

# Training Neural Networks on Non-Differentiable Losses

**Yash Patel**

Supervisor: Professor Jiří Matas

Visual Recognition Group, Czech Technical University in Prague

# Supervised Deep Learning: Training

Three main components of Supervised Deep Learning:

Training data



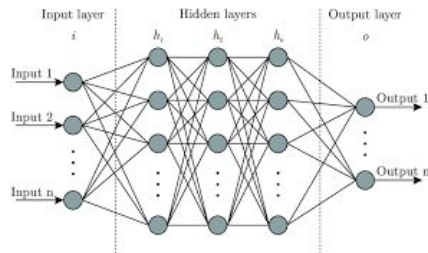
# Supervised Deep Learning: Training

Three main components of Supervised Deep Learning:

## Training data



## Model



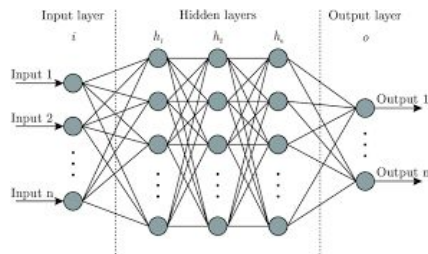
# Supervised Deep Learning: Training

Three main components of Supervised Deep Learning:

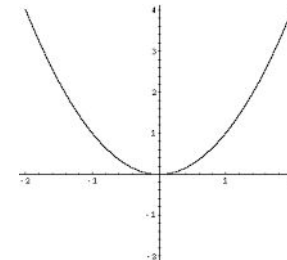
## Training data



## Model



## Loss



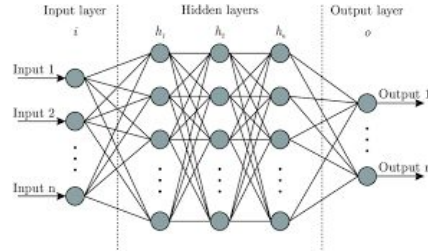
# Supervised Deep Learning: Training

Three main components of Supervised Deep Learning:

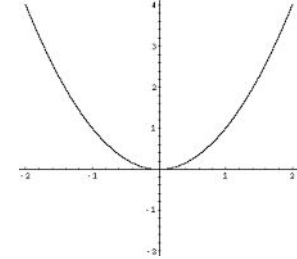
Training data



Model



Loss



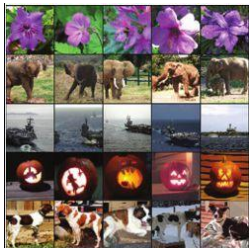
Using backpropagation, weights of the model are updated.

Note: The loss function needs to be differentiable for the use of chain-rule to obtain gradients with respect to the weights.

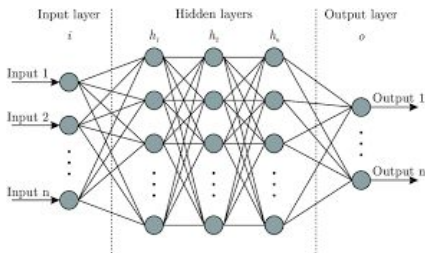
# Supervised Deep Learning: Evaluation

The evaluation metric is task dependent, not algorithm dependent.

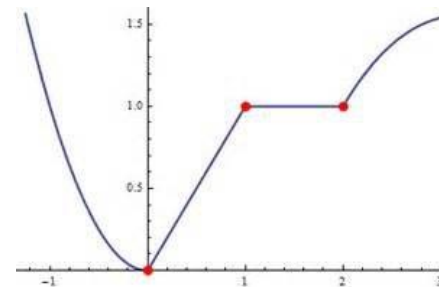
## Testing data



## Trained Model



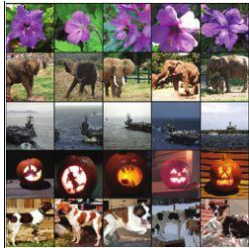
## Evaluation Metric



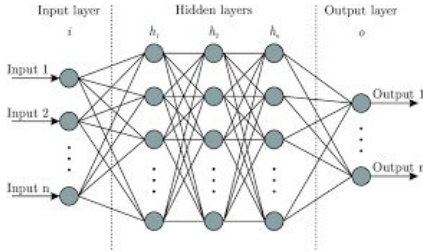
# Supervised Deep Learning: Evaluation

The evaluation metric is task dependent, not algorithm dependent.

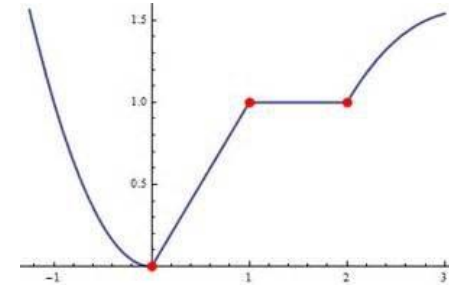
Testing data



Trained Model



Evaluation Metric

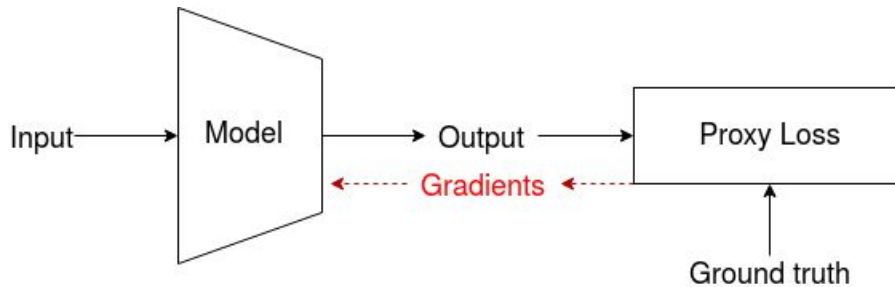


Two cases:

1. The evaluation metric is **differentiable**, therefore, can be used as a loss function.
2. The evaluation metric is **non-differentiable**, therefore, can not be used as a loss function.

# Proxy Losses

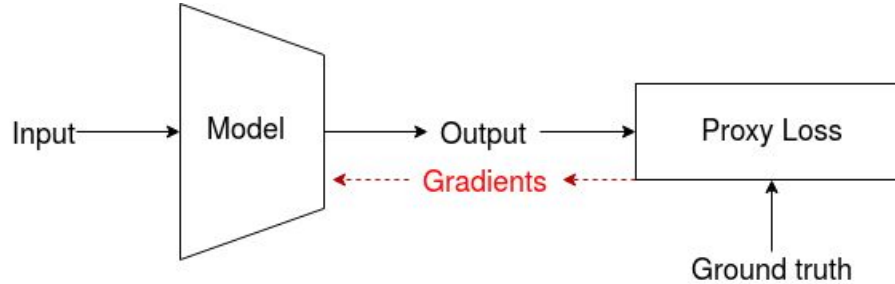
Proxy losses may not always align with the evaluation metric.





# Proxy Losses

Proxy losses may not always align with the evaluation metric.

