

Learning Generalized Intersection Over Union for Dense Pixelwise Prediction

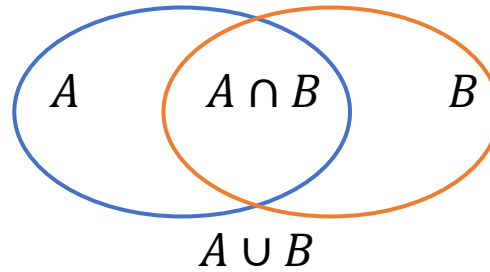
Jiaqian Yu¹, Jingtao Xu¹, Yiwei Chen¹, Weiming Li¹, Qiang Wang¹,
Byung In Yoo², Jae-Joon Han²

¹Samsung Research China - Beijing (SRCB)

²Samsung Advanced Institute of Technology (SAIT)

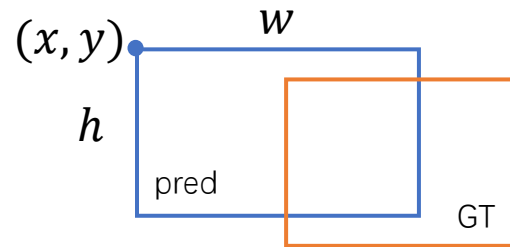
Introduction

- Intersection over Union (IoU)



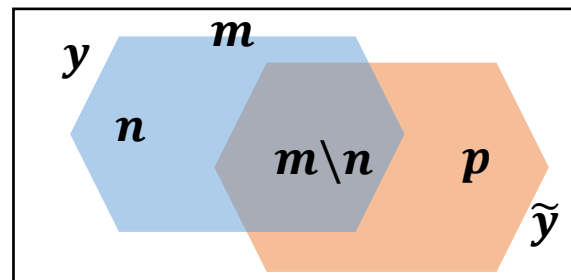
$$\text{IoU} = \frac{A \cap B}{A \cup B}$$

- Bounding box regression



$$\text{IoU} = \frac{\square \cap \square}{\square \cup \square} = \text{IoU}(x, y, w, h)$$

- Dense pixelwise prediction



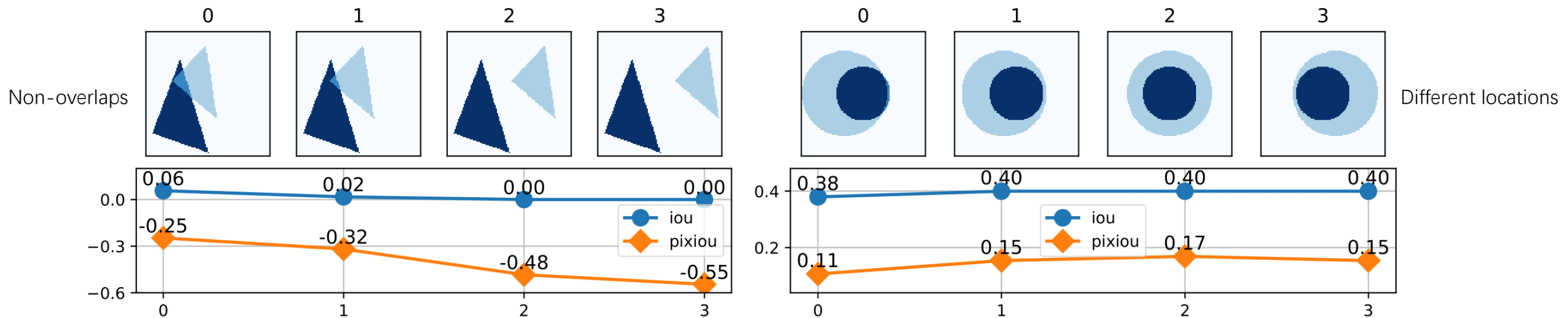
$$\text{IoU} = \frac{\text{hex} \cap \text{hex}}{\text{hex} \cup \text{hex}} = \frac{|m \setminus n|}{|m \cup p|} = \frac{|m| - |n|}{|m| + |p|}$$

$$\begin{aligned} m &= \{y = c\} \\ n &= \{y = c, \tilde{y} \neq c\} \\ p &= \{y \neq c, \tilde{y} = c\} \end{aligned}$$

Introduction

- Learning IoU: $L_{\text{iou}} = 1 - \text{IoU}$

- Issue: zero-gradient for learning if (i) no-overlaps, (ii) only different locations**



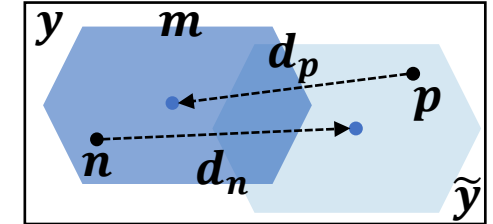
- Human cognition can judge that the optimization should be further performed.
- Learning over IoU yields suboptimal performance and leads to slower convergence.
- Our solution: PixIoU, steeper gradients for this cases.**

