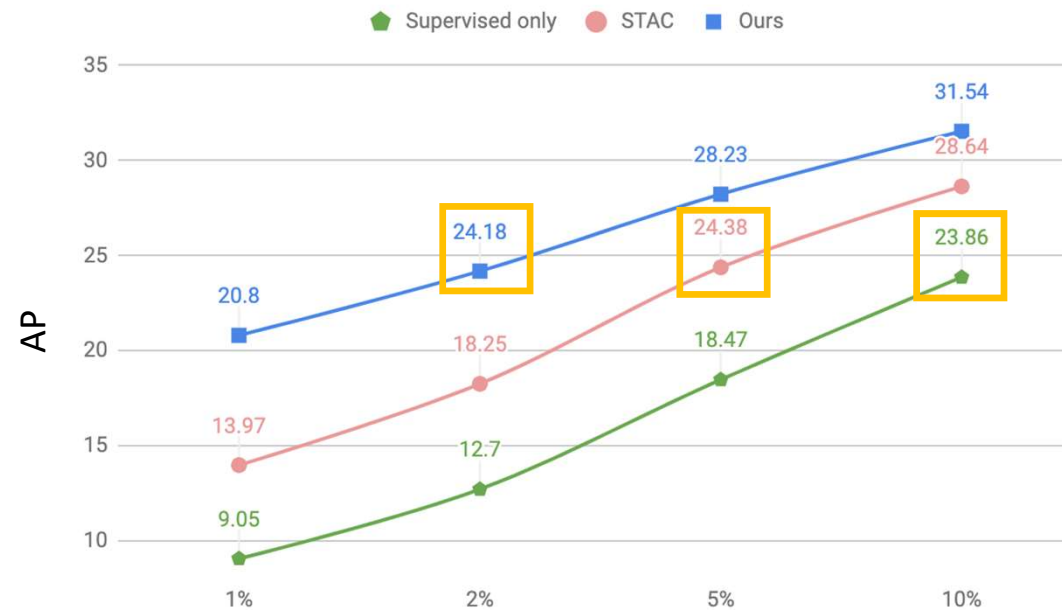[2]"A Simple Semi-Supervised Learning Framework for Object Detection", Sohn et al., arXiv May 2020

# Experiments



| | 1% | 2% | 5% | 10% |
|---|---|---|---|---|
| Supervised only | 9.05 | 12.70 | 18.47 | 23.86 |
| STAC [2] (SOTA from Google) | 13.97 (+4.92) | 18.25 (+5.55) | 24.38 (+5.91) | 28.64 (+4.78) |
| Ours | 20.80 (+11.75) | 24.18 (+11.48) | 28.23 (+9.76) | 31.54 (+7.68) |