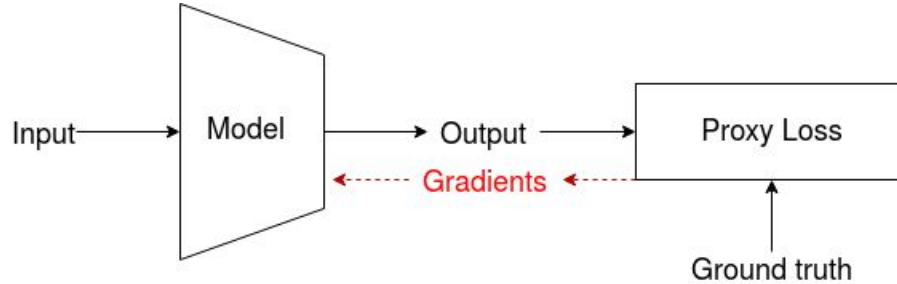


Evaluation metric in green, proxy loss functions in blue.

- Image compression (human perception of similarity): structural similarity index, peak signal to noise ratio, mean squared error, etc.

# Proxy Losses

Proxy losses may not always align with the evaluation metric.

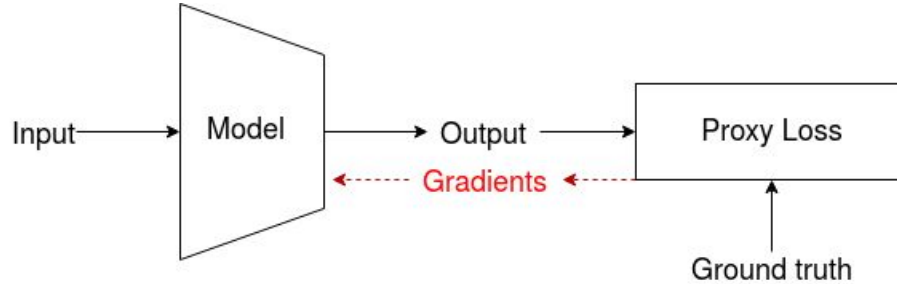


Evaluation metric in green, proxy loss functions in blue.

- Image compression (human perception of similarity): structural similarity index, peak signal to noise ratio, mean squared error, etc.
- Object detection (intersection over union): smooth-L1 distance, L2 distance.

# Proxy Losses

Proxy losses may not always align with the evaluation metric.

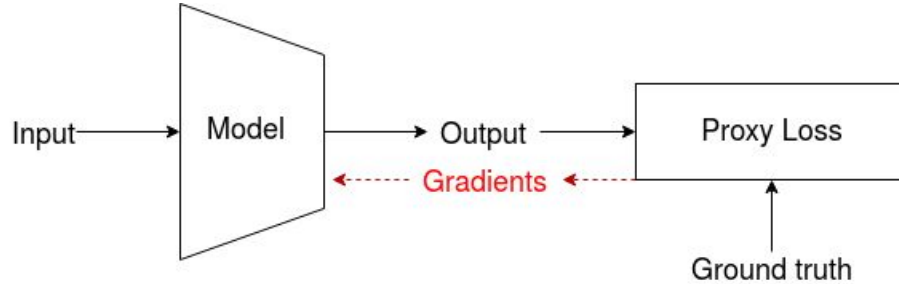


Evaluation metric in green, proxy loss functions in blue.

- Image compression (human perception of similarity): structural similarity index, peak signal to noise ratio, mean squared error, etc.
- Object detection (intersection over union): smooth-L1 distance, L2 distance.
- Scene text recognition (edit distance): per-character cross entropy, connectionist temporal classification.

# Proxy Losses

Proxy losses may not always align with the evaluation metric.



Evaluation metric in green, proxy loss functions in blue.

- Image compression (human perception of similarity): structural similarity index, peak signal to noise ratio, mean squared error, etc.
- Object detection (intersection over union): smooth-L1 distance, L2 distance.
- Scene text recognition (edit distance): per-character cross entropy, connectionist temporal classification.
- Image retrieval (mean average precision, recall@k): contrastive loss, triplet loss, proxy NCA, etc.

And many more...



# Saliency Driven Perceptual Image Compression

Yash Patel <sup>1\*</sup> Srikar Appalaraju <sup>2</sup> R. Manmatha <sup>2</sup>

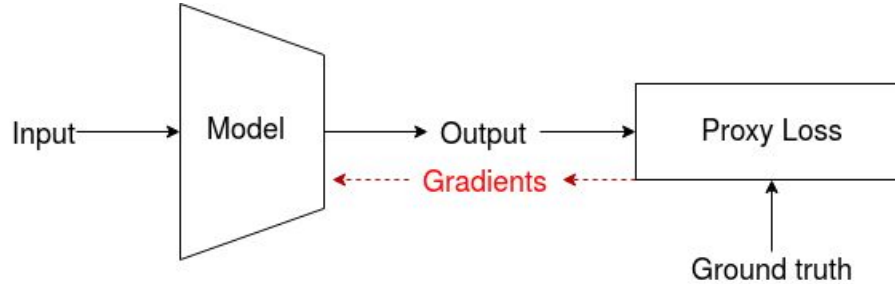
<sup>1</sup> Visual Recognition Group, Czech Technical University in Prague  
<sup>2</sup> Amazon Web Services, Palo Alto

WACV 2021

\* This research was conducted during Yash Patel's internship at AWS.

# Proxy Losses

Proxy losses may not always align with the evaluation metric.



Evaluation metric in green, proxy loss functions in blue.

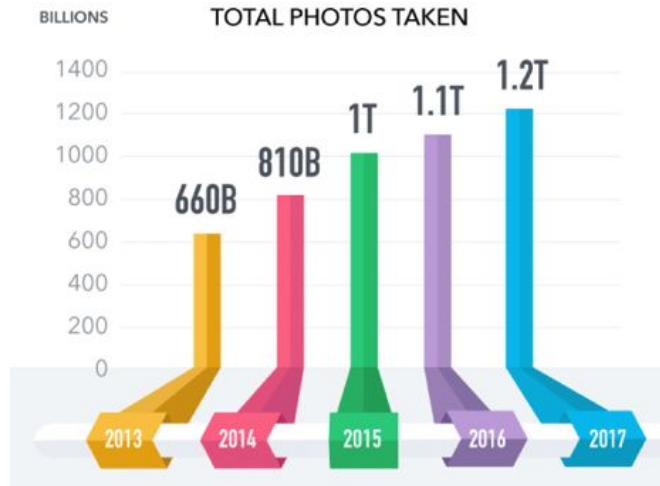
- Image compression (human perception of similarity): structural similarity index, peak signal to noise ratio, mean squared error, etc.
- 
- **Note: Can you even express this as a mathematical function?**
- Image retrieval (mean average precision, recall@k): contrastive loss, triplet loss, proxy NCA, etc.

And many more...

# Motivation



## Lower Storage Requirements



1. Snapchat users share 527,760/min photos.
2. Instagram users post 46,740/min photos.

[1] How Much Data Do We Create Every Day? The Mind-Blowing Stats Everyone Should Read, Bernard Marr, Forbes 2018.