Method

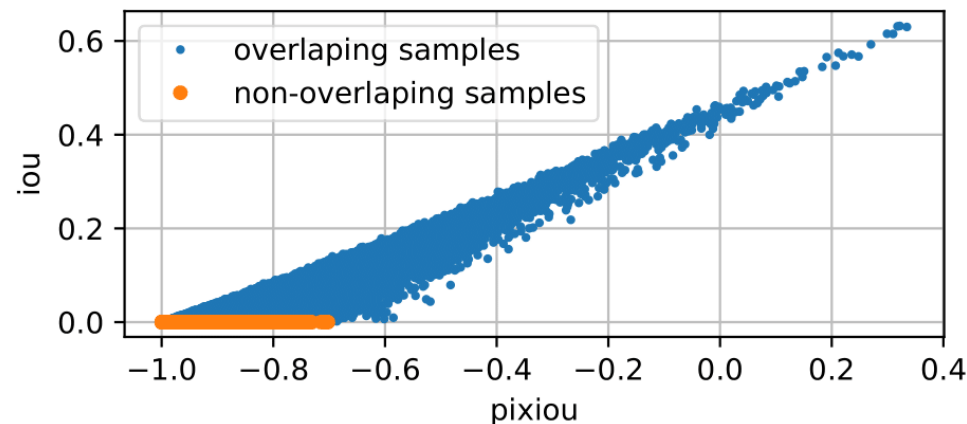
- Ingrate the coordinates of mispredicted pixels into the calculation:

$$\text{IoU} = \frac{|m| - |n|}{|m| + |p|}$$

Definition: $\text{PixIoU} = \frac{|m| - \langle d_n, \mathbf{1}_n \rangle}{|m| + \langle d_p, \mathbf{1}_p \rangle} + \text{IoU} - 1$,
 d_n, d_p is the normalized Euclidean distances to the centers.



- IoU: $O(N)$ for N pixels.
- PixIoU: $O(kN)$ with k additionally: a mean for the centers, an Euclidean, a dot product.



Properties:

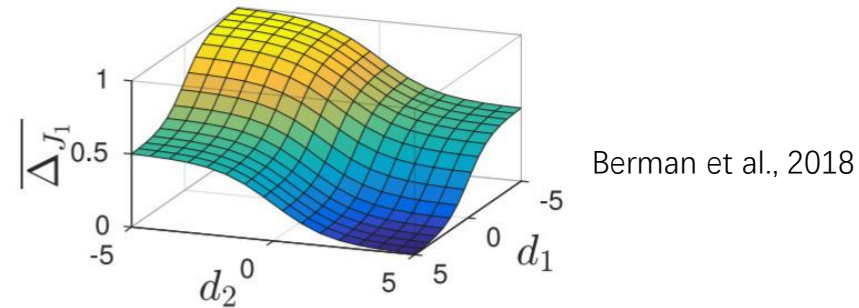
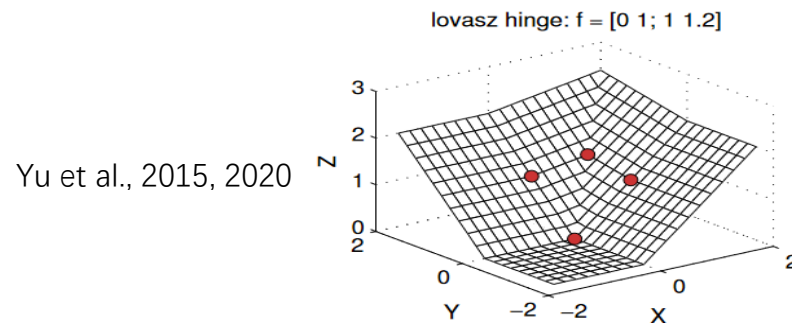
- PixIoU is invariant to the scale of the problem.
- PixIoU is always a lower bound of IoU; it becomes tighter when the predictions get better.
- PixIoU is well-bounded.

Method

- Learning PixIoU: $L_{\text{pix}} = 1 - \text{PixIoU}$

Proposition: Given a groundtruth, $L_{\text{pix}}(\cdot, A)$, is **submodular** w.r.t. the set of mispredictions of A to the groundtruth.

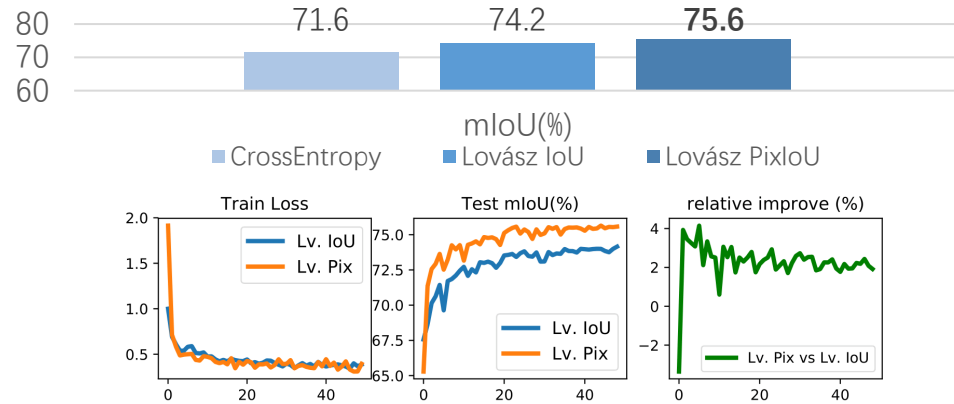
- Submodular functions: diminishing returns property.
- Efficient surrogate function: the **Lovász surrogate**, yields a convex surface, provides a polynomial computation complexity.
- Lovász Softmax (Berman et al., 2018) is proposed for learning L_{iou} .