[2]"A Simple Semi-Supervised Learning Framework for Object Detection", Sohn et al., arXiv May 2020
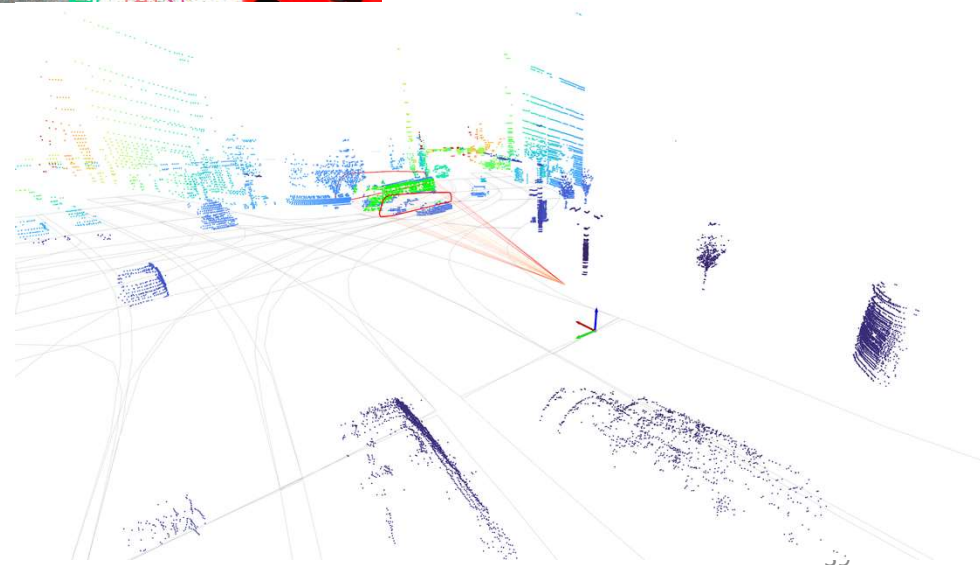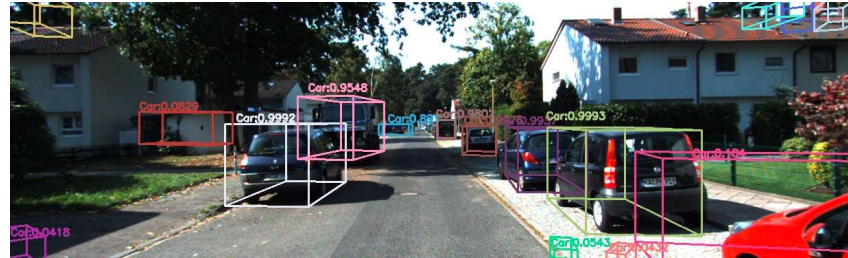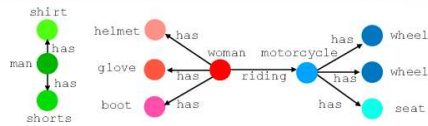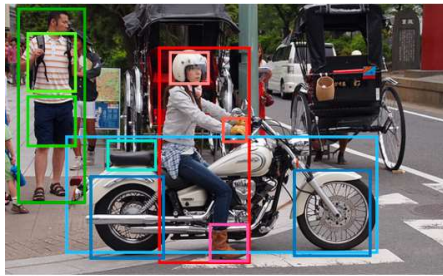
# Experiments



| | 1% | 2% | 5% | 10% |
|---|---|---|---|---|
| Supervised only | 9.05 | 12.70 | 18.47 | 23.86 |
| STAC [2] (SOTA from Google) | 13.97 (+4.92) | 18.25 (+5.55) | 24.38 (+5.91) | 28.64 (+4.78) |
| Ours | 20.80 (+11.75) | 24.18 (+11.48) | 28.23 (+9.76) | 31.54 (+7.68) |

[2]"A Simple Semi-Supervised Learning Framework for Object Detection", Sohn et al., arXiv May 2020

