



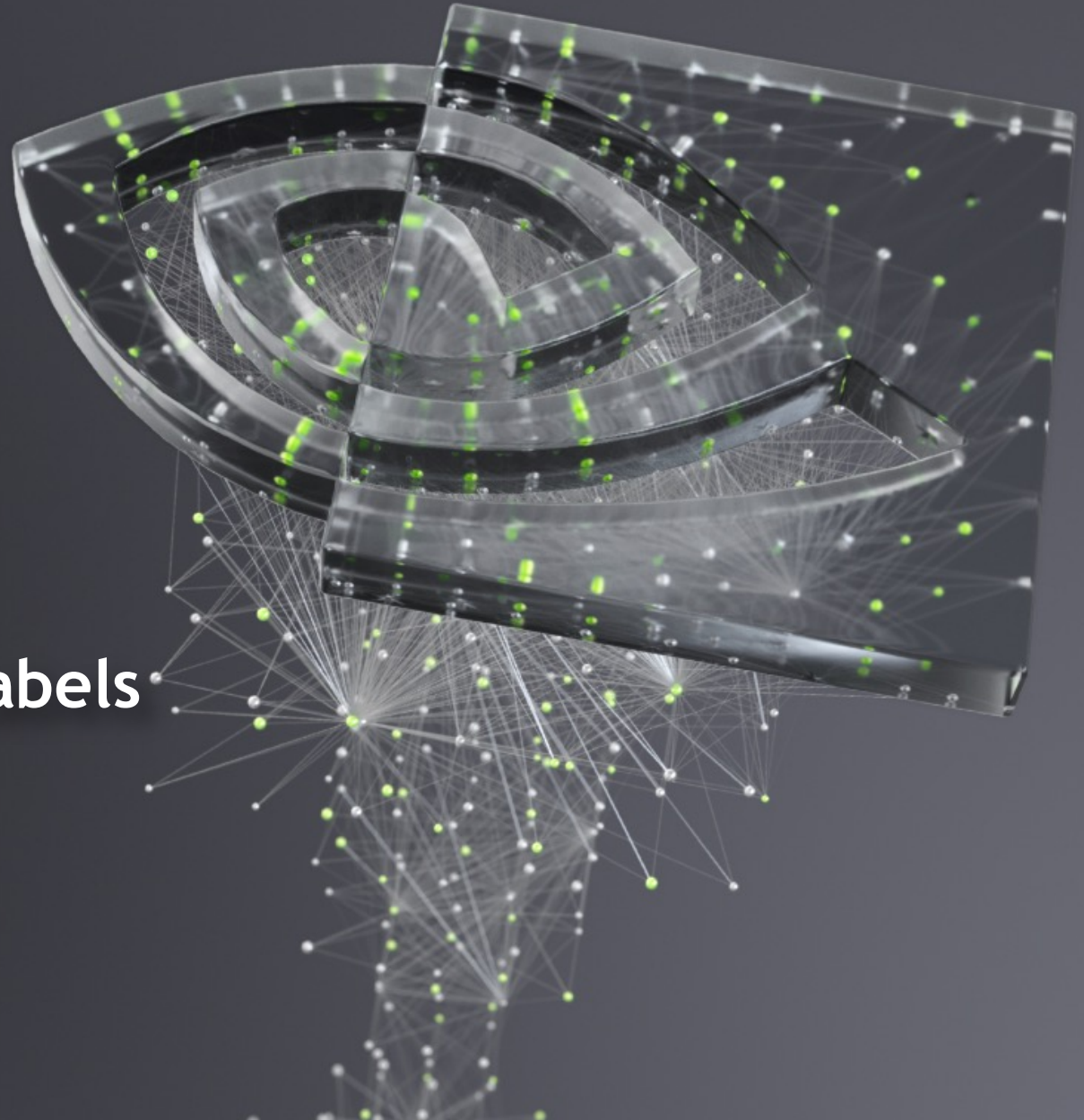
NVIDIA



Caltech

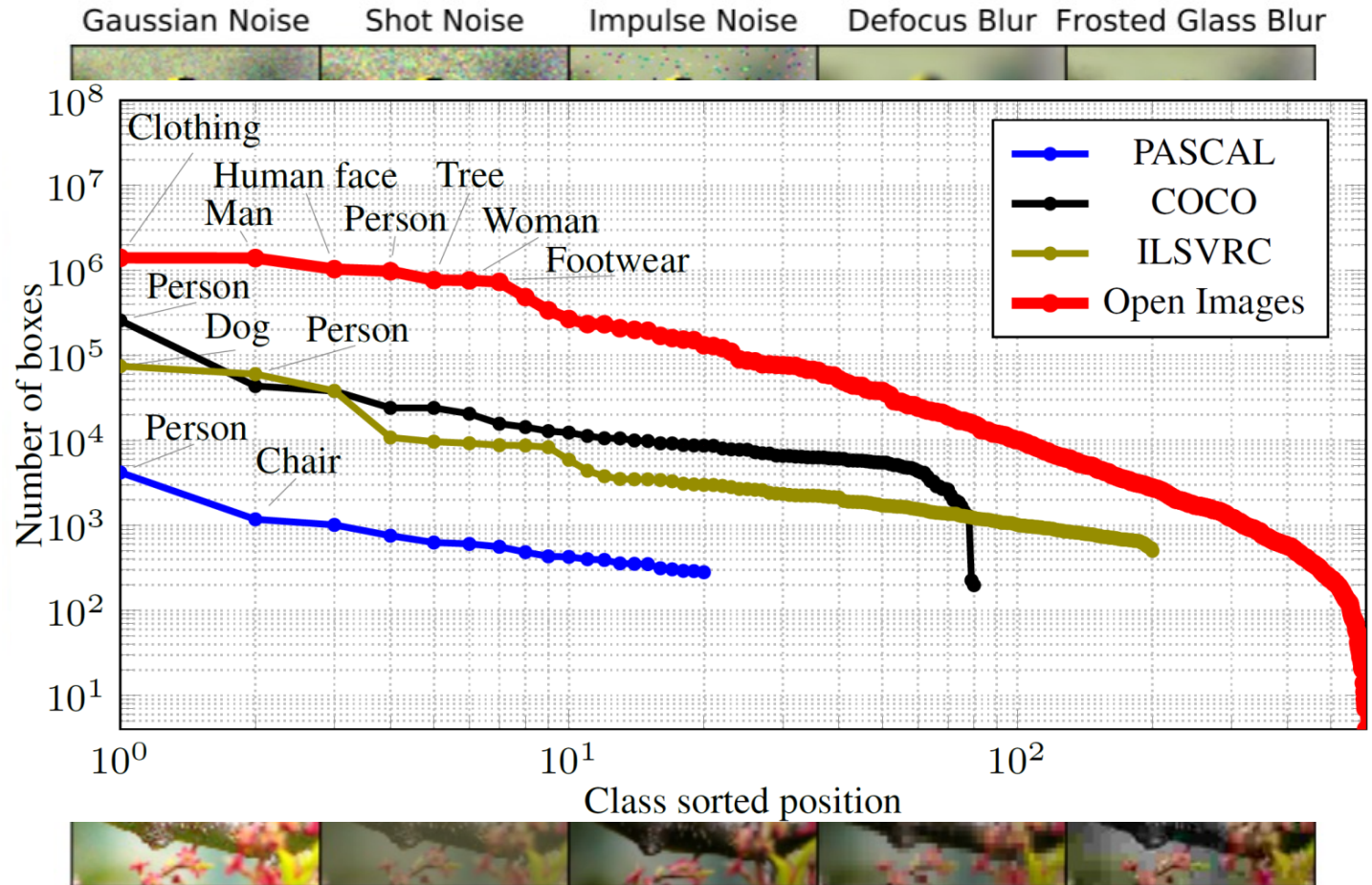
Learning with Imperfect Labels and Visual Data