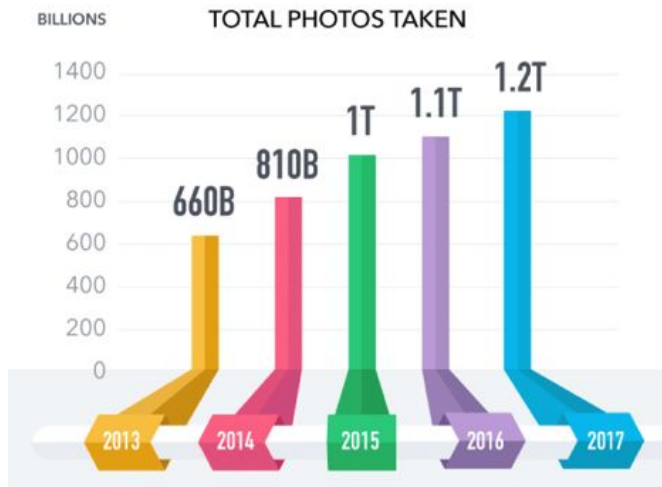
[2] Here's How Many Digital Photos Will Be Taken in 2017, Tech Today, 2016.

[3] Towards Image Understanding from Deep Compression without Decoding, Torfason et al. ICLR 2018.

# Motivation

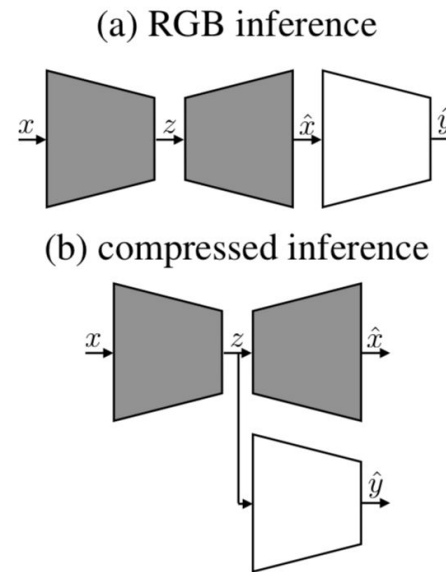


## Lower Storage Requirements



1. Snapchat users share 527,760/min photos.
2. Instagram users post 46,740/min photos.

## Faster Inference for subsequent tasks



Learning based compression methods lead to faster inference for subsequent tasks such as classification, detection and semantic segmentation.

[1] How Much Data Do We Create Every Day? The Mind-Blowing Stats Everyone Should Read, Bernard Marr, Forbes 2018.

[2] Here's How Many Digital Photos Will Be Taken in 2017, Tech Today, 2016.

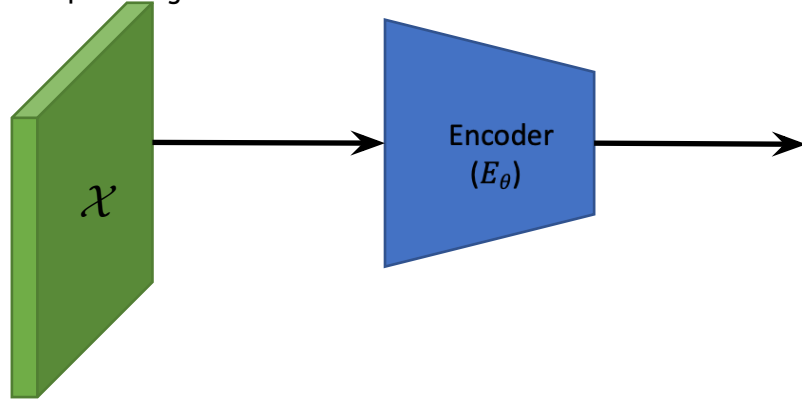
[3] Towards Image Understanding from Deep Compression without Decoding, Torfason et al. ICLR 2018.



# Learning Image Compression



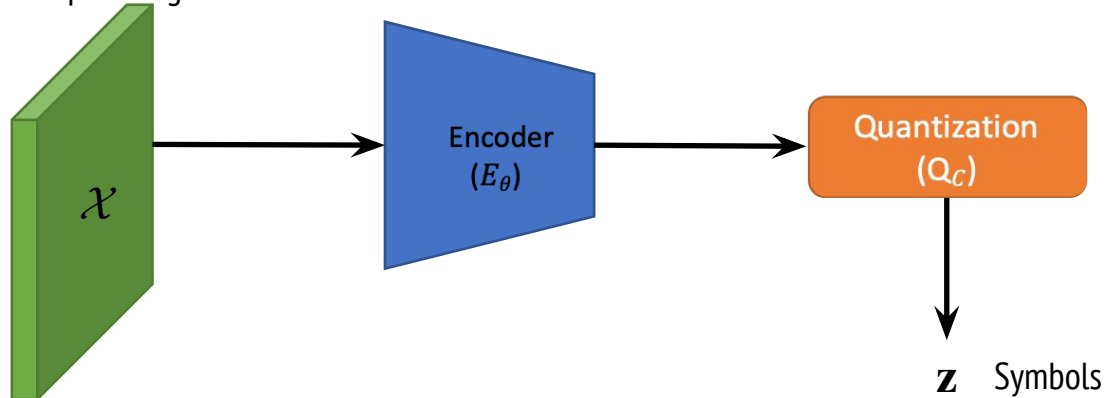
Input Image



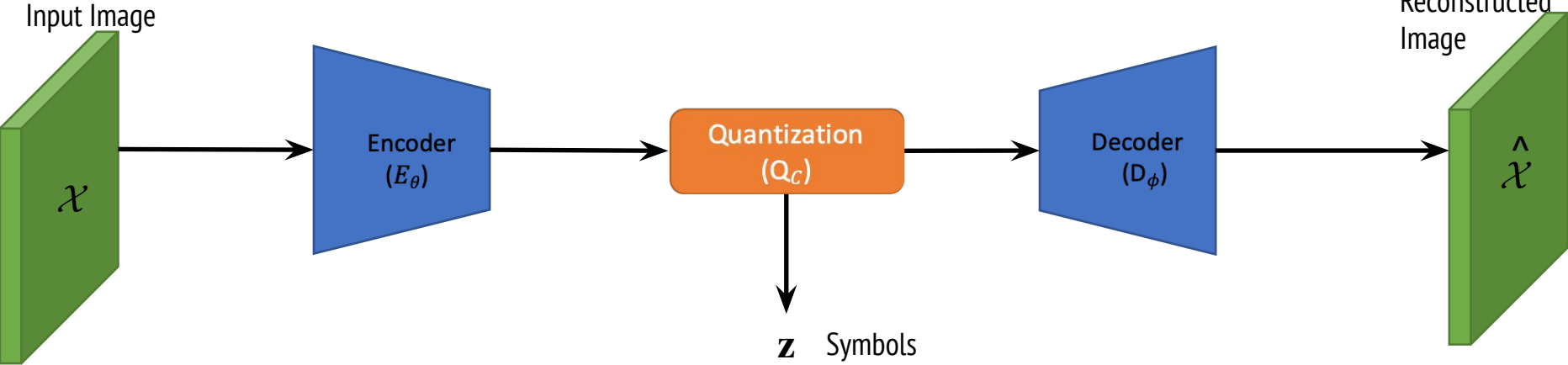
# Learning Image Compression



Input Image



# Learning Image Compression



# Evaluation Metrics / Proxy Losses



## 1. Multi-scale Structural Similarity (MS-SSIM)

For a sliding window on original and reconstructed images.

