

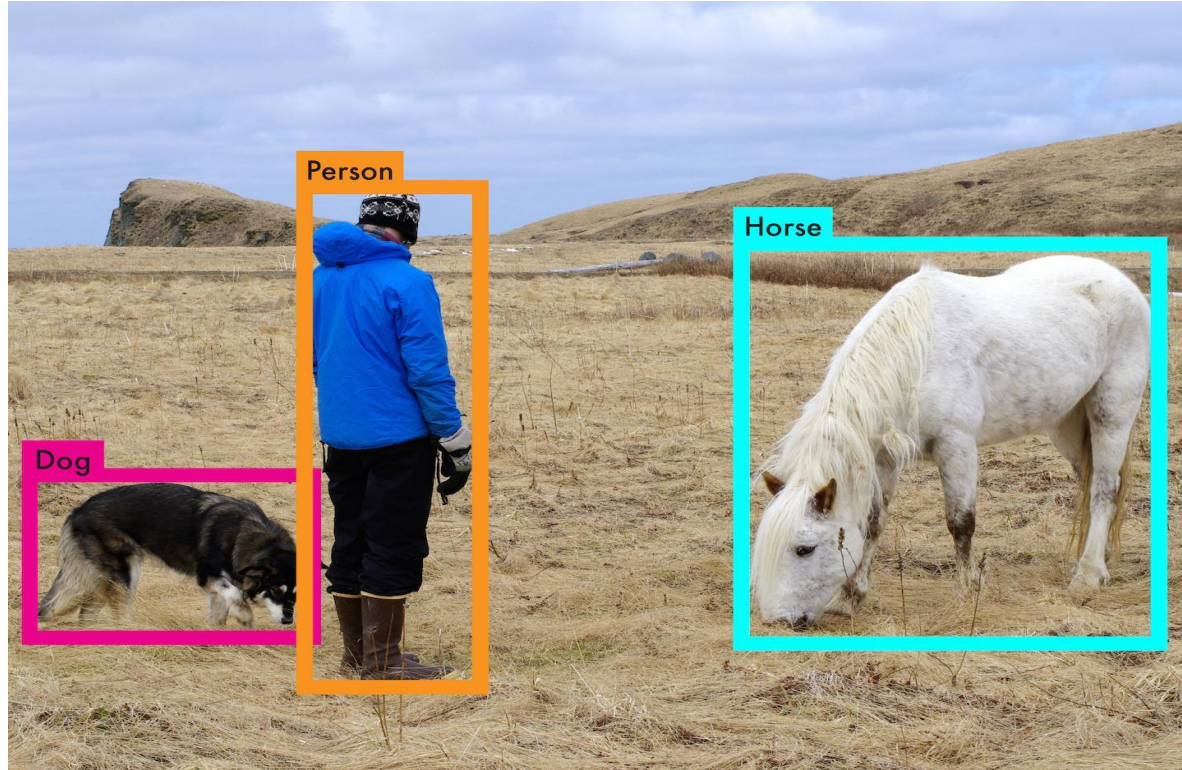
Locating Objects Without Bounding Boxes

Ribera et al. (CVPR 2019)

Utkarsh Ojha, Aron Sarmasi, Rafael A. Rivera-Soto



Object detection



Object detection



Object detection



- Bounding boxes are hard to collect

Object detection



- Bounding boxes are hard to collect
- Why do we even need them?

Object detection



- Bounding boxes are hard to collect
- Why do we even need them?
- Smaller objects seem easier to be detected by **points** than boxes.

Object detection



- Bounding boxes are hard to collect
- Why do we even need them?
- Smaller objects seem easier to be detected by **points** than boxes.
- Points are sometimes enough for weaker localization, and counting instances.

Goal of the paper



Input

Goal of the paper



Input

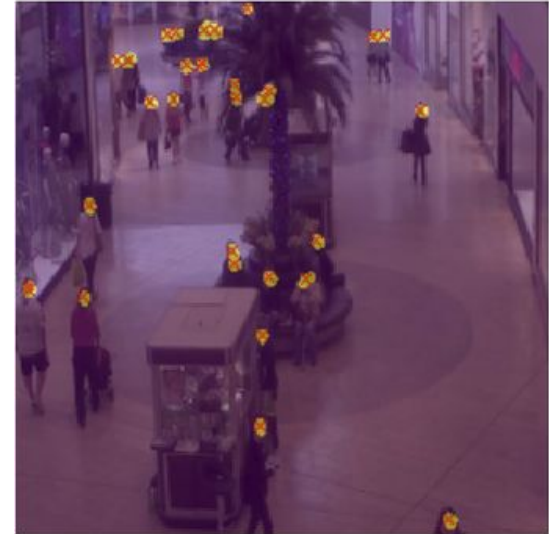
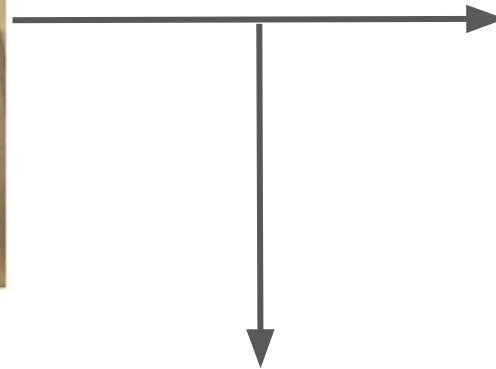


Predicted locations

Goal of the paper



Input



Predicted locations

29
(Object count)

Technical contributions

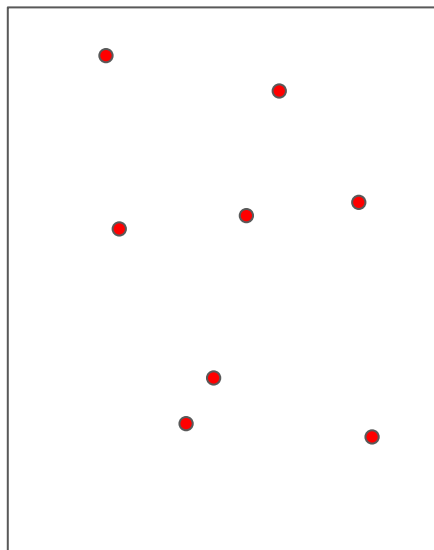
- Modified **Hausdorff loss** to train a *Fully Convolutional Neural Network* (FCN) for object localization.
- Joint estimation of location and number of objects without access to bounding boxes.
- Outperforms state-of-the-art generic object detectors; achieves comparable results for crowd counting.