Anima Anandkumar



Challenge I - Real World Data Are Imperfect

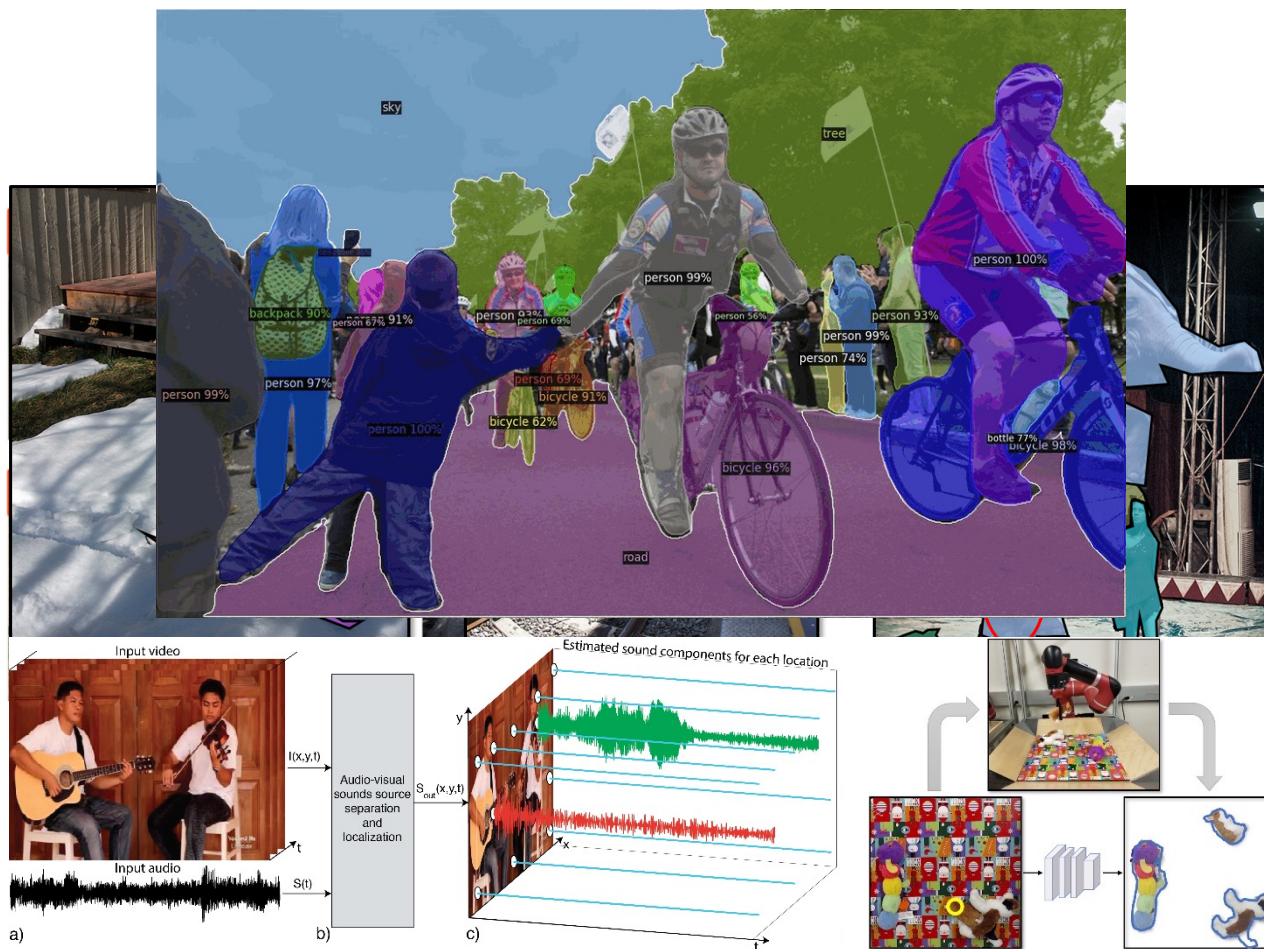
- Domain gap
- Data bias
- Data noise
- Can be *Long tail*
- Can contain *Occlusions*
- Can be *Cluttered*
- Can be *Ambiguous*
- Can be *Multisensory*
- ...



Challenge II - Real World Labels Are Imperfect

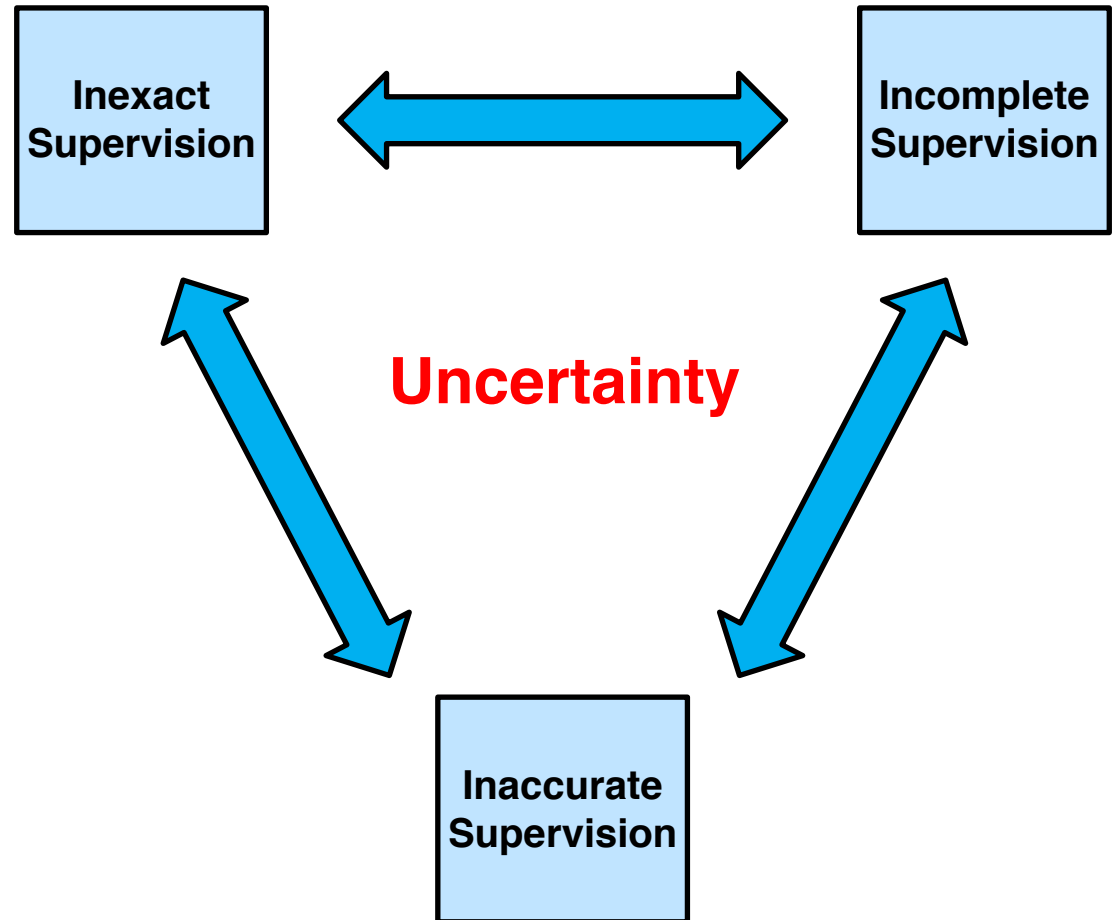
- Inexact (indirect) supervision
- Incomplete (limited) supervision
- Inaccurate (noisy) supervision
- Can be *Versatile*
- Can be *Multimodal*
- Can be *Sequential*
- Can be *Sparse*
- Can be *Interactive*

•••

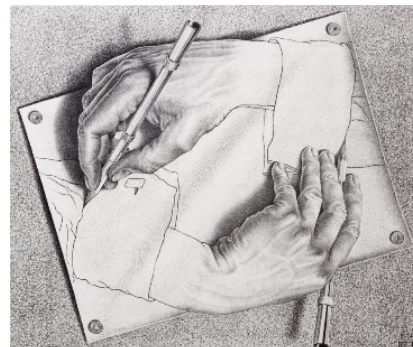
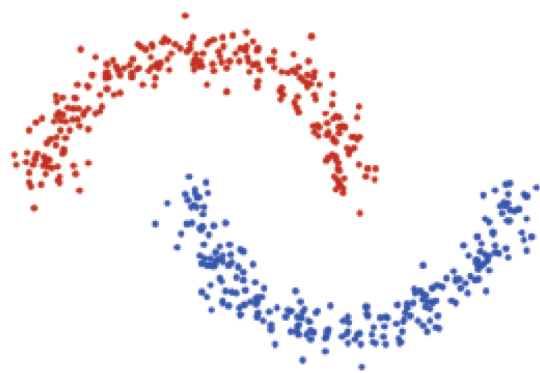
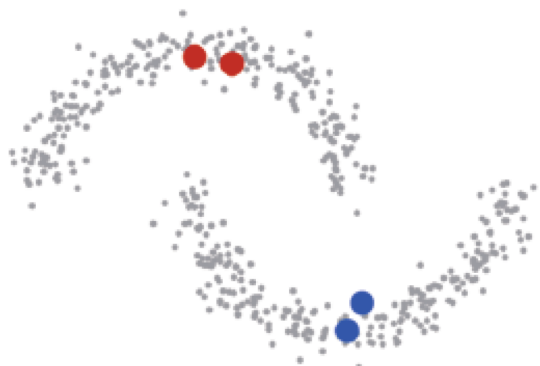


Goal and Challenge

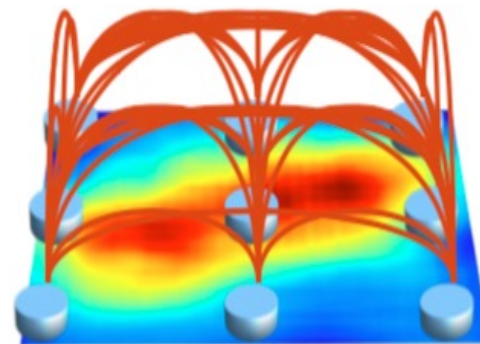
- The ability to improve generalization by learning continuously on new data
- The ability to leverage diverse forms of weak supervision, or self-improving on unlabeled data
- The ability to overcome the **uncertainty** (ill-posed nature) due to the lack of constraints under partial supervision



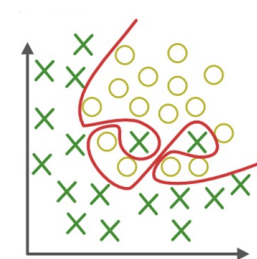
Overcoming Uncertainties from Imperfect Labels



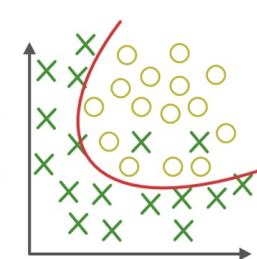
Self-supervision



Structuredness

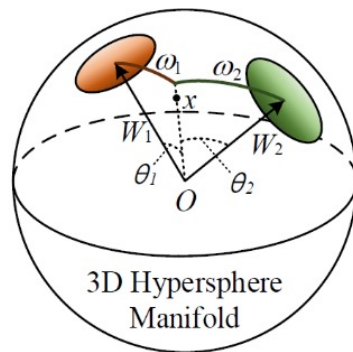


Over-fitting
(forcefitting--too good to be true)