$$SSIM(x, \hat{x}) = \frac{(2\mu_x\mu_{\hat{x}} + c_1)(2\sigma_{x\hat{x}} + c_2)}{(\mu_x^2 + \mu_{\hat{x}}^2 + c_1)(\sigma_x^2 + \sigma_{\hat{x}}^2 + c_2)}$$

Aggregating SSIM at multiple scales is MS-SSIM

## 2. Peak Signal to Noise Ratio (PSNR)

$$MSE(x, \hat{x}) = \frac{1}{mn} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} [x(i, j) - \hat{x}(i, j)]^2$$

$$PSNR = 10 \log_{10} \left( \frac{MAX_x^2}{MSE(x, \hat{x})} \right)$$

# Evaluation Metric Problems



Original Image



MS-SSIM/SSIM cannot distinguish slightly blurred and not-blurred versions.

Higher MS-SSIM



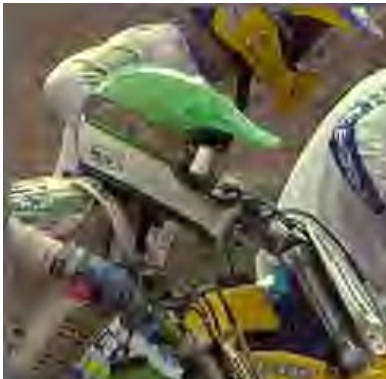
Mentzer et al. CVPR'18



Ballé et al. ICLR'17

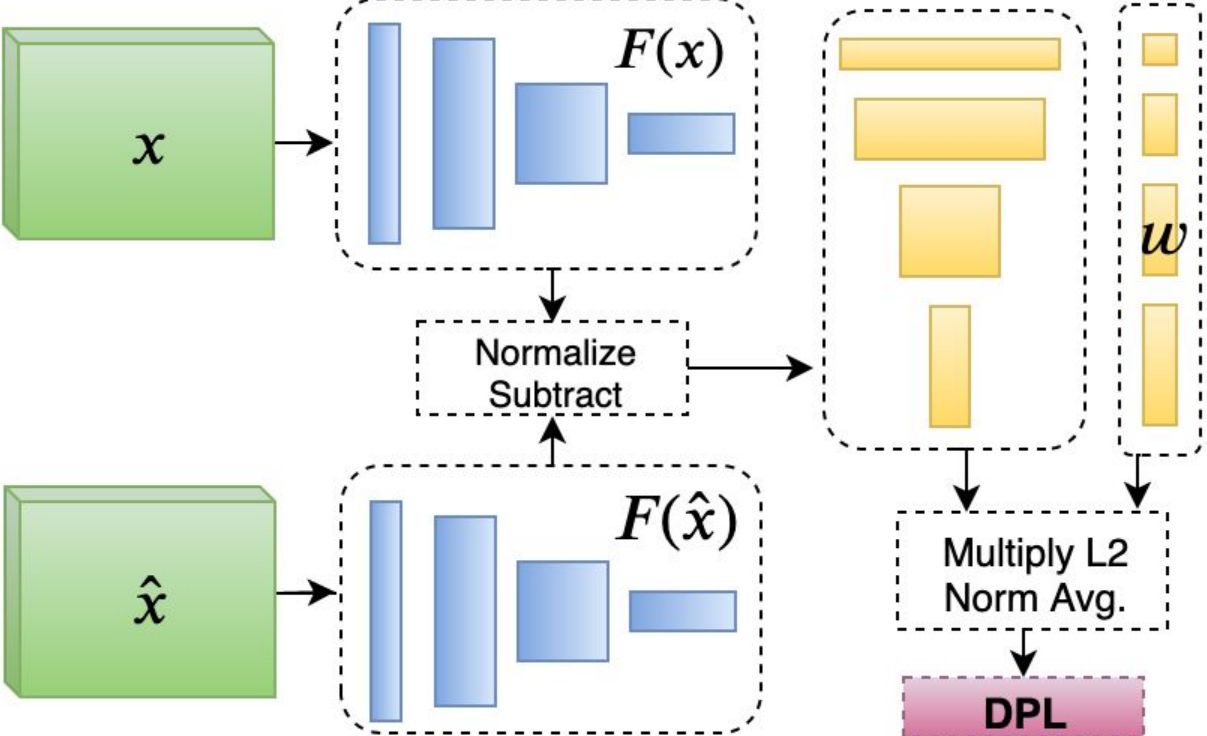


BPG



JPEG-2000

# LPIPS-Comp



# Human Evaluations



Full Images



Synchronized magnifying glass

Note: We first test the evaluators on internally annotated golden set.

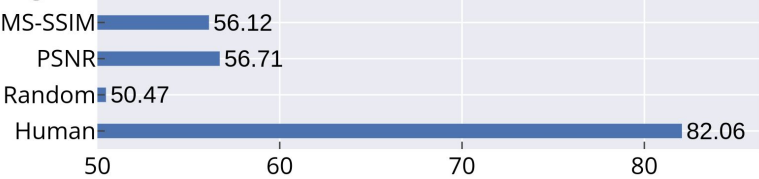
Image A is more similar

Image B is more similar

# Evaluation Metrics



Hand-crafted metrics {



Patel et al., Saliency Driven Perceptual Image Compression, WACV 2021.