Localizing points

Object localization

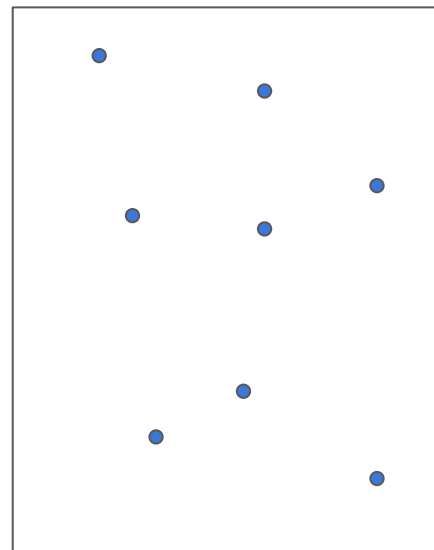


Estimating point locations



Model predictions

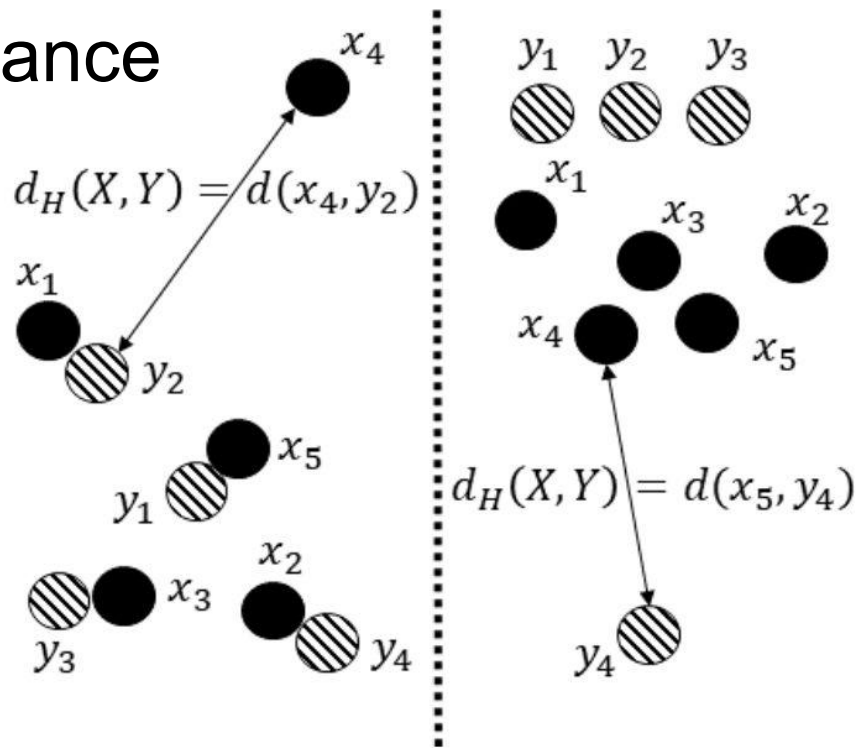
Compute similarity



Ground truth

Visualizing Hausdorff Distance

- Largest smallest distance between points in X and Y
- Intuition: measure of distance of worst outlier
- Not a very good measure for point localization
- Not differentiable w.r.t the FCN output



$$d_H(X, Y) = \max \left\{ \sup_{x \in X} \inf_{y \in Y} d(x, y), \sup_{y \in Y} \inf_{x \in X} d(x, y) \right\}$$

Improvements to Hausdorff Distance

- We need a distance measure that is differentiable w.r.t the FCN output p
- Every output pixel/activation needs to contribute to loss
- High activations near ground truth should have little penalty, and low activations far from the closest ground truth should have little penalty

$$d_{\text{WH}}(p, Y) = \underbrace{\frac{1}{\mathcal{S} + \epsilon} \sum_{x \in \Omega} p_x \min_{y \in Y} d(x, y)}_{\text{High loss to activations far from ground truth}} + \underbrace{\frac{1}{|Y|} \sum_{y \in Y} M_{\alpha} [p_x d(x, y) + (1 - p_x) d_{\text{max}}]}_{\text{Discourages all-zero activations, as term inside generalized mean is maximized by all-zeros}}$$