# Summary

- Perform the state-of-art on semi-supervised object detection
- Different tasks based on object detection can benefit from this model
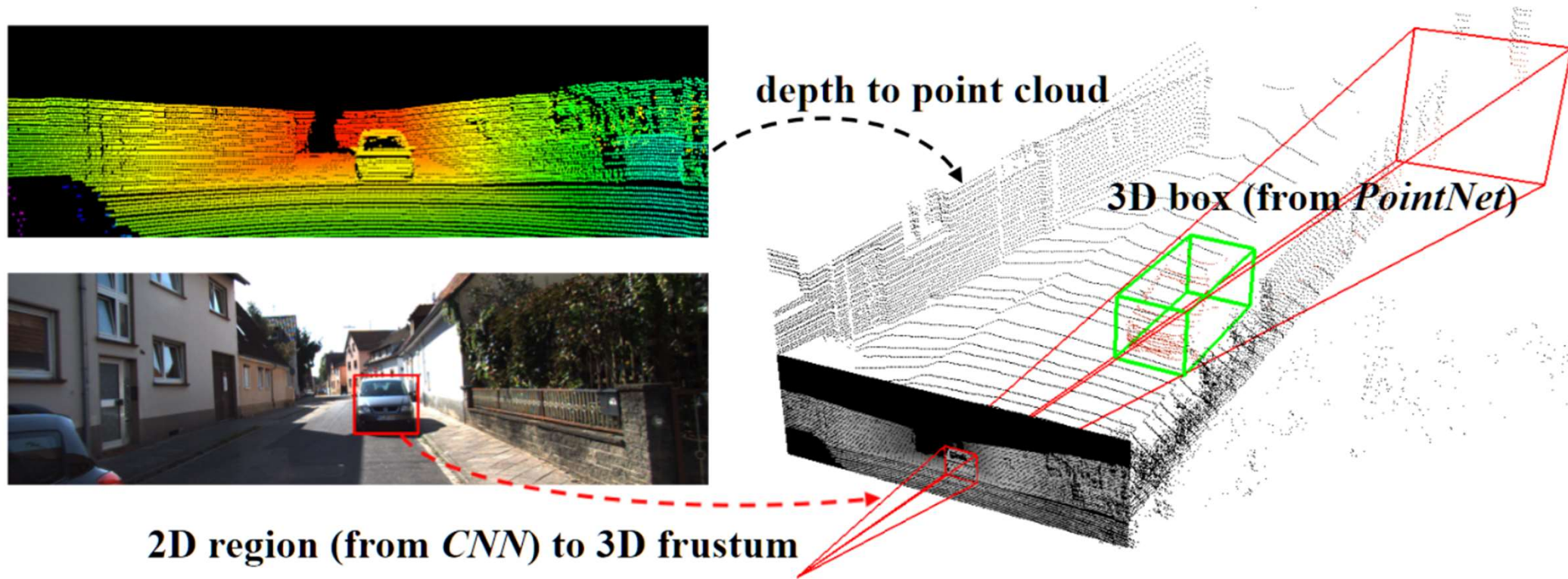


What about 3D?

3D Object Detection Goal:
Recover the *amodal* 3D spatial extent and heading of objects in a scene.



Can we "inflate" 2D instance segmentations into 3D cuboids?

Wilson et al., 3D for Free: Crossmodal Transfer Learning using HD Maps, https://arxiv.org/abs/2008.10592

Frustum PointNets for 3D Object Detection from RGB-D Data.
Charles R. Qi, Wei Liu, Chenxia Wu, Hao Su, Leonidas J. Guibas. 2017.