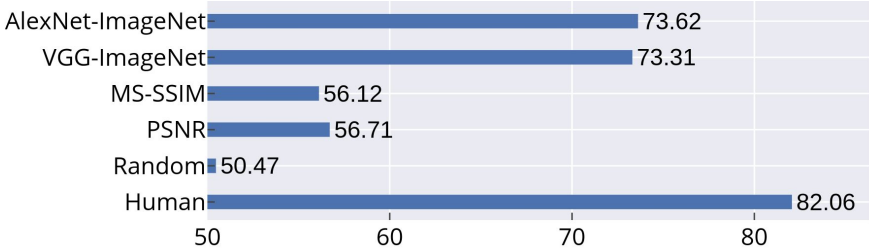
2AFC score %



# Evaluation Metrics



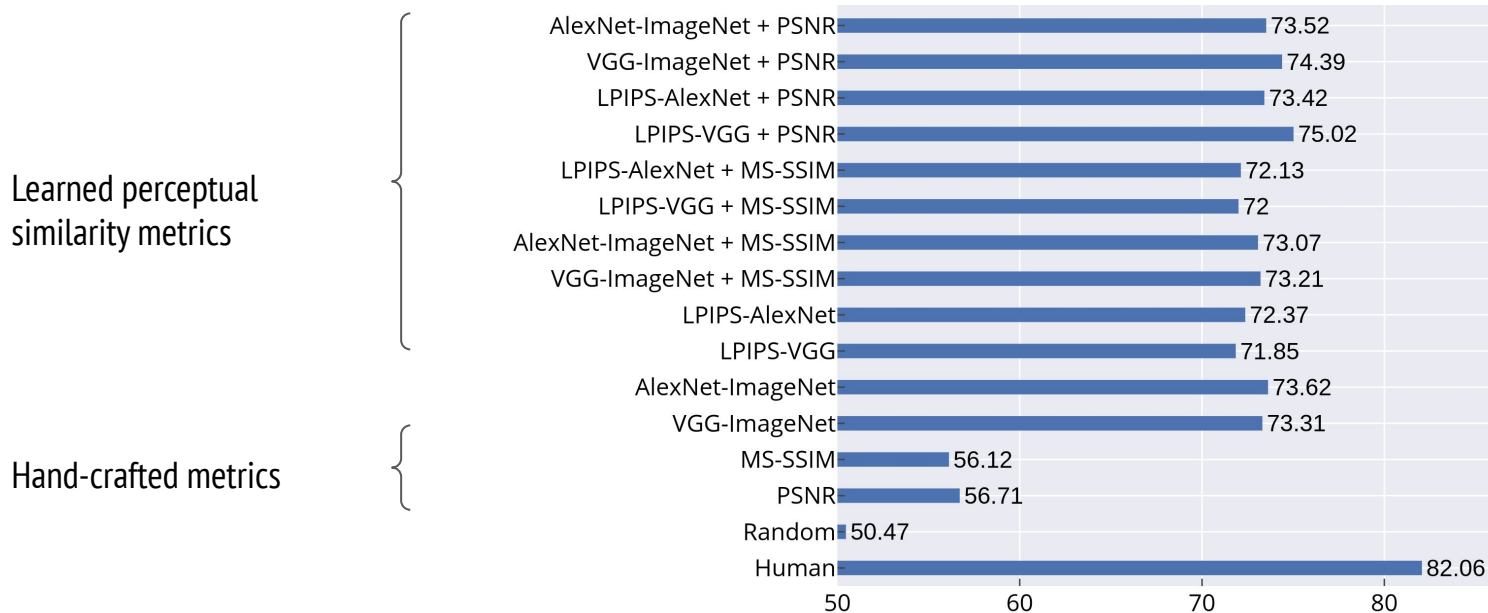
Hand-crafted metrics {



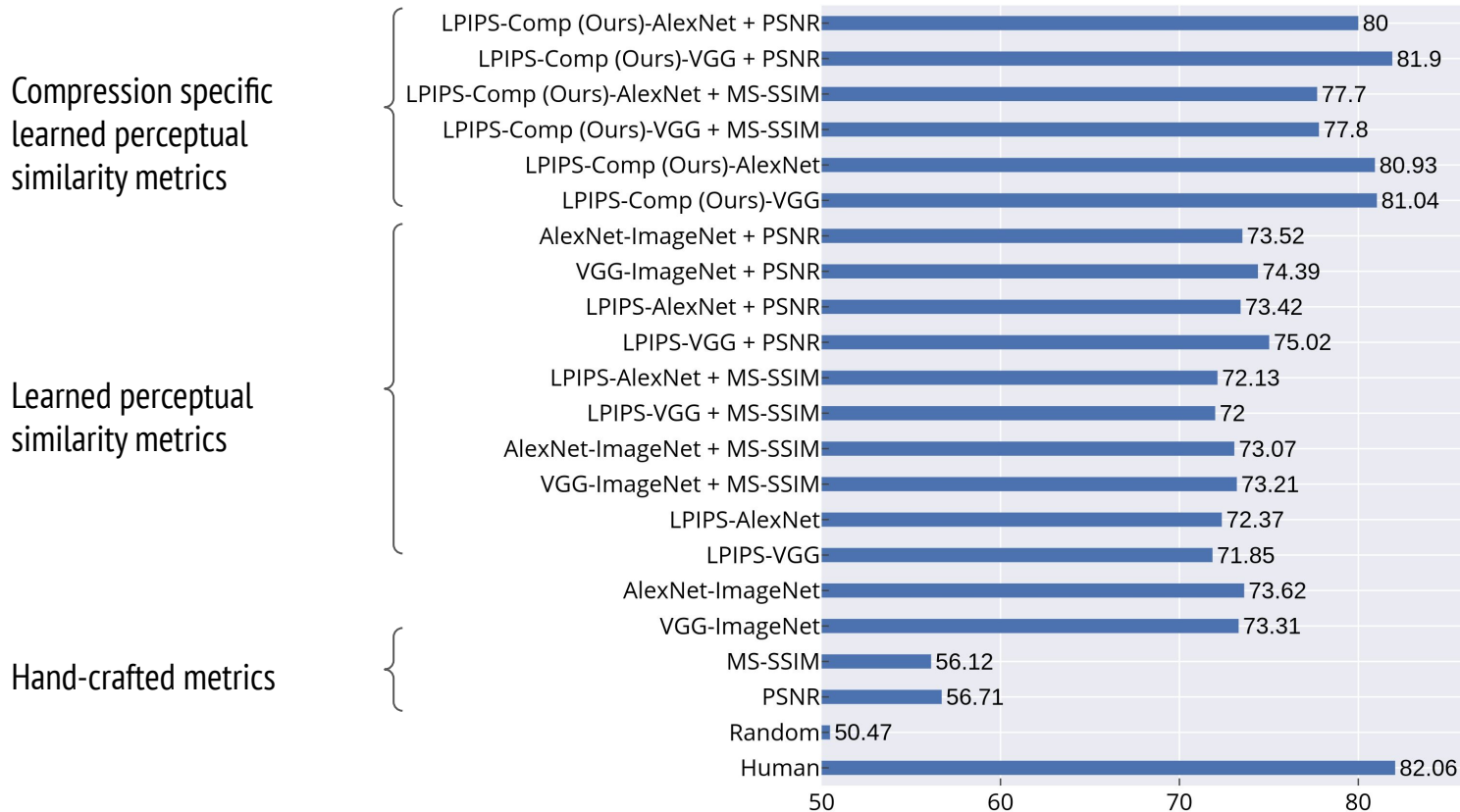
Patel et al., Saliency Driven Perceptual Image Compression, WACV 2021.