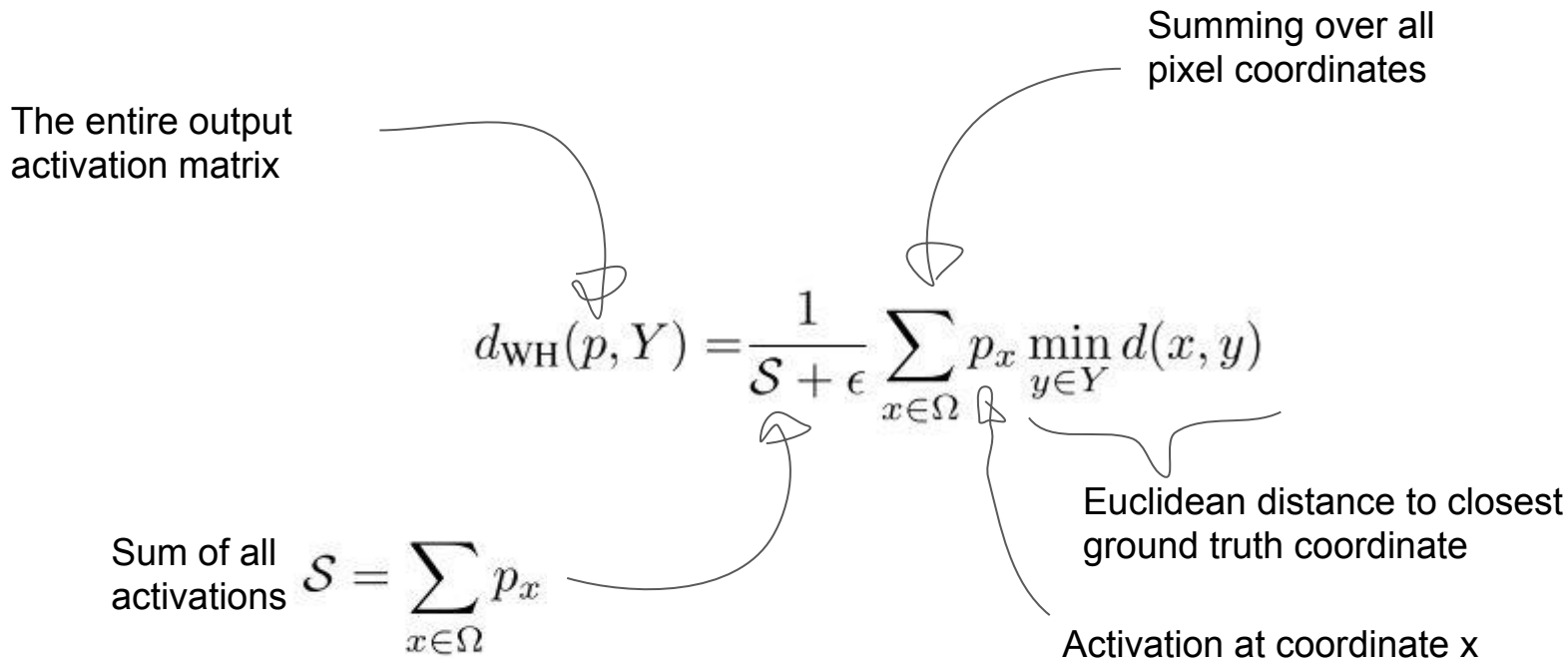
High loss to activations far from ground truth

Discourages all-zero activations, as term inside generalized mean is maximized by all-zeros

Intuition: penalize high activations far from ground truth.

$$d_{\text{WH}}(p, Y) = \frac{1}{\mathcal{S} + \epsilon} \sum_{x \in \Omega} p_x \min_{y \in Y} d(x, y)$$

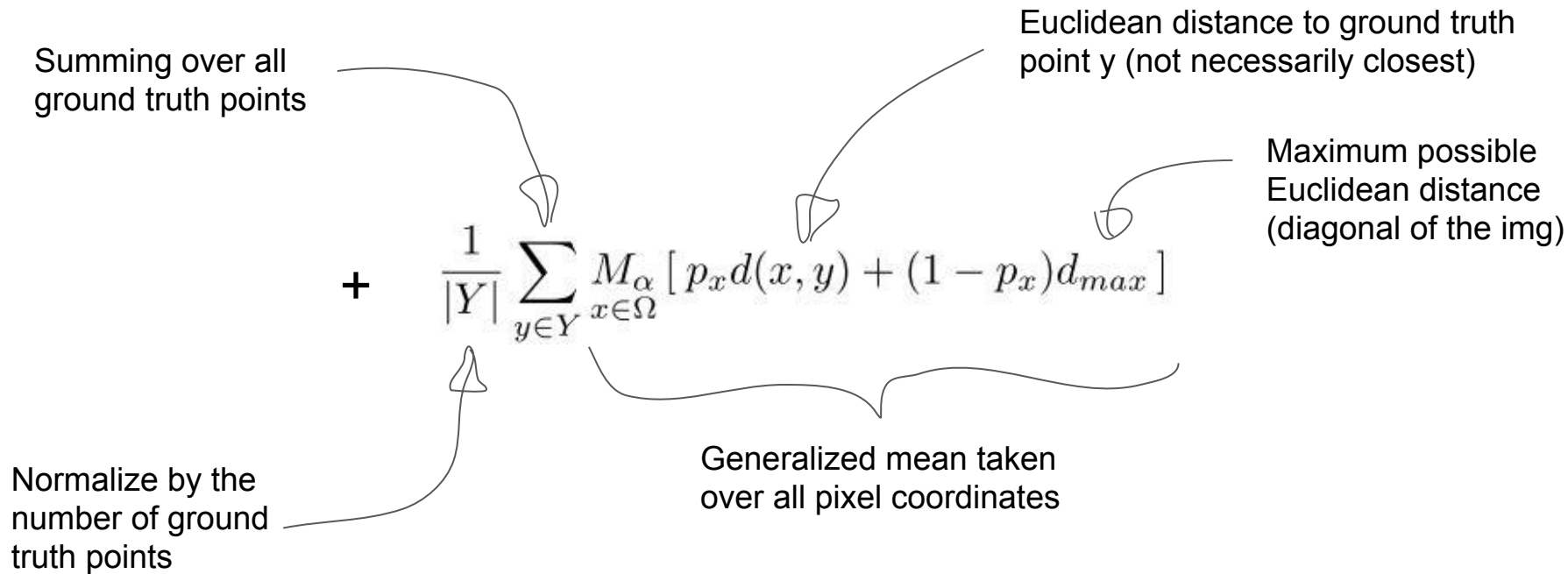
Intuition: penalize high activations far from ground truth.