Slide Credit: Ben Wilson & James Hays / Argo AI

Detections from: Youngwan Lee and Jong
Park. CenterMask: Real-Time Anchor-Free
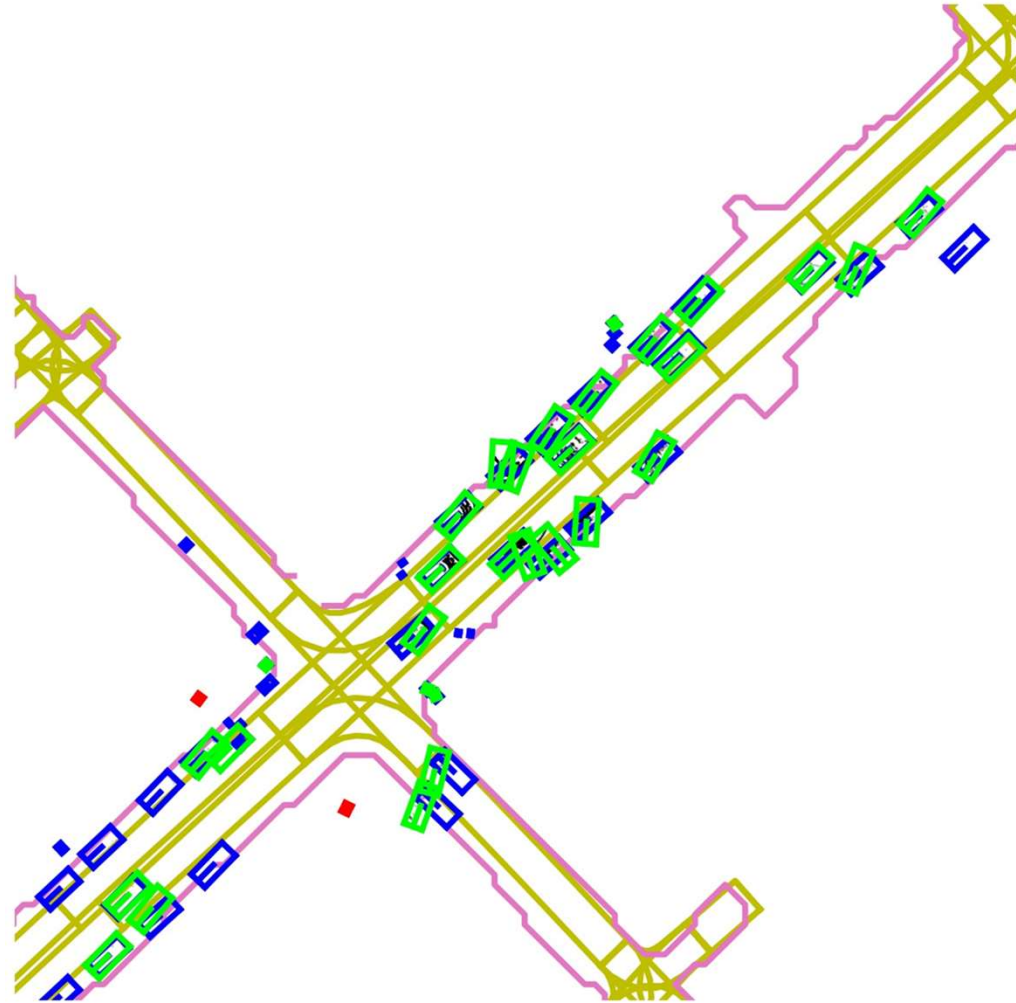Instance Segmentation. November 2019.
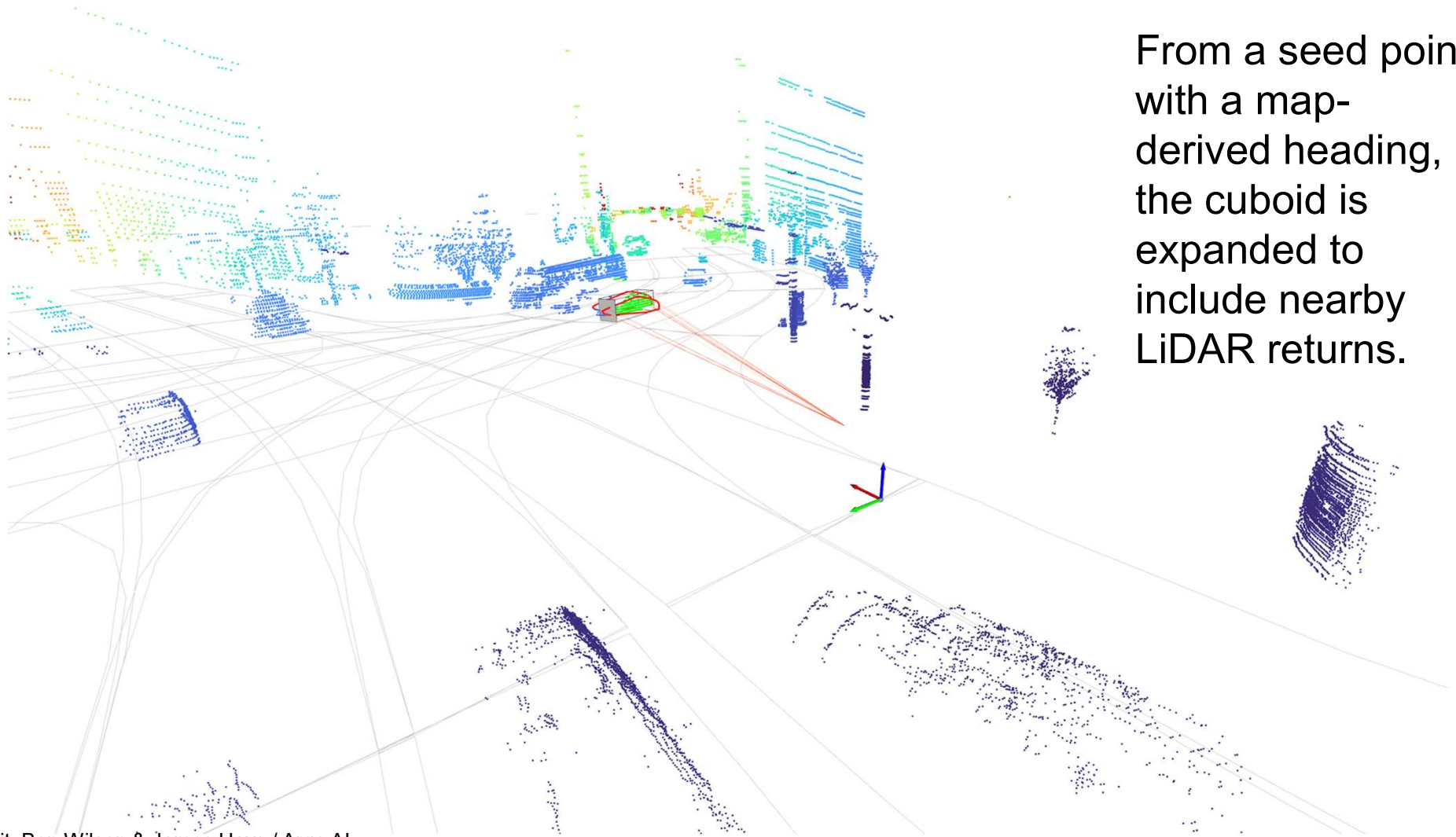
Slide Credit: Ben Wilson & James Hays / Argo AI