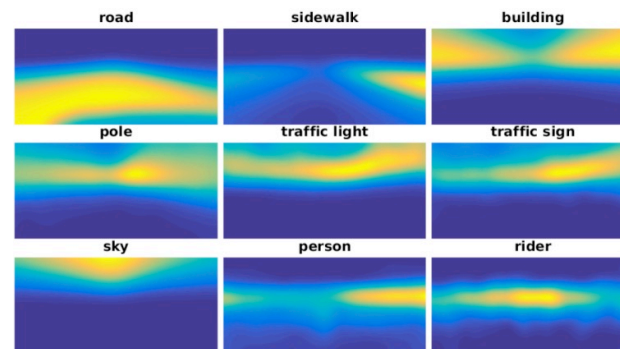


Appropriate-fitting

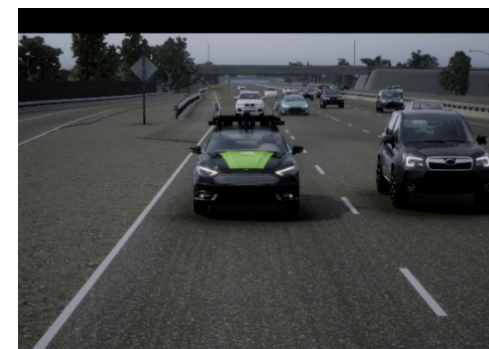
Regularization



Inductive Bias



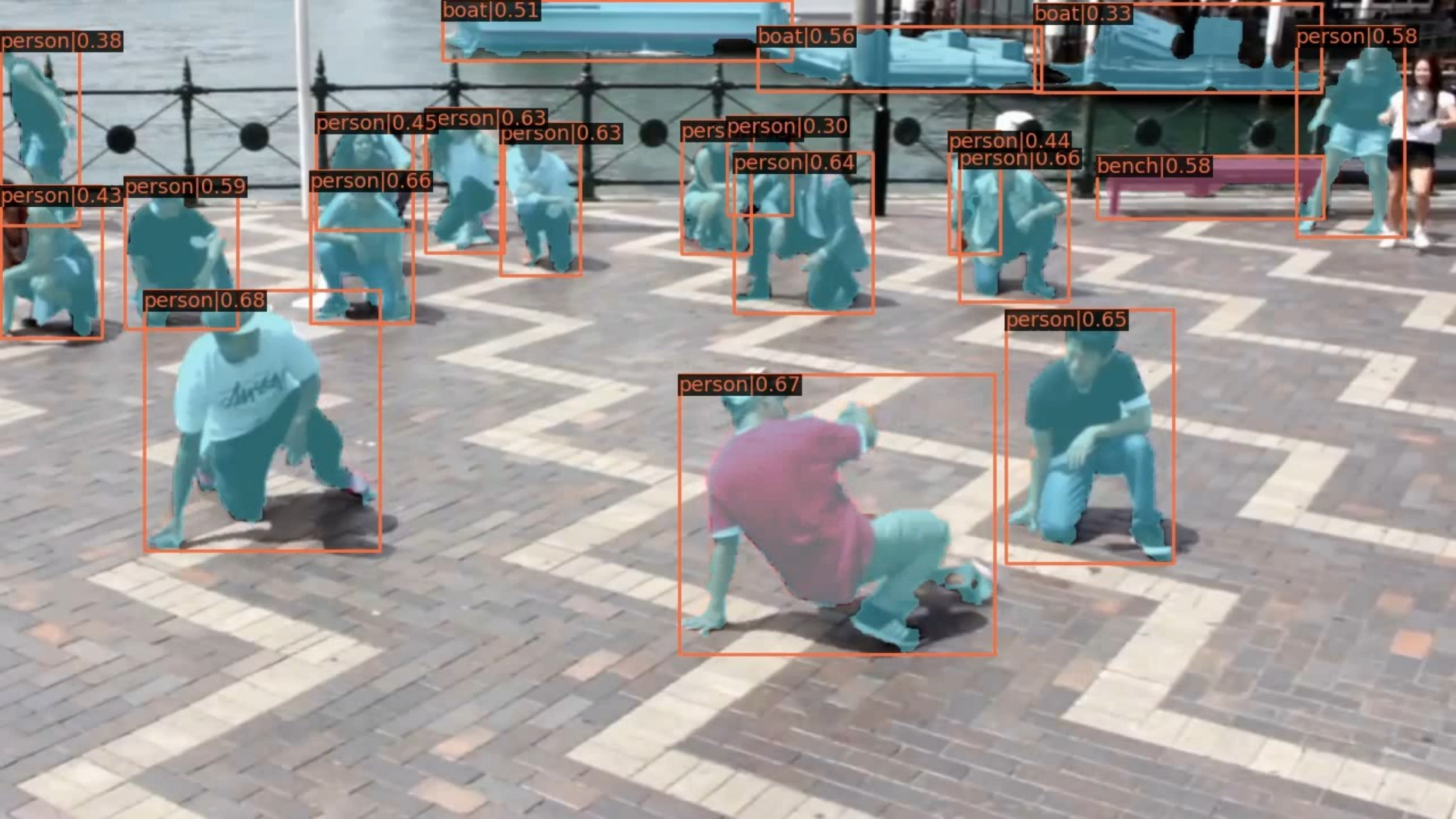
Domain Prior



Synthetic Data



DISCOBOX



person|0.38

boat|0.51

boat|0.56

boat|0.33

person|0.58

person|0.45

person|0.63

person|0.63

person|0.30

person|0.44

person|0.66

bench|0.58

person|0.43

person|0.59

person|0.66

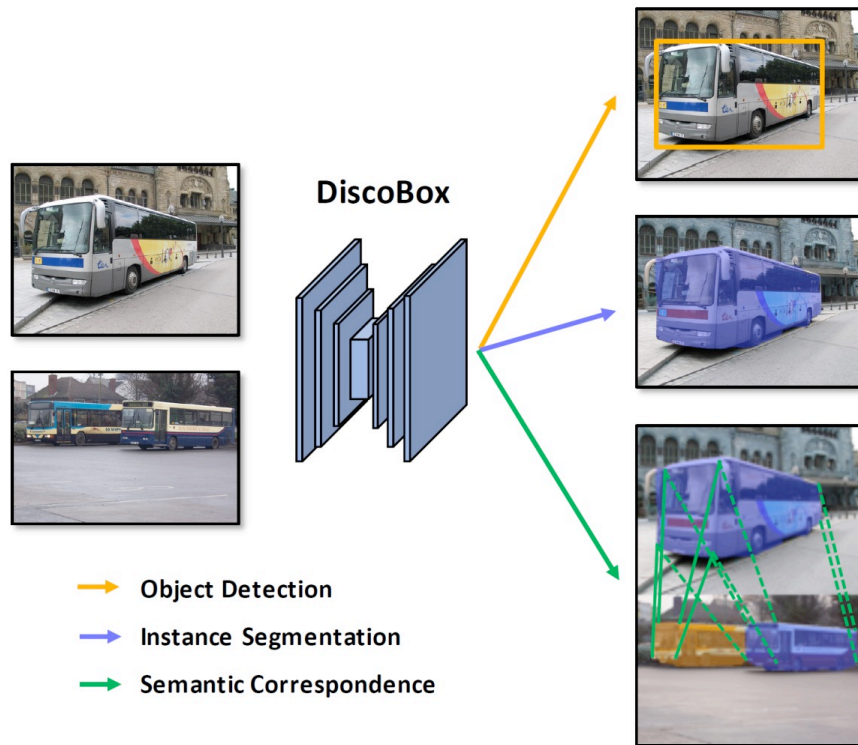
person|0.64

person|0.68