Intuition: discourage all zero activations

$$+ \frac{1}{|Y|} \sum_{y \in Y} M_{\alpha} [p_x d(x, y) + (1 - p_x) d_{max}]$$

Intuition: discourage all zero activations



$$p_x d(x, y) + (1 - p_x) d_{max}$$

Px = 1, close to gt	True positive	low	low	✓
Px = 1, far from gt	False positive	high	low	✓
Px = 0, close to gt	False negative	low	dmax	✓
Px = 0, far from gt	True negative	low	dmax	?

Generalized Mean
to the Rescue!

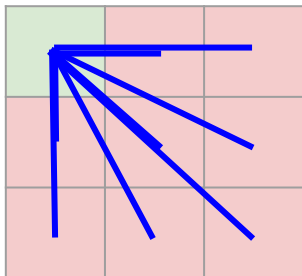
$$M_{\alpha} [f(a)] = \left(\frac{1}{|A|} \sum_{a \in A} f^{\alpha}(a) \right)^{\frac{1}{\alpha}}$$

Num
pixels = n

where $a = p_x d(x, y) + (1 - p_x) d_{max}$

How can the harmonic mean help?

Ground truth (Y)



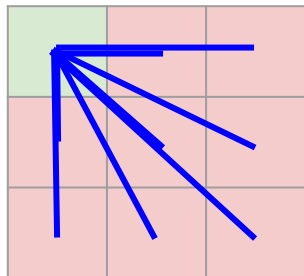
$D_{\max} = 2.83$

Predictions (px)

0	1	0
0	0	0
0	0	0

How can the harmonic mean help?

Ground truth (Y)



Predictions (p_x)

0	1	0
0	0	0
0	0	0

$d(x, y)$ $p_x d(x, y) + (1 - p_x) d_{max}$

1	1
1	2.83
1.41	2.83
2	2.83
2	2.83
2.23	2.83
2.23	2.83
2.83	2.83

$D_{max} = 2.83$

minimum = 1 ($\alpha = -\infty$)

harmonic mean = 1.78 ($\alpha = -1$)

geometric mean = 2.52 ($\alpha = 0$)

arithmetic mean = 2.63 ($\alpha = 1$)

maximum = 2.83 ($\alpha = +\infty$)

$$p_x d(x, y) + (1 - p_x) d_{max}$$

$P_x = 1$, close to gt	True positive	low	low	✓
$P_x = 1$, far from gt	False positive	high	low	✓
$P_x = 0$, close to gt	False negative	low	d_{max}	✓
$P_x = 0$, far from gt	True negative	low	d_{max}	?

Harmonic mean **greatly weighted** towards the lower values of $p_x d(x, y) + (1 - p_x) d_{max}$

most penalty

$$\frac{1}{|Y|} \sum_{y \in Y} M_{\alpha} [p_x d(x, y) + (1 - p_x) d_{max}]$$

		Penalty Amount
Px = 1, close to gt	True positive	least
Px = 1, far from gt	False positive	most
Px = 0, close to gt	False negative	most
Px = 0, far from gt	True negative	middle

least penalty

penalty

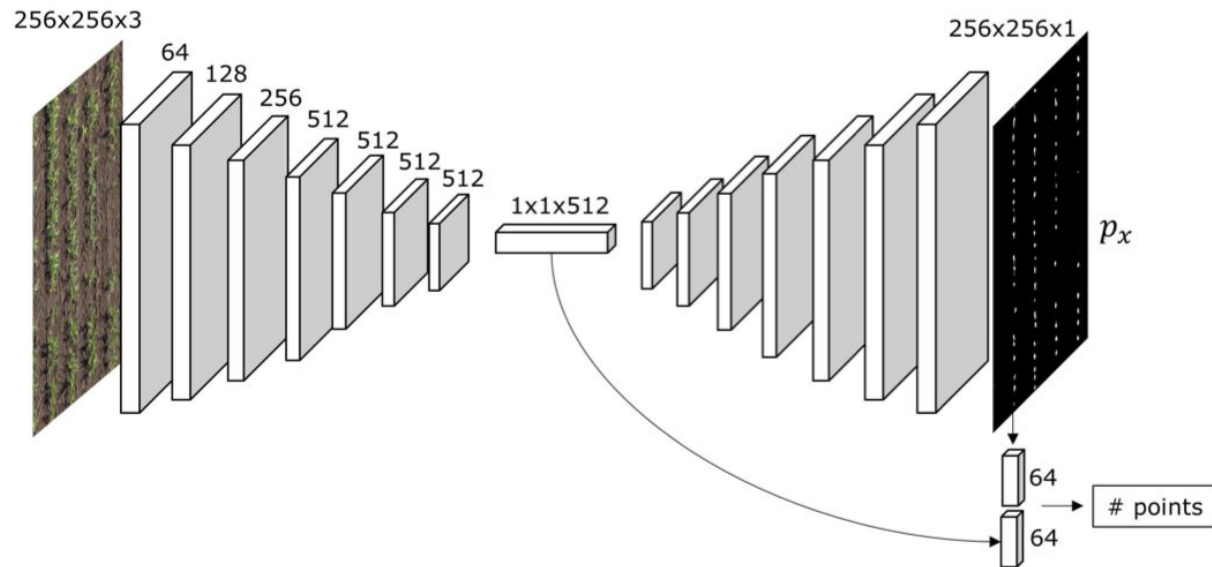
$$\frac{1}{S + \epsilon} \sum_{x \in \Omega} p_x \min_{y \in Y} d(x, y)$$

To tie it all together:

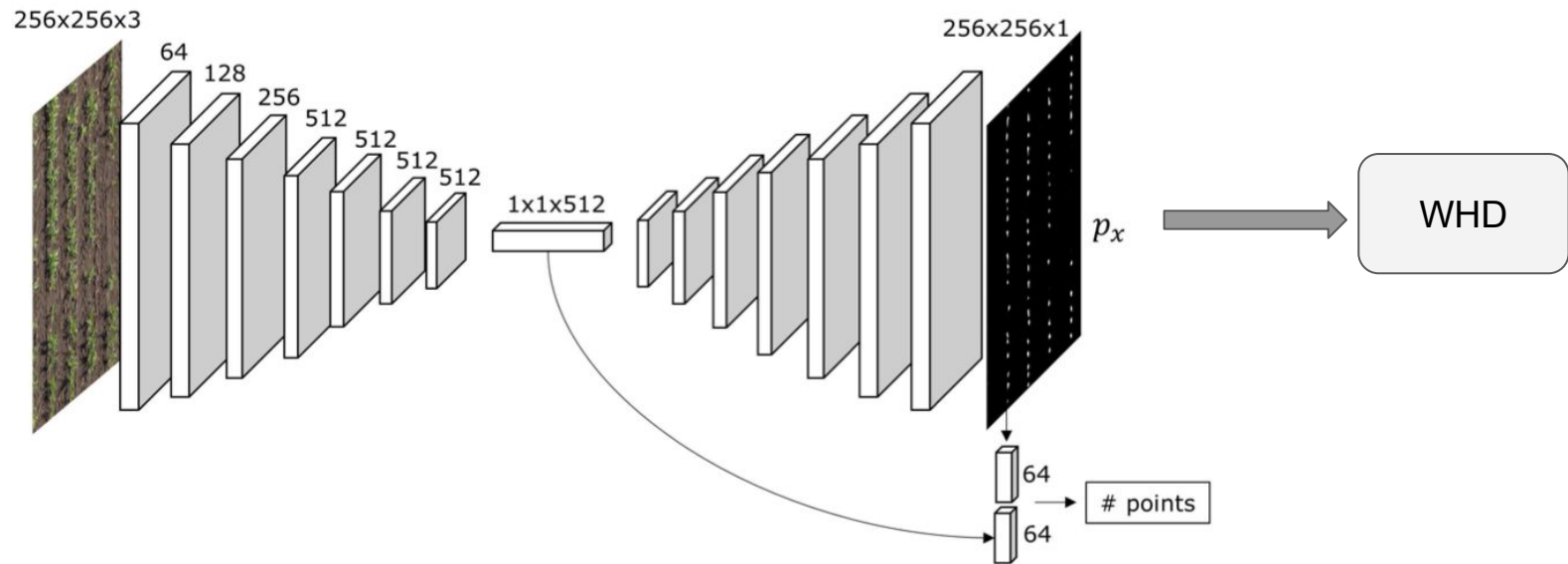
$$d_{\text{WH}}(p, Y) = \frac{1}{S + \epsilon} \sum_{x \in \Omega} p_x \min_{y \in Y} d(x, y) + \frac{1}{|Y|} \sum_{y \in Y} M_\alpha [p_x d(x, y) + (1 - p_x) d_{\text{max}}]$$

- Fully differentiable w.r.t output of FCN
- Converges to maximize true positives true negatives, and minimize all else