2AFC score %

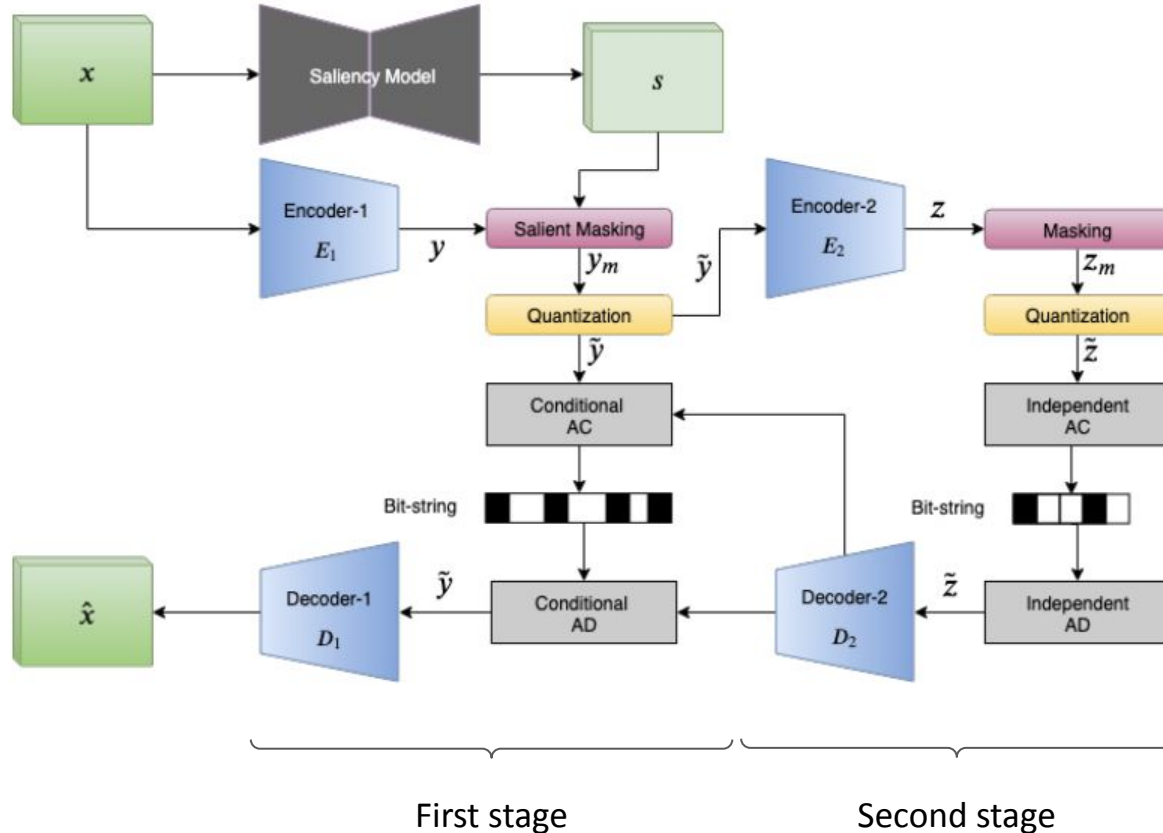
# Evaluation Metrics



# Evaluation Metrics



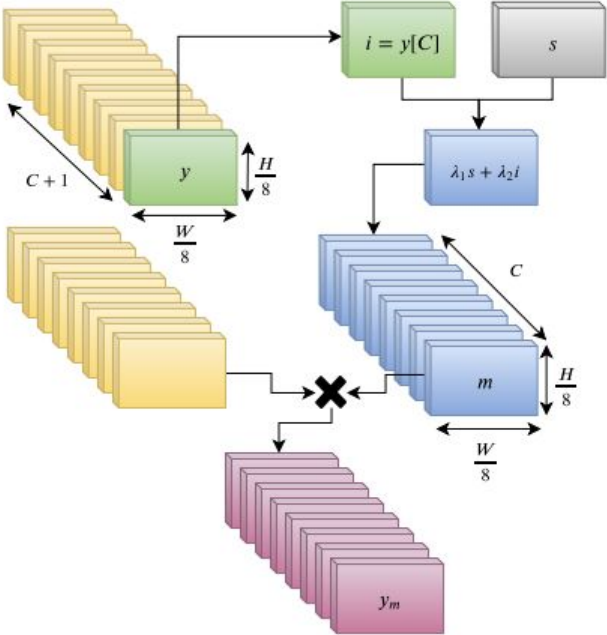
# Hierarchical Auto-Regressive Model



# Saliency Matters



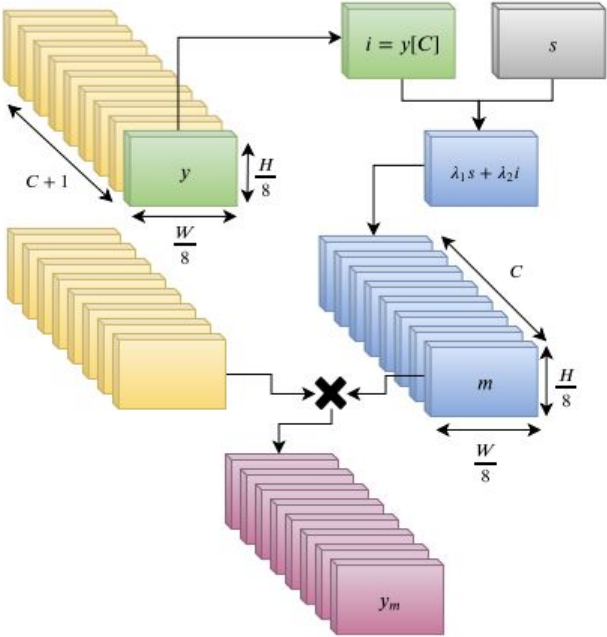
## Saliency Masking



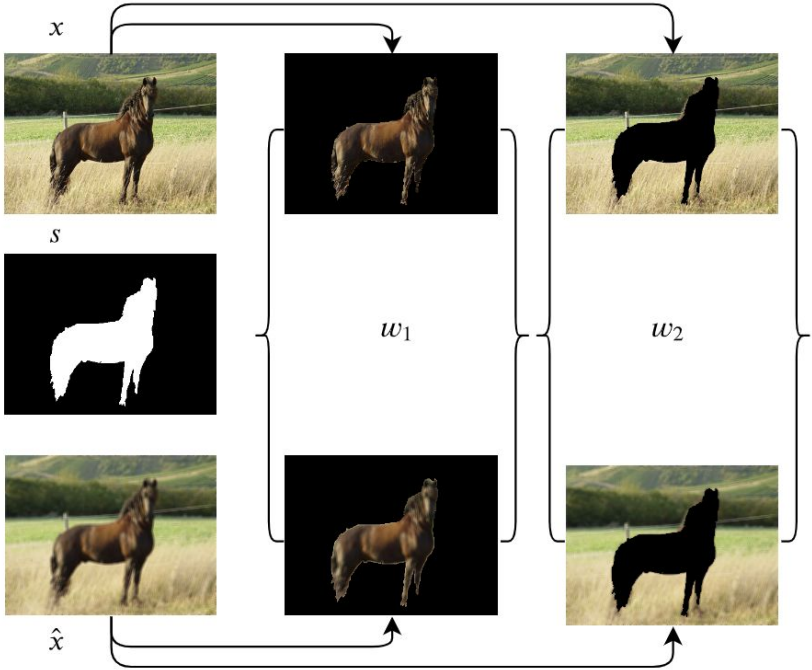
# Saliency Matters



Saliency Masking



Weighted distortion loss

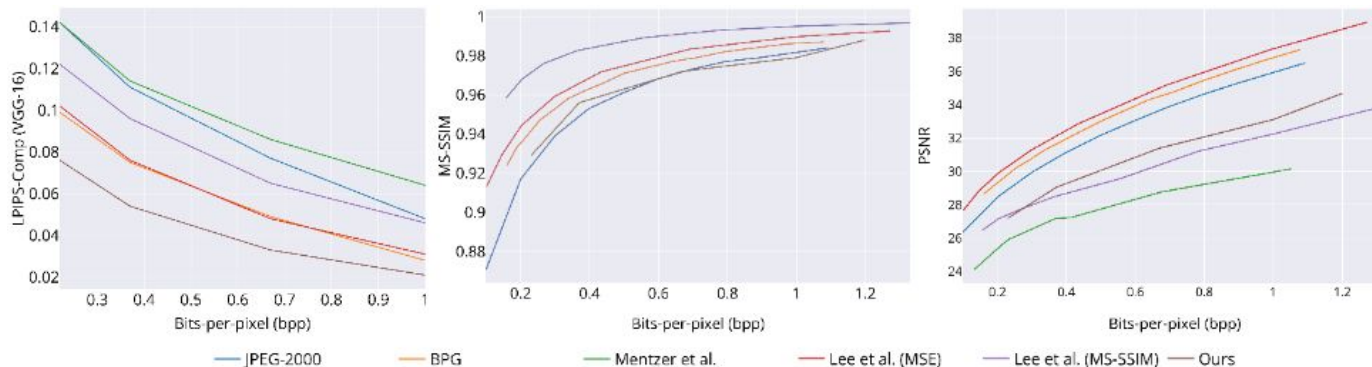


Patel et al., Saliency Driven Perceptual Image Compression, WACV 2021.

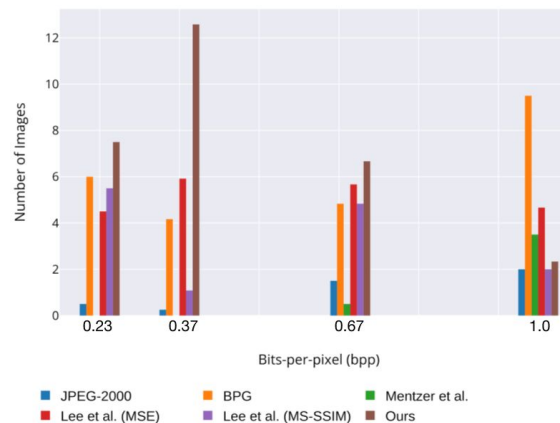
# Compression Results



## Distortion metrics



## Human Evaluations



# Compression Results - Object Detection



## Object Detection

Method	0.23	0.37	0.67	1.0
JPEG-2000 [37]	23.2	29.1	34.4	36.8
BPG [5]	25.2	<u>32.5</u>	35.4	<u>37.7</u>
Mentzer <i>et al.</i> [26]	25.5	30.2	34.5	36.6
Lee <i>et al.</i> [22] (MSE)	<u>28.3</u>	-	<u>36.2</u>	37.6
Lee <i>et al.</i> [22] (MS-SSIM)	27.2	<u>32.5</u>	-	37.6
Ours (MSE + DPL)	<b>29.3</b>	<b>33.7</b>	<b>36.6</b>	<b>37.9</b>

Object Detection on MS-COCO 2017 validation set using Faster-RCNN.

## Instance Segmentation

Method	0.23	0.37	0.67	1.0
JPEG-2000 [37]	20.2	25.4	30.1	<u>32.2</u>
BPG [5]	22.0	28.5	30.8	<u>32.2</u>
Mentzer <i>et al.</i> [26]	9.3	10.5	11.9	22.0
Lee <i>et al.</i> [22] (MSE)	<u>25.4</u>	-	<u>32.2</u>	<b>33.2</b>
Lee <i>et al.</i> [22] (MS-SSIM)	25.1	<u>28.9</u>	-	<b>33.2</b>
Ours (MSE + DPL)	<b>26.1</b>	<b>30.0</b>	<b>32.3</b>	<b>33.2</b>

Instance segmentation on MS-COCO 2017 validation set using Mask-RCNN.