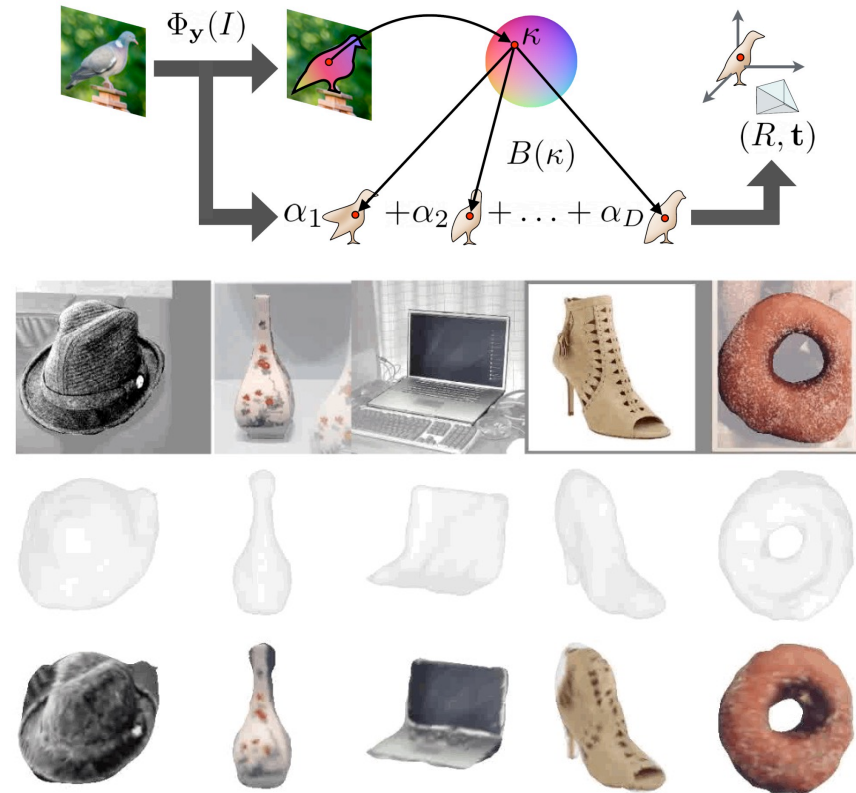
person|0.65

person|0.67

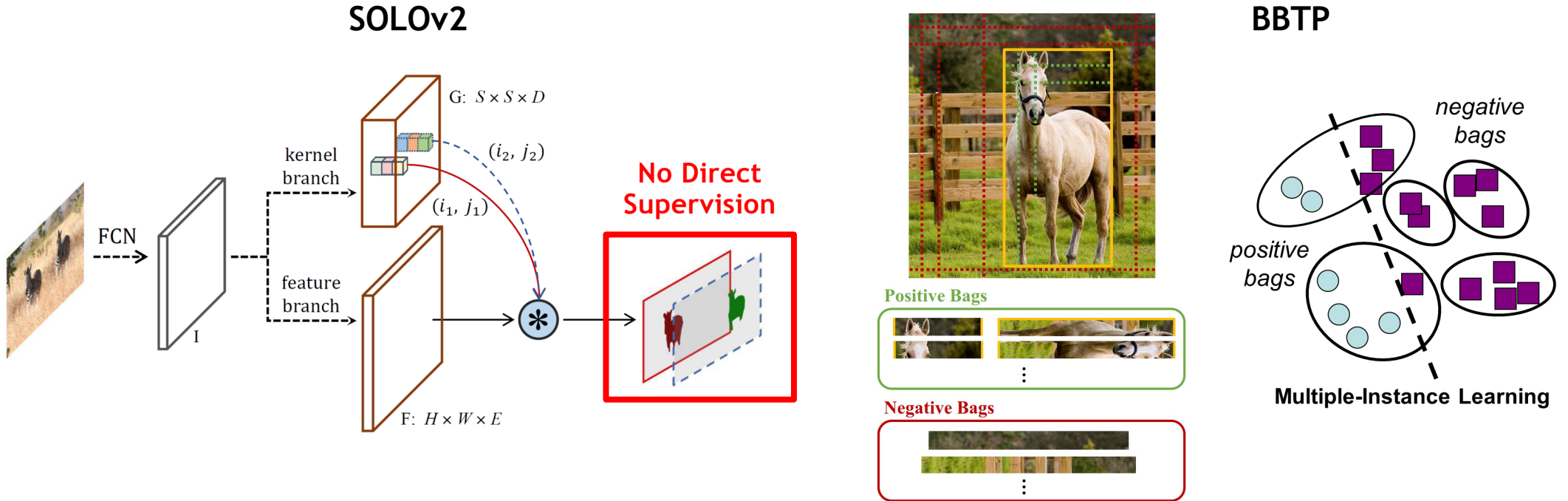
Method Overview



- Object Detection
- Instance Segmentation
- Semantic Correspondence



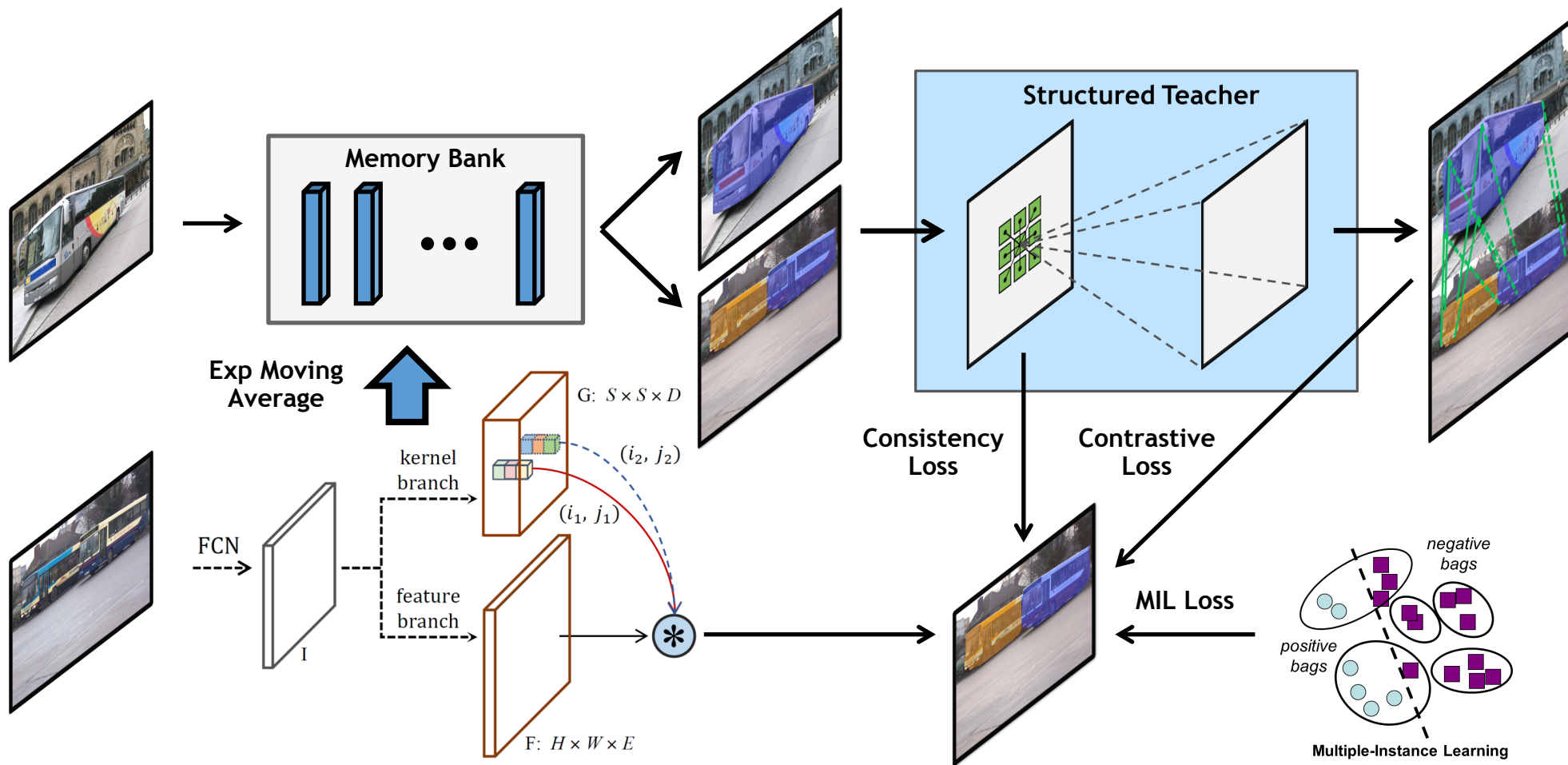
Task Network Design



Wang et al., SOLOv2: Dynamic and Fast Instance Segmentation, NeurIPS20

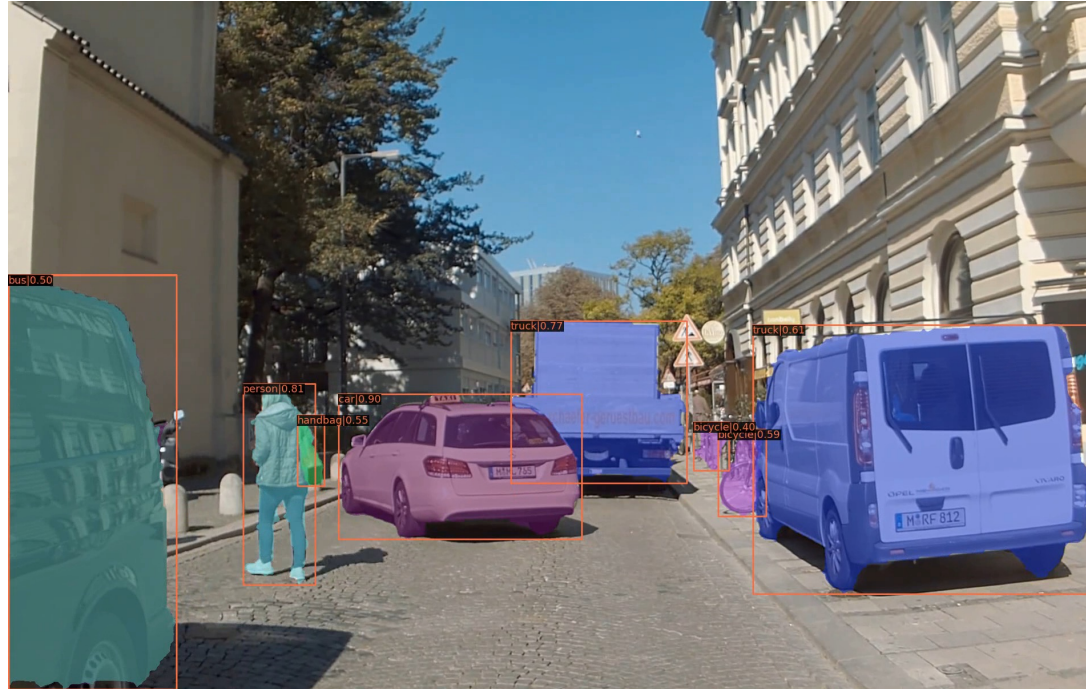
Hsu et al., Weakly supervised instance segmentation using the bounding box tightness prior, NeurIPS19