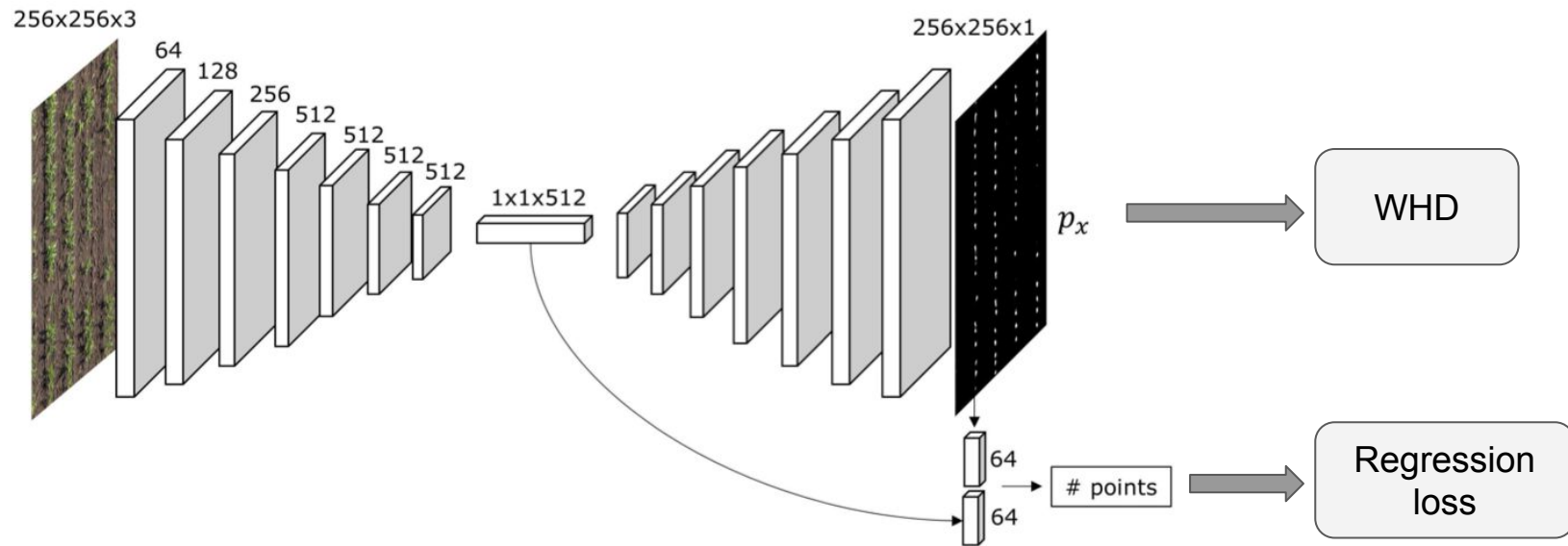
Model Architecture



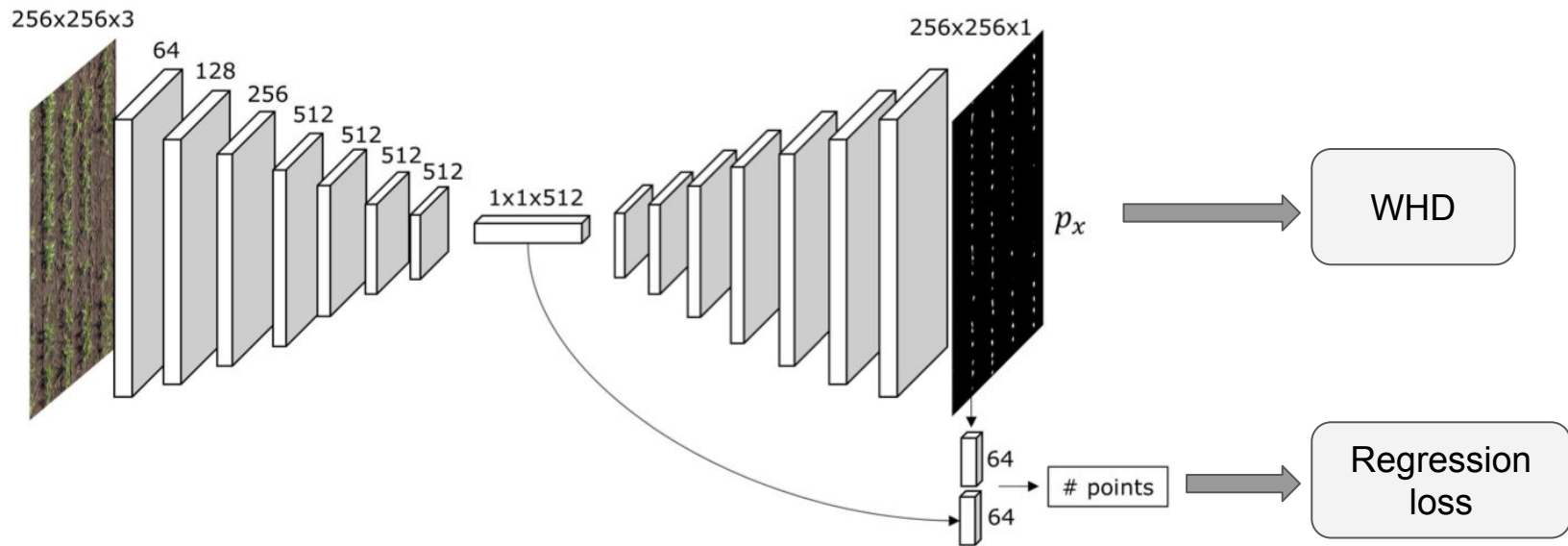
Model Architecture



Model Architecture

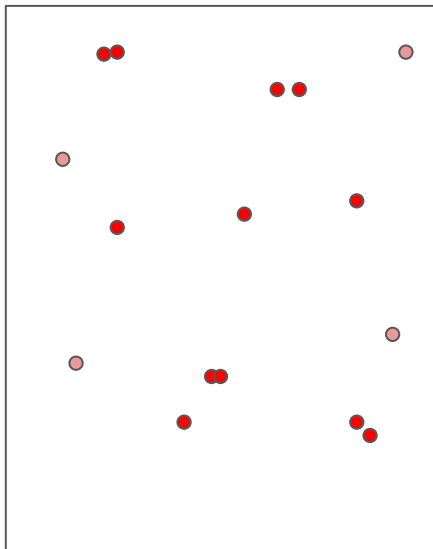


Model Architecture



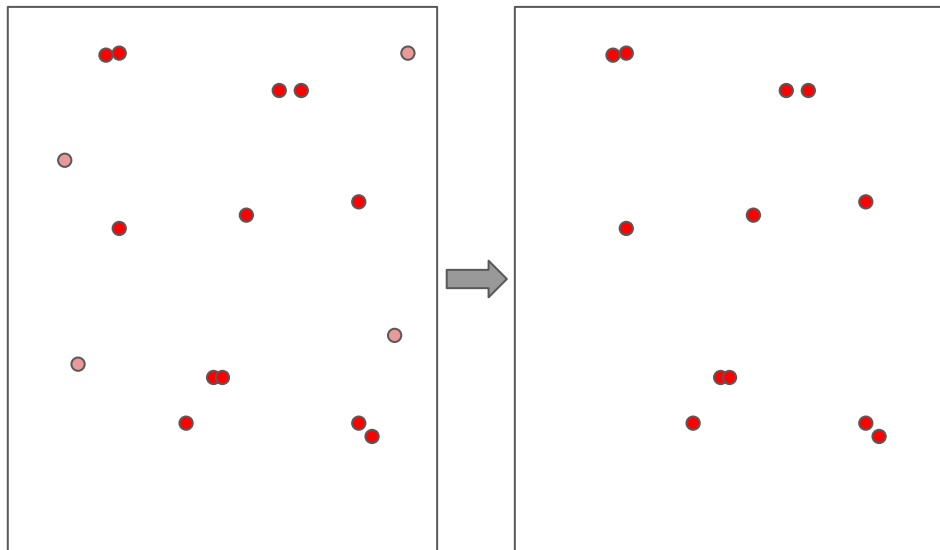
$$L_{reg}(x) = \begin{cases} 0.5x^2 & \text{for } |x| < 1 \\ |x| - 0.5 & \text{for } |x| \geq 1 \end{cases}$$