All LiDAR points in the instance segmentation frustum

# Using HD Maps (Centerlines)

- We use *centerlines* to improve orientation estimates



Slide Credit: Ben Wilson & James Hays / Argo AI

From a seed point, with a map-derived heading, the cuboid is expanded to include nearby LiDAR returns.

Slide Credit: Ben Wilson & James Hays / Argo AI

If the cuboid is too small, its *amodal* extent is estimated by growing the cuboid

Slide Credit: Ben Wilson & James Hays / Argo AI

Slide Credit: Ben Wilson & James Hays / Argo AI

Slide Credit: Ben Wilson & James Hays / Argo AI

Slide Credit: Ben Wilson & James Hays / Argo AI

Slide Credit: Ben Wilson & James Hays / Argo AI

LiDAR baseline is Benjin Zhu, Zhengkai Jiang, Xiangxin Zhou, Zeming Li, and Gang Yu. Class-balanced Grouping and Sampling for Point Cloud 3D Object Detection. arXiv, August 2019



| Method | Human Annotation | ↑ mAP | ↓ mATE | ↓ mASE | ↓ mAOE | ↑ CDS |
|---|---|---|---|---|---|---|
| Supervised Baseline | ✓ | 23.76 | 0.39 | **0.22** | 0.86 | 18.48 |
| Inflating with Rotating Calipers | | 23.62 | 0.55 | 0.32 | 1.68 | 15.18 |
| Inflating with with Map Mined Data | | 27.66 | 0.44 | 0.30 | **0.77** | 19.72 |

# If we *mine* 1,151 unlabeled vehicle logs...

# Deep Learning is Robust to Massive Label Noise

David Rolnick [*1]  Andreas Veit [*2]  Serge Belongie [2]  Nir Shavit [3]

94v3 [cs.LG] 26 Feb 2018

## Abstract

Deep neural networks trained on large supervised datasets have led to impressive results in image classification and other tasks. However, well-annotated datasets can be time-consuming and expensive to collect, lending increased interest to larger but noisy datasets that are more easily obtained. In this paper, we show that deep neural networks are capable of generalizing from training data for which true labels are massively outnumbered by incorrect labels. We demonstrate remarkably high test performance after training on corrupted data from MNIST, CIFAR, and ImageNet. For example, on MNIST we obtain test accuracy above 90 percent even after each clean training example has been diluted with 100 randomly-labeled examples. Such behavior holds across multiple patterns of label noise, even when erroneous labels are biased towards confusing classes. We show that training in this regime requires a significant but manageable increase in dataset size

Thus, annotation can be expensive and, for tasks requiring expert knowledge, may simply be unattainable at scale.