Self-Ensembling with Structured Teacher



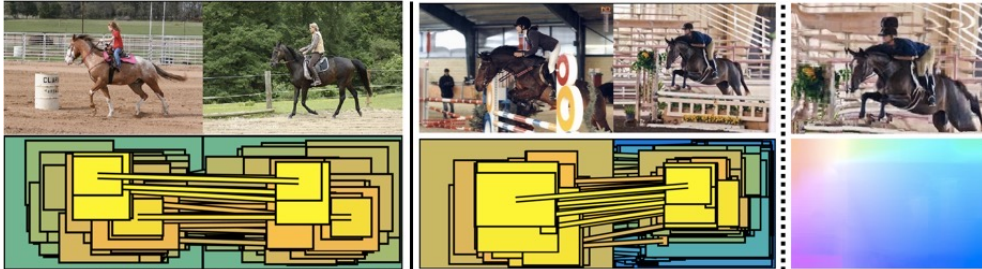
Instance Segmentation

NVIDIA DriveAV Data

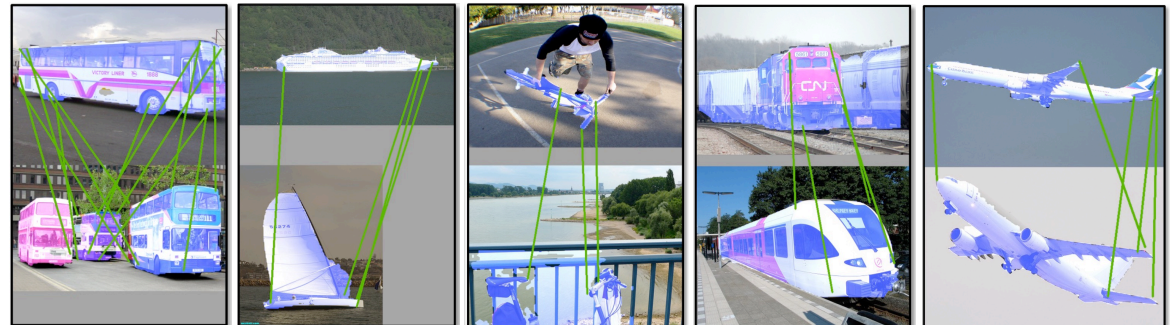
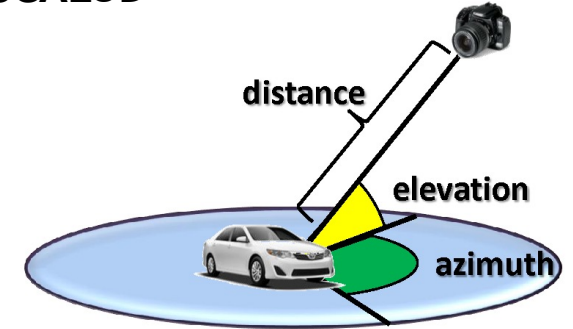
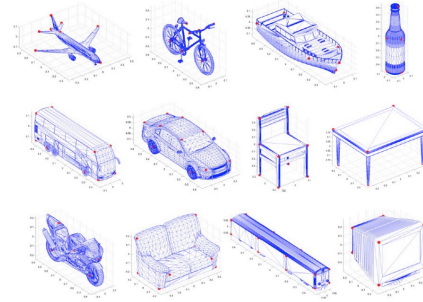


Semantic Correspondence

PF-PASCAL



PASCAL3D+



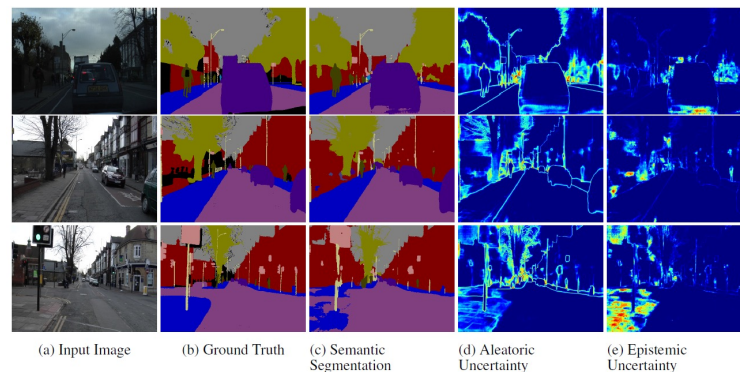
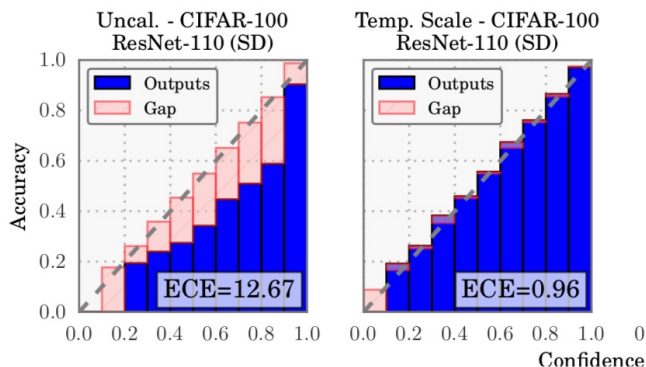
Ham et al., Proposal flow, CVPR16

Xiang et al., Beyond PASCAL: A benchmark for 3d object detection in the wild, WACV14

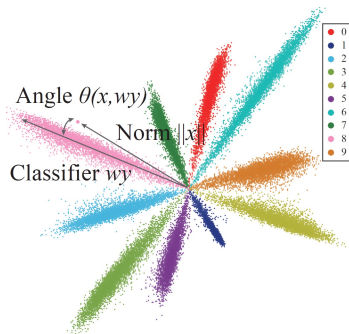


Syn-to-real Adaptation

Reliable Uncertainty Estimation Under Domain Shift

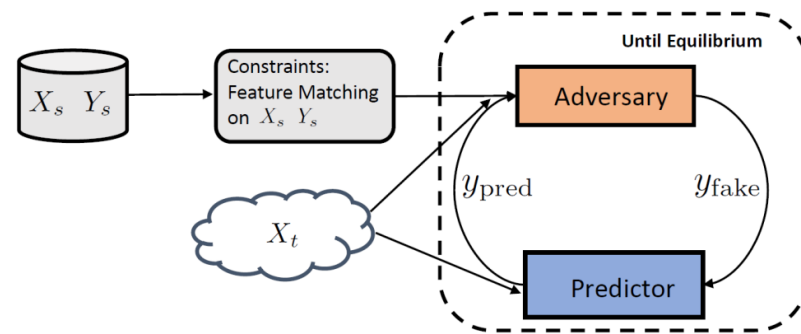


Temperature Scaling



Angular Distance

Bayesian Deep Learning



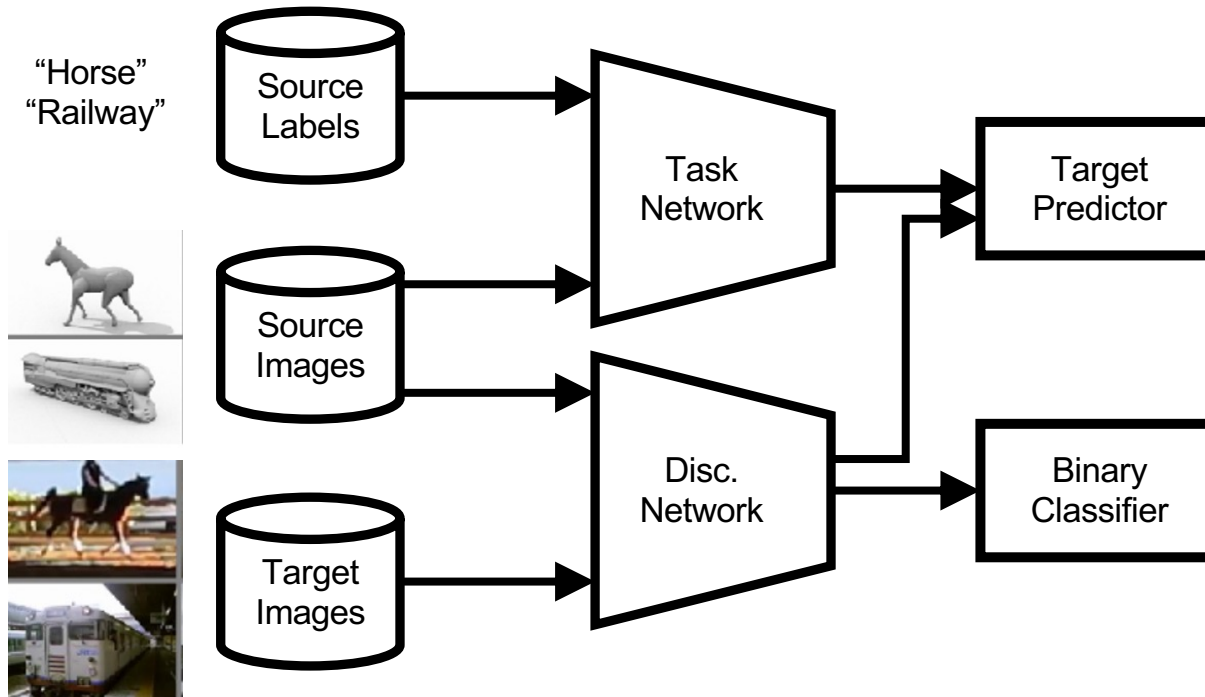
Distributionally Robust Learning

Credit: Li et al., Improving Confidence Estimates for Unfamiliar Examples, CVPR20