Computing model's predictions



Predicted probability
map

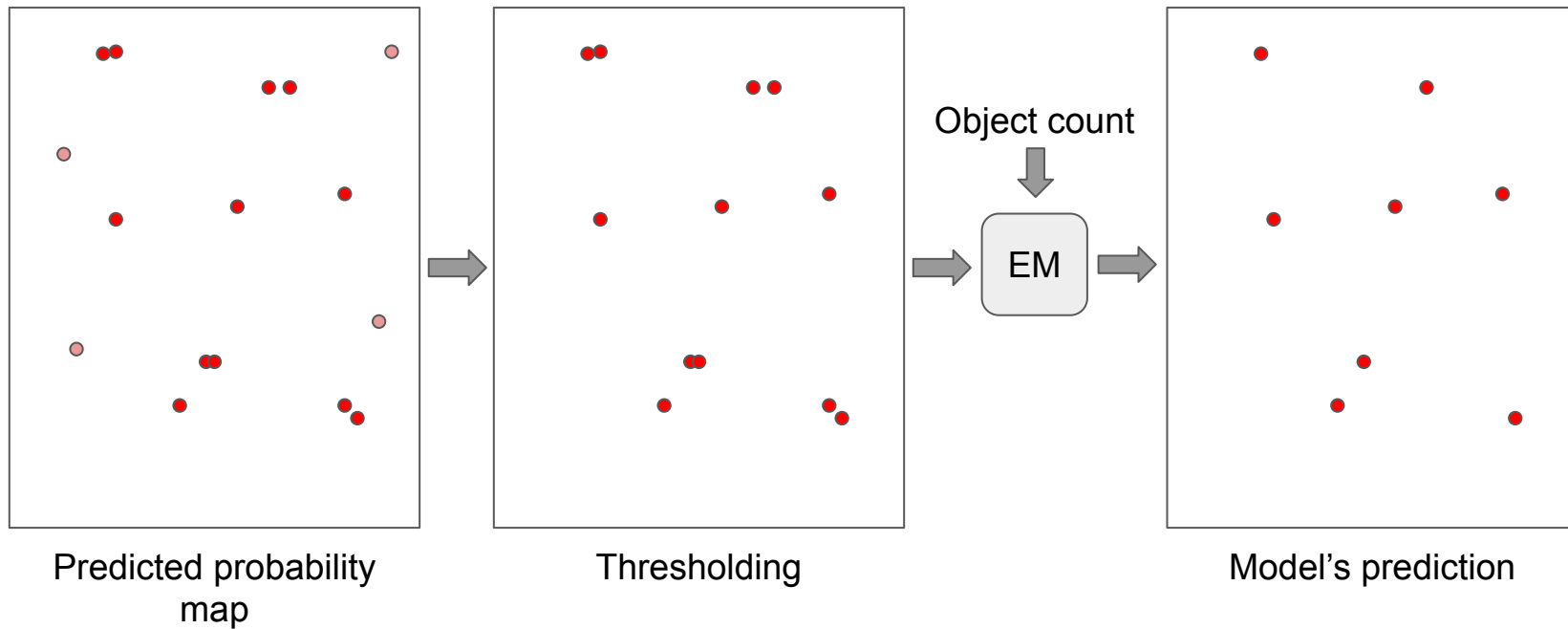
Computing model's predictions



Predicted probability
map

Thresholding

Computing model's predictions



Overall training objective

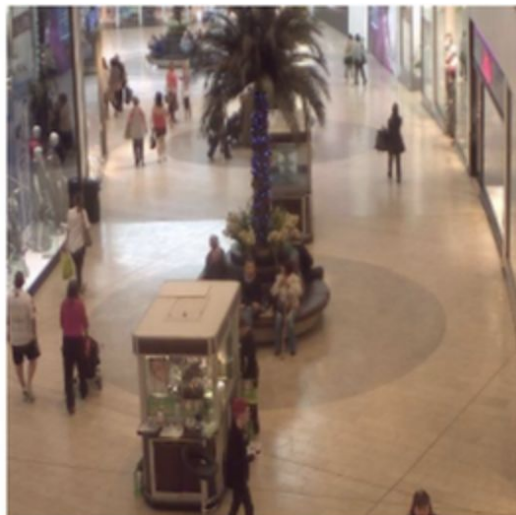
$$L(p, Y) = d_{wh}(p, Y) + L_{reg}(C - \hat{C}(p))$$

Overall training objective

$$L(p, Y) = d_{wh}(p, Y) + L_{reg}(C - \hat{C}(p))$$



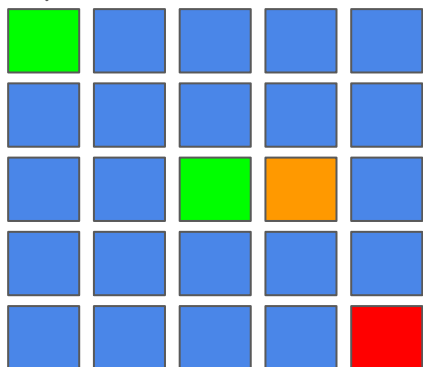
Datasets



- 80/10/10 train, validation and test split for each dataset
- Images resized to 256x256
- Augmented with random horizontal flips