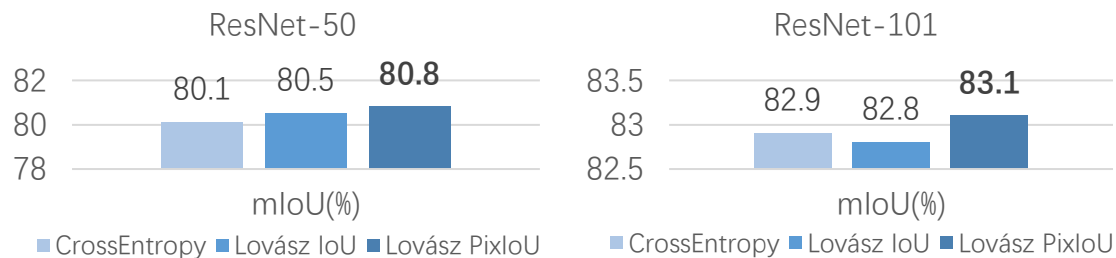
- Lovász surrogate for learning L_{pix} : additionally $\langle d_n, \mathbf{1}_n \rangle$ and $\langle d_p, \mathbf{1}_p \rangle$.

Results

- Pixelwise Object Tracking on VOT2020

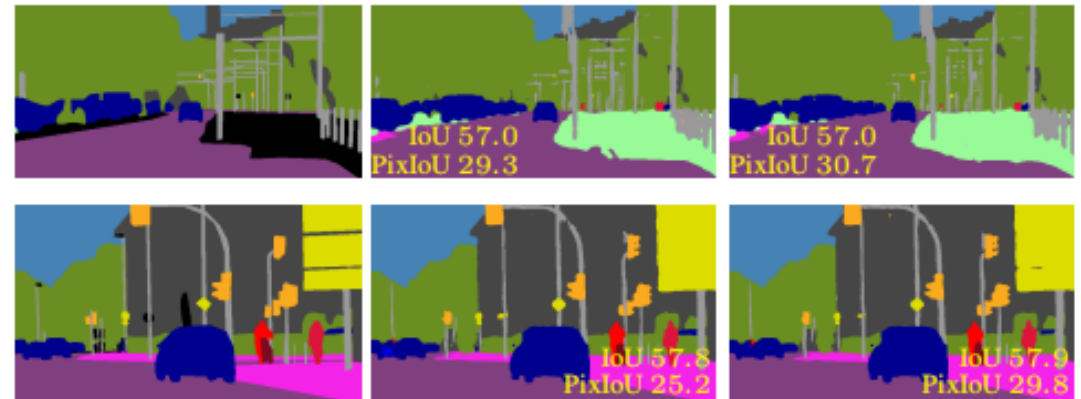
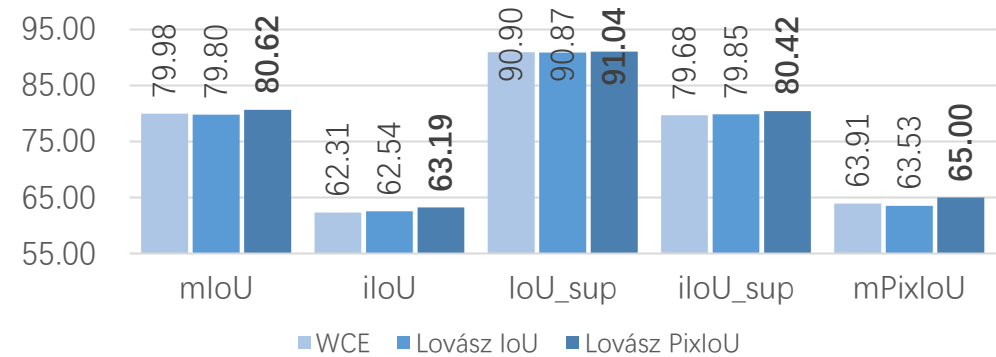


- Semantic segmentation on Pascal VOC



Groundtruth CrossEntropy Lovász IoU Lovász PixIoU

- Semantic segmentation on Cityscapes



PixIoU provides larger gradients than IoU.
Predictions with larger PixIoU provides better qualitative results.

More details please refer to our paper.

ICML | 2021

Thirty-eighth International Conference on
Machine Learning



Thank you!

Learning Generalized Intersection Over Union
for Dense Pixelwise Prediction

Paper: 2497