# Compression Results - Qualitative



# Conclusions



This paper makes following contributions:

1. An adequate compression specific perceptual similarity metric.
2. Incorporating saliency for image compression.
3. A hierarchical auto-regressive model for image compression.

Results:

1. The proposed perceptual similarity metric aligns well with human perception of similarity.
2. The method generates image that are visually better and are useful for subsequent vision tasks such as object detection and image segmentation.
3. Object detection and image segmentation as an evaluation metric aligns with human perception of similarity.

Link to the paper: <https://arxiv.org/abs/2002.04988>

Link to the supplementary material: [https://yash0307.github.io/SDPIC\\_WACV2021\\_Supplementary\\_Material.pdf](https://yash0307.github.io/SDPIC_WACV2021_Supplementary_Material.pdf)

*Patel et al., Saliency Driven Perceptual Image Compression, WACV 2021.*

# Learning Surrogates via Deep Embedding

Yash Patel   Tomas Hodan   Jiri Matas

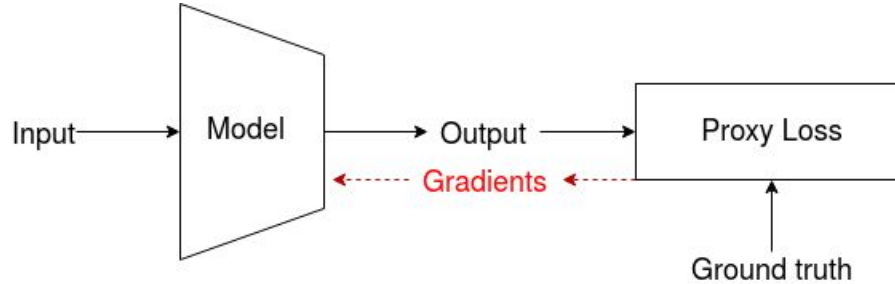
Visual Recognition Group, Czech Technical University in Prague

European Conference on Computer Vision (ECCV), 2020

Project webpage: [https://yash0307.github.io/LS\\_ECCV2020](https://yash0307.github.io/LS_ECCV2020)

# Proxy Losses

Proxy losses may not always align with the evaluation metric.

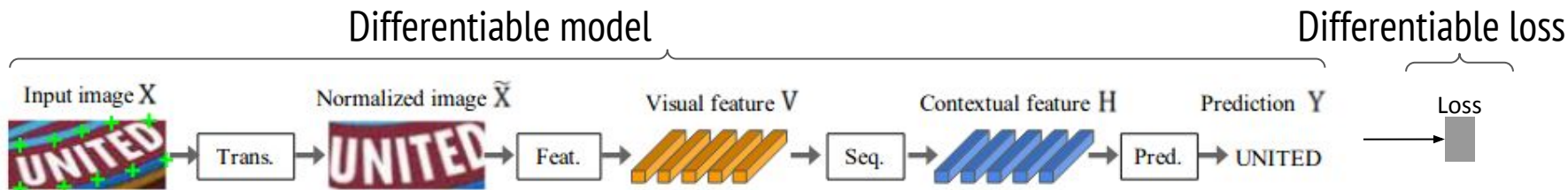


Evaluation metric in green, proxy loss functions in blue.