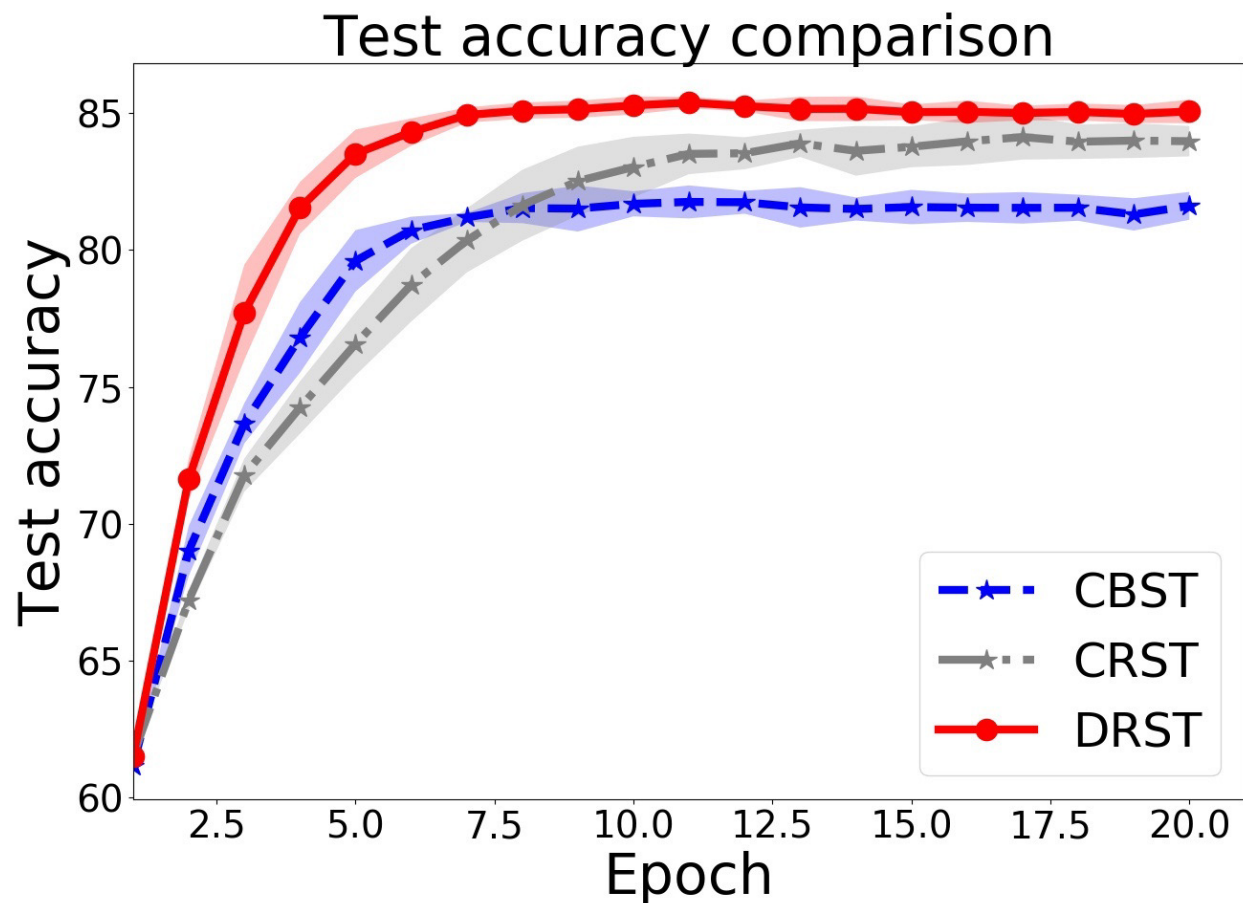
Want et al., Deep distributionally robust learning for calibrated uncertainties under domain shift, under submission at NeurIPS21

DENSITY RATIO ESTIMATION

$$\hat{P}(y|x) \propto \exp\left(\frac{P_{\text{src}}(x)}{P_{\text{trg}}(x)} \theta \cdot \phi(x, y)\right)$$



ACCURACY OF DOMAIN ADAPTION (VISDA2017)



CBST: Class-Balanced
CRST: Confidence-Regularized
DRST: Ours

We achieve the
SOTA results
in self-training
domain adaptation.

INTERPRETABLE DENSITY RATIOS

Noisy target data has lower source support

Source



A less typical “railway”



Density Ratio: 1.004

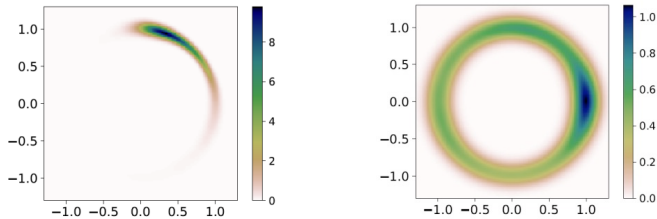
A more typical “railway”



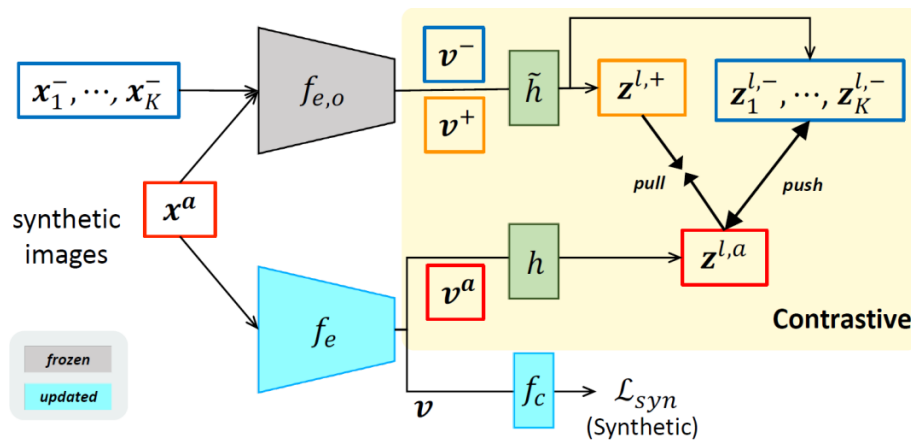
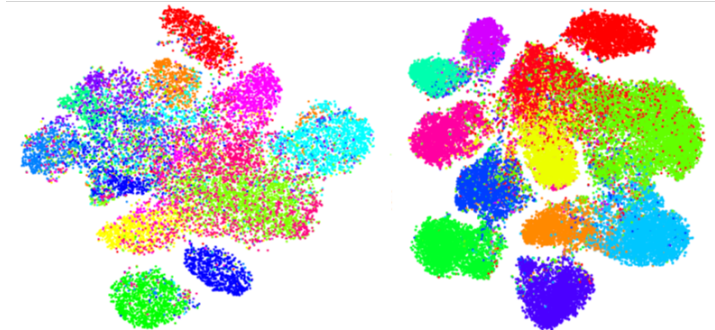
Density Ratio: 2.223

Automated Syn-to-Real Generalization

Synthetic training leads to collapsed representations.



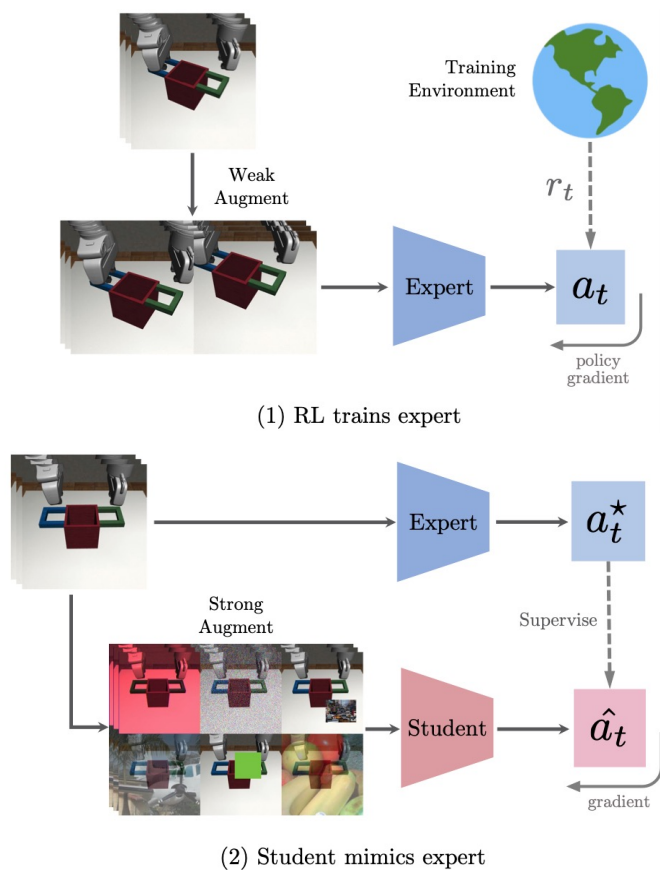
Improved representation and accuracy



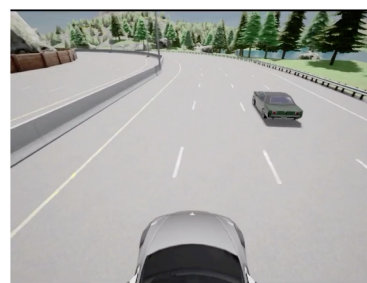
Contrastive knowledge distillation

Chen et al., Automated Synthetic-to-Real Generalization, ICML20.
Chen et al., Contrastive Syn-to-Real Generalization, ICLR21.