To address this limitation, other training paradigms have been investigated to alleviate the need for expensive annotations, such as unsupervised learning (Le, 2013), self-supervised learning (Pinto et al., 2016; Wang & Gupta, 2015) and learning from noisy annotations (Joulin et al., 2016; Natarajan et al., 2013; Veit et al., 2017). Very large datasets (e.g., Krasin et al. (2016); Thomee et al. (2016)) can often be obtained, for example from web sources, with partial or unreliable annotation. This can allow neural networks to be trained on a much wider variety of tasks or classes and with less manual effort. The good performance obtained from these large, noisy datasets indicates that deep learning approaches can tolerate modest amounts of noise in the training set.

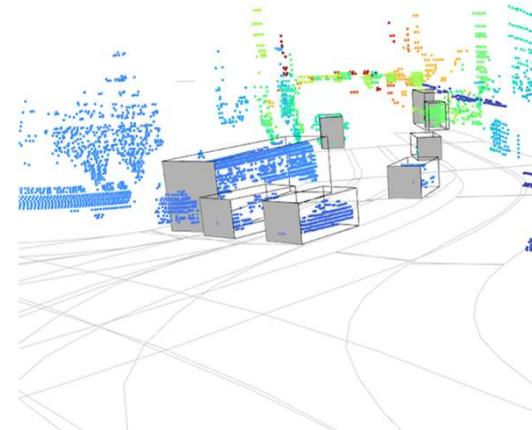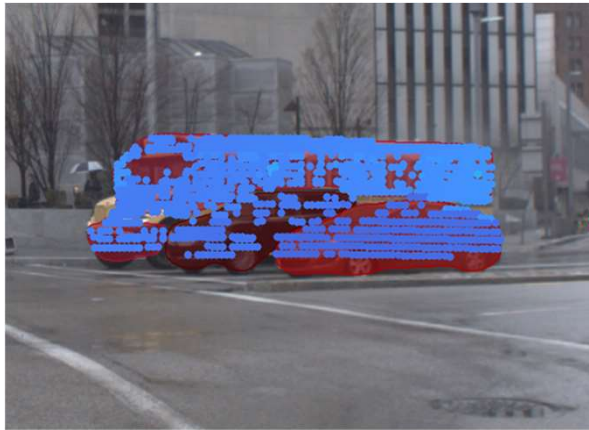In this work, we study the behavior of deep neural networks under extremely low label reliability, only slightly above chance. The insights from our study can help guide future settings in which arbitrarily large amounts of data are easily obtainable, but in which labels come without any guarantee.

David Rolnick, Andreas Veit, Serge Belongie, and Nir Shavit. **Deep Learning is Robust to Massive Label Noise**. arXiv, May 2017.

# Learning from mined, inflated data

| Method | Human Annotation | ↑ mAP | ↓ mATE | ↓ mASE | ↓ mAOE | ↑ CDS |
|---|---|---|---|---|---|---|
| Supervised Baseline | ✓ | 23.76 | 0.39 | **0.22** | 0.86 | 18.48 |
| Inflating with Rotating Calipers | | 23.62 | 0.55 | 0.32 | 1.68 | 15.18 |
| Inflating with with Map Mined Data | | 27.66 | 0.44 | 0.30 | **0.77** | 19.72 |
| Training with Map Mined Data (ours) | | **30.56** | 0.39 | 0.28 | 1.01 | **22.18** |



Slide Credit: Ben Wilson & James Hays / Argo AI

# Learning from mined, inflated data

| Method | Human Annotation | ↑ mAP | ↓ mATE | ↓ mASE | ↓ mAOE | ↑ CDS |
|---|---|---|---|---|---|---|
| Supervised Baseline | ✓ | 23.76 | 0.39 | **0.22** | 0.86 | 18.48 |
| Inflating with Rotating Calipers | | 23.62 | 0.55 | 0.32 | 1.68 | 15.18 |
| Inflating with with Map Mined Data | | 27.66 | 0.44 | 0.30 | **0.77** | 19.72 |
| Training with Map Mined Data (ours) | | **30.56** | 0.39 | 0.28 | 1.01 | **22.18** |

| Method | Vehicle | Bus | Motorcycle | Bicycle | Pedestrian | ↑ mAP | ↑ CDS |
|---|---|---|---|---|---|---|---|
| Supervised Baseline | **58.50** | **12.70** | 2.70 | 2.10 | **42.80** | 23.76 | 18.48 |
| Training with Map-Mined Data (ours) | 52.40 | 12.10 | **17.00** | **35.50** | 35.80 | **30.56** | **22.18** |