Metrics

False Negative



$r = 1$



Ground Truth



True Positive



False Positive

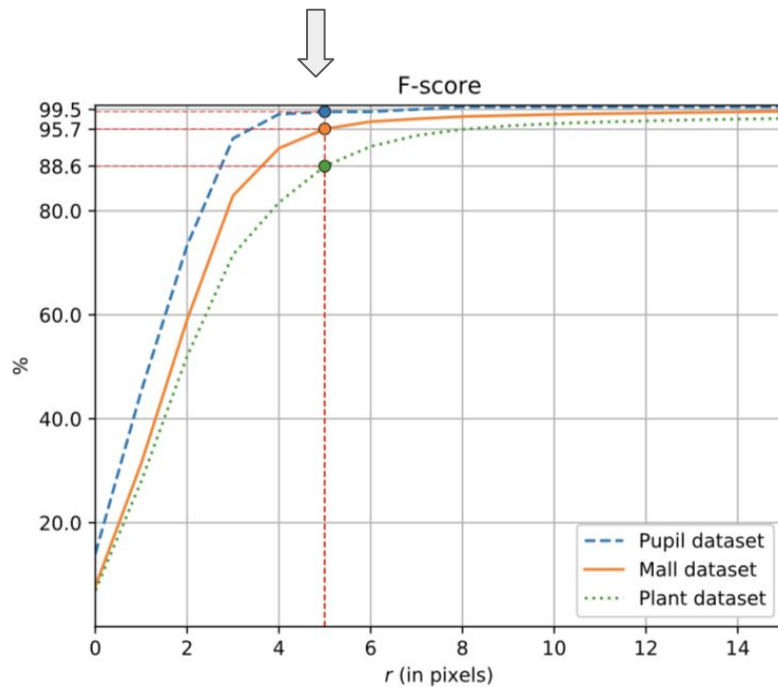
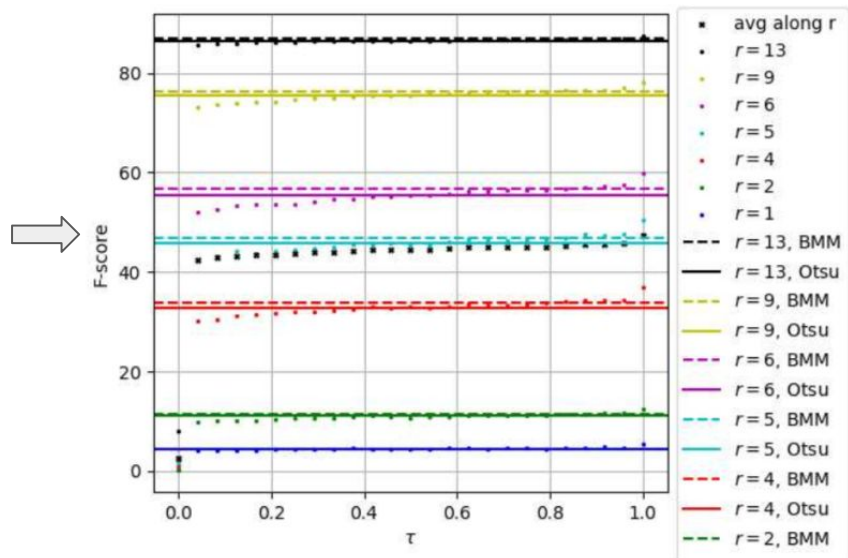
$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$FScore = 2 \frac{Precision * Recall}{Precision + Recall}$$

- Precision and Recall can be 100% even if the model estimates 2 object locations per ground truth point.
- MAE, RMSE and MAPE are reported to counteract this.

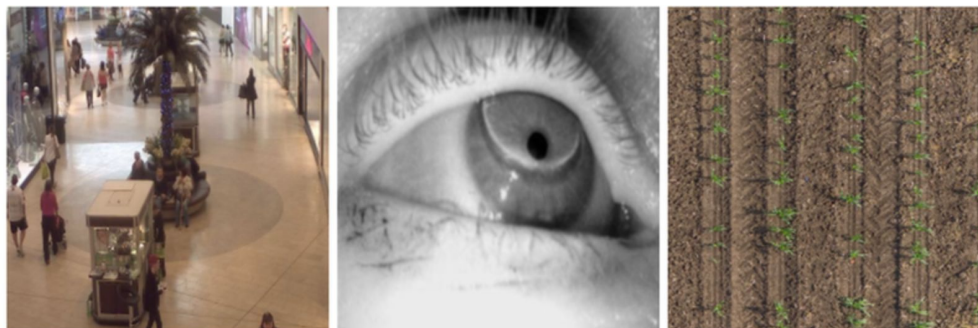
Experimental Evaluation



- The bigger “ r ” is, the easier the problem becomes

Experimental Evaluation

- Comparison against Faster-RCNN with bounding boxes of 20x20 centered at true point
- Model also evaluated on ShanghaiTech Part B achieving MAE of 19.9



Metric	Faster-RCNN	Ours
Precision	81.1%	95.2 %
Recall	76.7%	96.2 %
F-score	78.8 %	95.7 %
AHD	7.6 px	4.5 px
MAE	4.7	1.4
RMSE	5.6	1.8
MAPE	14.8%	4.4 %

Method	Precision	Recall	AHD
Swirski [53]	77 %	77 %	-
ExCuSe [13]	77 %	77 %	-
Faster-RCNN	99.5 %	99.5 %	2.7 px
Ours	99.5 %	99.5 %	2.5 px

Metric	Faster-RCNN	Ours
Precision	86.6 %	88.1 %
Recall	78.3 %	89.2 %
F-score	82.2 %	88.6 %
AHD	9.0 px	7.1 px
MAE	9.4	1.9
RMSE	13.4	2.7
MAPE	17.7 %	4.2 %

Strengths

- Dramatically reduces amount of work to annotate a dataset
- No major architectural constraints
- Tested on multiple datasets
- Re-formulation of the object localization problem as the minimization of the distances between a set of points

Weaknesses

- No indication of the size, orientation, occlusion, etc of the object predicted, only center position and instance count
- No comparison between the weighted Hausdorff Distance and other pixel-wise losses such as L2 or MSE.
- Each dataset contained only one type of object, does the method work when trying to detect a wide variety of objects?
- How does it perform with videos, where the objects can exhibit a wide variety of behaviors?
- Notation for generalized mean is misleading.
- Motivation of WHD is done assuming $\alpha = -\infty$. But in practice they use $\alpha = -1$.

Thank You!