Sequential Decision Making



CARLA: Driving in Diverse Weathers





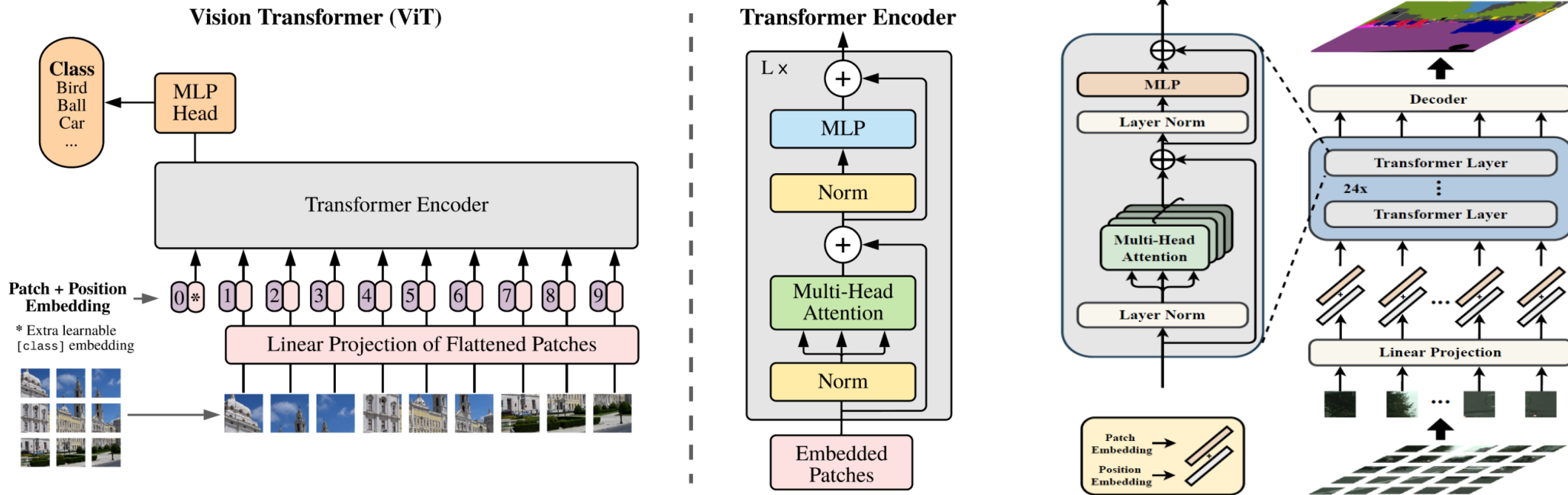
SegFormer: Semantic Segmentation with Vision Transformers

SegFormer: Simple and Efficient Design for Semantic Segmentation with Transformers

Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M. Alvarez, Ping Luo

The University of Hong Kong Nanjing University NVIDIA Caltech

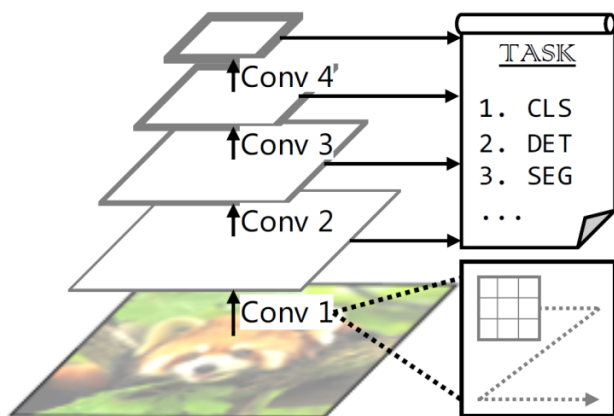
Vision Transformer (ViT)



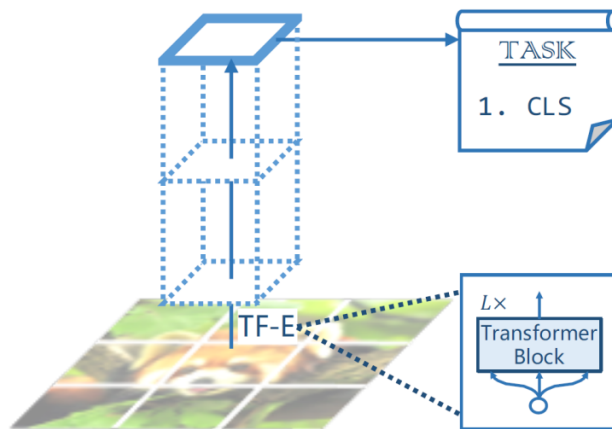
Dosovitskiy et al., An image is worth 16x16 words: Transformers for image recognition at scale, ICLR21.

Zheng et al., Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers, CVPR21.

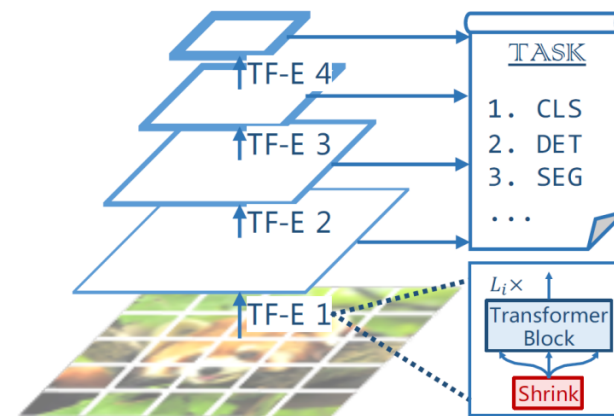
Pyramid Vision Transformer (PVT)



(a) CNNs: VGG [41], ResNet [15], etc.

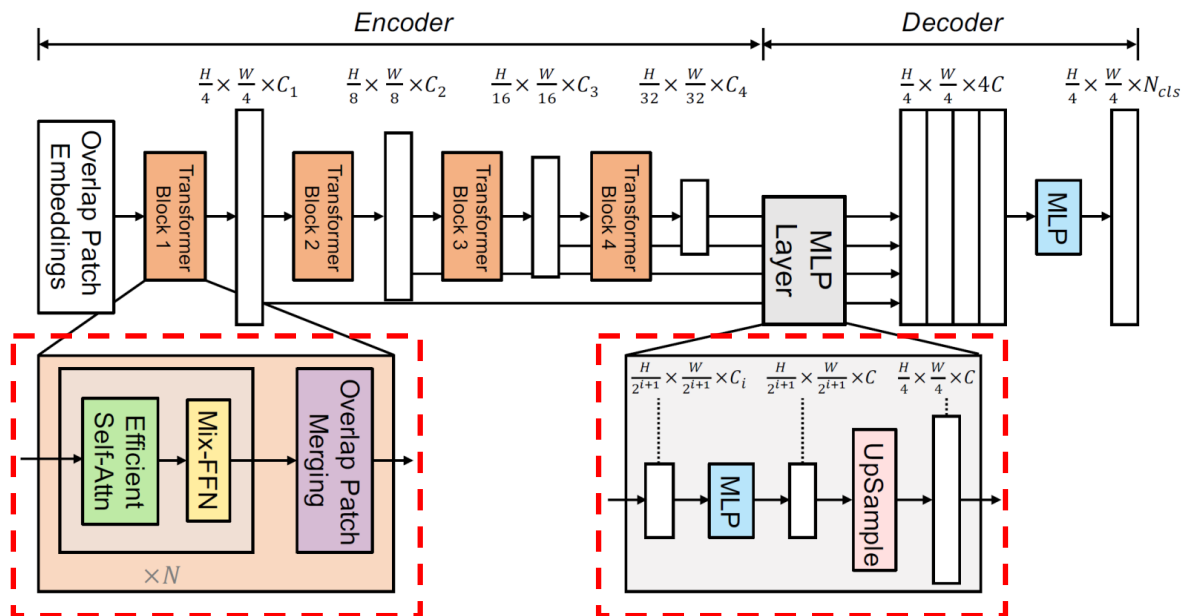


(b) Vision Transformer [10]



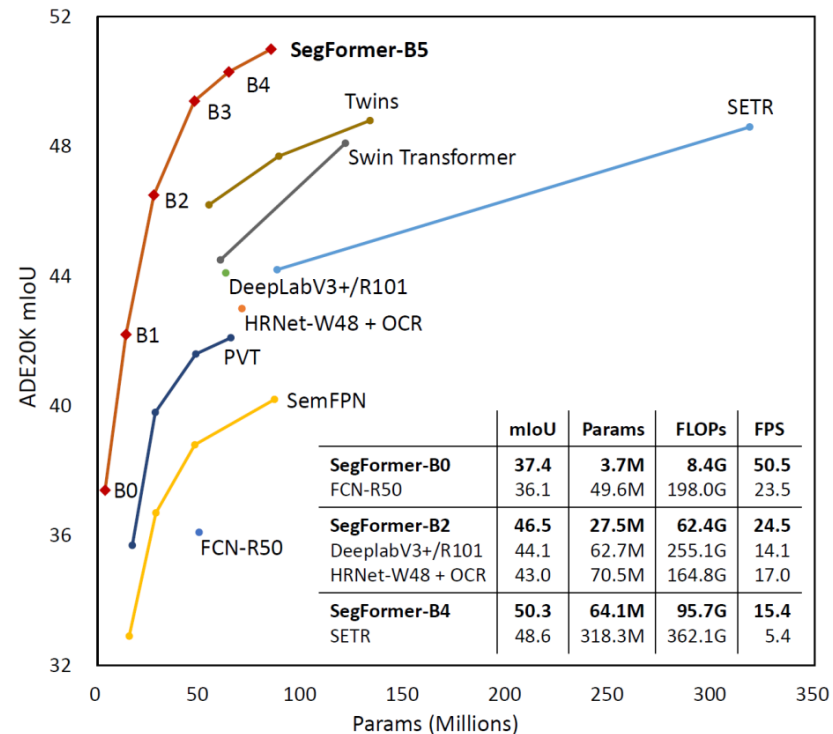
(c) Pyramid Vision Transformer (ours)