Slide Credit: Ben Wilson & James Hays / Argo AI
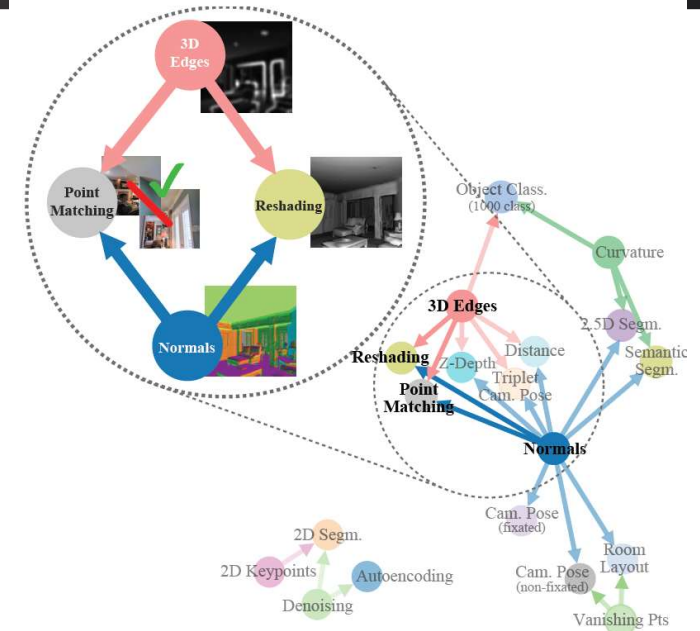
CDS = mAP  (mATE + mASE + mAOE)

Slide Credit: Ben Wilson & James Hays / Argo AI
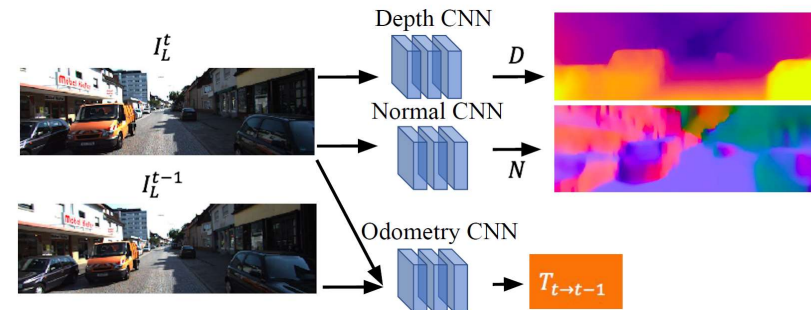
# The Methods are Surprisingly Simple

- A handful of common techniques:
  - **Data augmentation**
  - **Pseudo-labeling** / distillation
  - Surrogate tasks / contrastive losses
  - Temperature scaling / Entropy maximization
  - Cosine/metric learning
  - **Prototypes**
  - **Graph neural networks**
  - Meta-learning

- Autonomous vehicles: **Multi-modal and multi-task!**

# Conclusions

- Amazing gains have been made across learning with limited labels

- **Data augmentation** a crucial aspect; we develop methods for:
  - **Complex feature-space augmentation** using graphs and leveraging manifold structure

- Move **beyond image classification** for autonomous vehicles
  - Object Detection
  - 3D
  - Leverage large amount of unlabeled data though *many tasks*



Zamir et al., Taskonomy, CVPR 2018



Zhan et al., Self-supervised Learning for Single View Depth and Surface Normal Estimation, CVPR 2018

# Acknowledgement and Questions?

Chia-Wen Kuo

Ben Wilson

Chih-Yao Ma

Yen-Chang Hsu

Yen-Cheng Liu