SegFormer: Architecture Overview

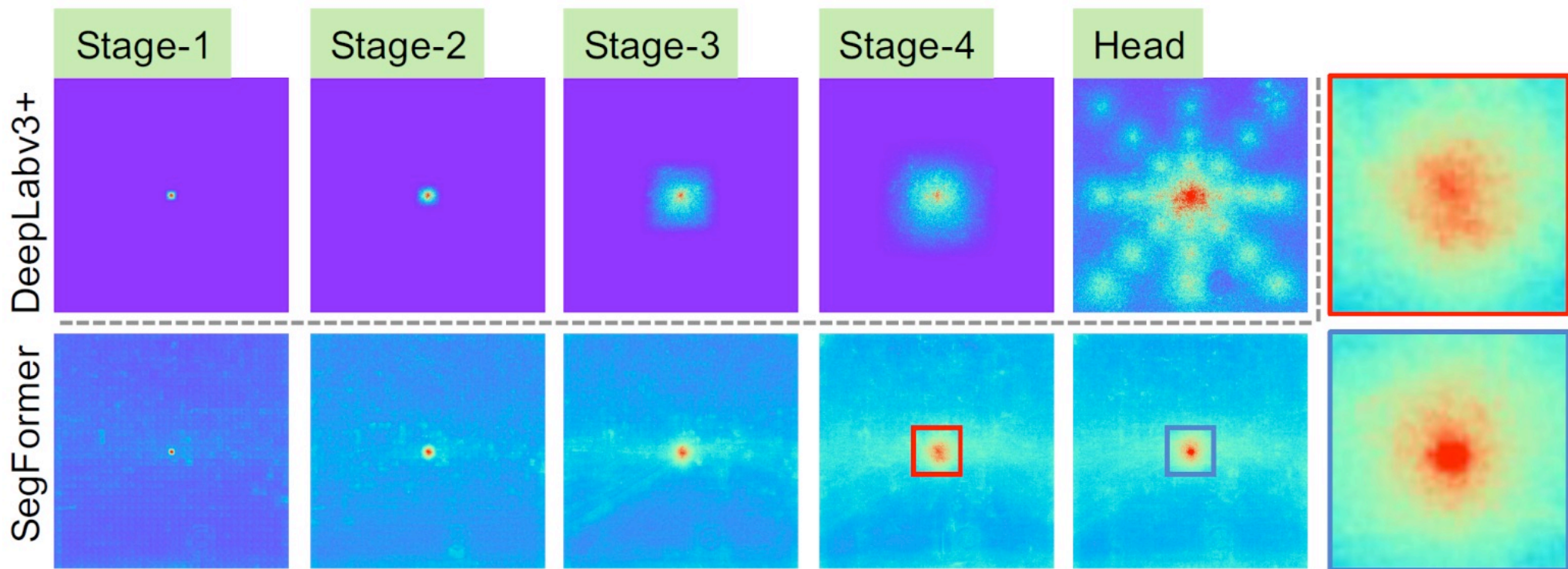


Positional encoding free,
pyramid structured encoder

Lightweight all-MLP decoder

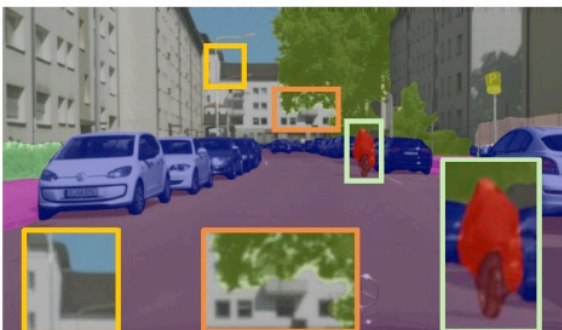


SegFormer: Decoder Design

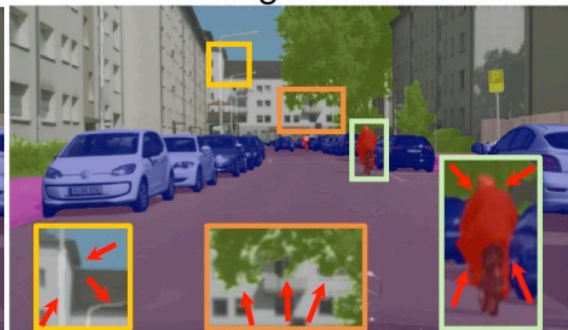


Experiment: Qualitative Results

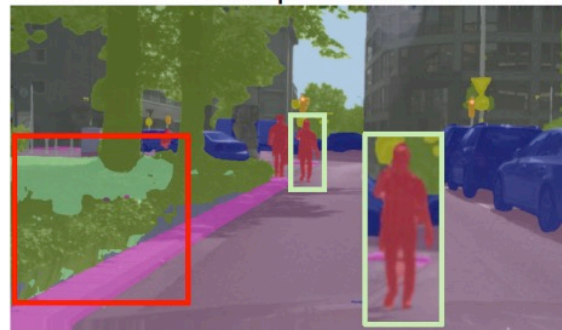
SETR



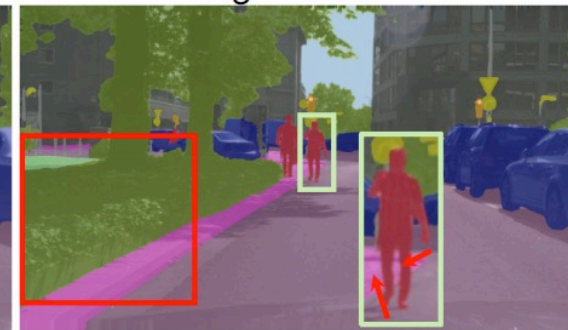
SegFormer



DeepLabv3+



SegFormer



SegFormer



SETR



DeepLabV3+

Experiment: Robustness to Corruptions

Gaussian Noise



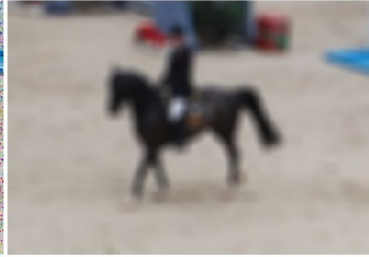
Shot Noise



Impulse Noise



Defocus Blur



Glass Blur



Motion Blur



Zoom Blur



Snow



Frost



Fog



Brightness



Contrast



Elastic Transform



Pixelate

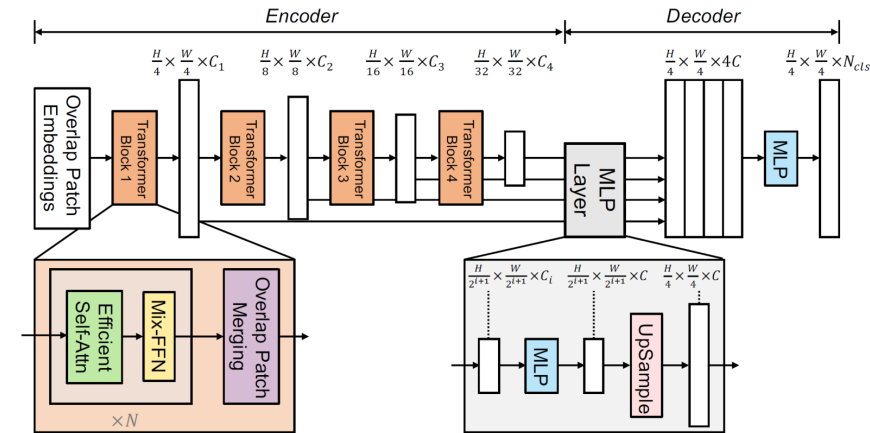
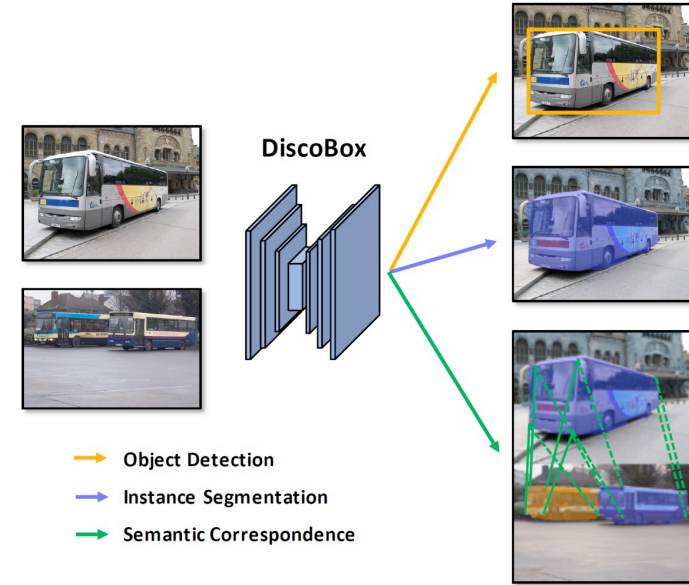


JPEG Compression



Conclusion

- Mask information probably can be totally removed in future for instance segmentation problems.
- Auto-labeling is promising for many dense prediction tasks.
- The emergence of Transformers will probably add even more to the above directions.
- Synthetic data can be seamlessly adapted to real-world tasks.





Thank You!