- Image compression (human perception of similarity): structural similarity index, peak signal to noise ratio, mean squared error, etc.
- Object detection (intersection over union): smooth-L1 distance, L2 distance.
- Scene text recognition (edit distance): per-character cross entropy, connectionist temporal

● Note: These functions are decomposable, *i.e.*, for a prediction there is a fixed target.

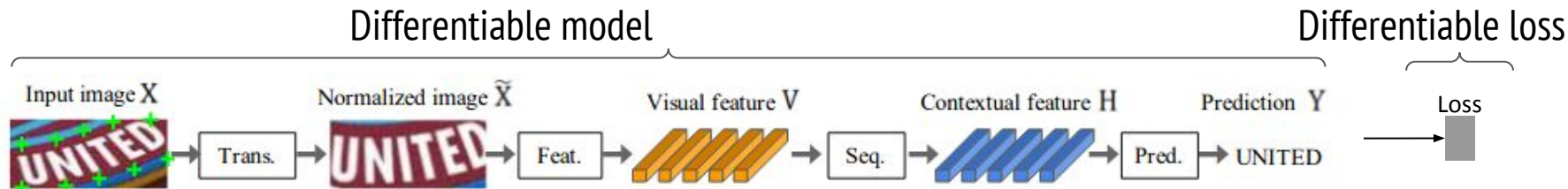
# Proxy Loss for Scene Text Recognition



*Baek et al., What Is Wrong With Scene Text Recognition Model Comparisons? Dataset and Model Analysis, ICCV 2019.*

The loss used for the end-to-end training of scene text recognition models is per-character **cross-entropy**.

# Proxy Loss for Scene Text Recognition






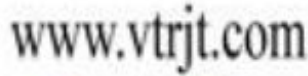
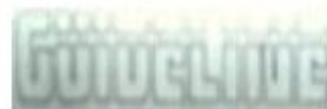
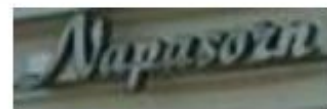
*Baek et al., What Is Wrong With Scene Text Recognition Model Comparisons? Dataset and Model Analysis, ICCV 2019.*

The loss used for the end-to-end training of scene text recognition models is per-character **cross-entropy**.

		<u>k</u>	<u>i</u>	<u>t</u>	<u>t</u>	<u>e</u>	<u>n</u>
	0	1	2	3	4	5	6
<u>s</u>	1	<u>1</u>	2	3	4	5	6
<u>i</u>	2	2	<u>1</u>	2	3	4	5
<u>t</u>	3	3	2	<u>1</u>	2	3	4
<u>t</u>	4	4	3	2	<u>1</u>	2	3
<u>i</u>	5	5	4	3	2	<u>2</u>	3
<u>n</u>	6	6	5	4	3	3	<u>2</u>
<u>g</u>	7	7	6	5	4	4	<u>3</u>

The evaluation metric for scene text recognition is **edit distance** computed using dynamic programming.

# Proxy Loss for Scene Text Recognition

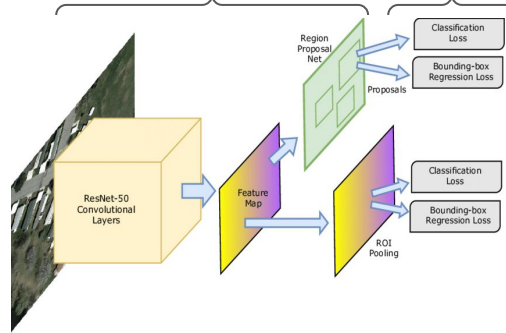
 <p>GT: wwwshutterstockcom</p> <p>M1: wonnishultersock Acc.=0, ED=9 M2: wnwishuttersockcom Acc.=0, ED=3</p>	 <p>GT: lilliput</p> <p>M1: umlout Acc.=0, ED=5 M2: lilidut Acc.=0, ED=2</p>	 <p>GT: frikkie</p> <p>M1: exikive Acc.=0, ED=4 M2: erikkie Acc.=0, ED=1</p>
 <p>GT: wwwvtrjtcom</p> <p>M1: wwwitcom Acc.=0, ED=4 M2: wwwvtritcom Acc.=0, ED=1</p>	 <p>GT: guiucliinc</p> <p>M1: guidglitus Acc.=0, ED=5 M2: guiuclituc Acc.=0, ED=2</p>	 <p>GT: napasorn</p> <p>M1: napasozin Acc.=0, ED=3 M2: napasozn Acc.=0, ED=1</p>

M1: Total Accuracy = 0, Total ED = 30

M2: Total Accuracy = 0, Total ED = 10

# Proxy Loss for Object Detection

Differentiable model    Differentiable losses



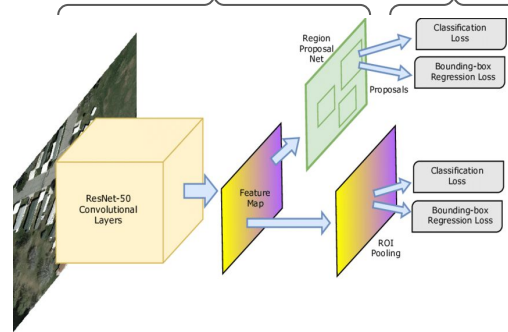
*Ren et al, Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, NeurIPS 2015.*

The regression loss used for the training of Faster R-CNN model is **smooth-L1**.




# Proxy Loss for Object Detection

Differentiable model    Differentiable losses



*Ren et al, Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, NeurIPS 2015.*

The regression loss used for the training of Faster R-CNN model is **smooth-L1**.

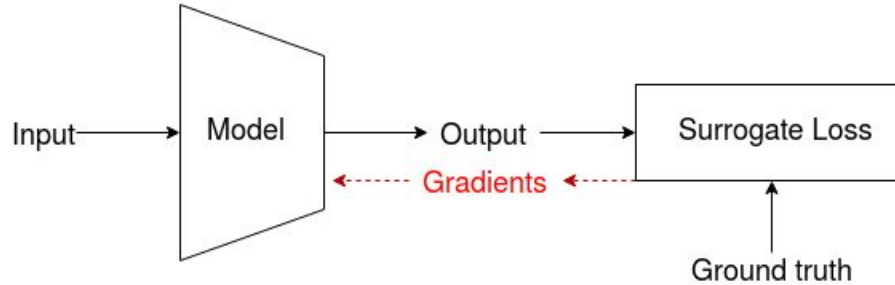
$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$


The bounding boxes are evaluated using **IoU**, which is not differentiable if the intersection cannot be expressed as an explicit function of predicted and ground truth bounding boxes.

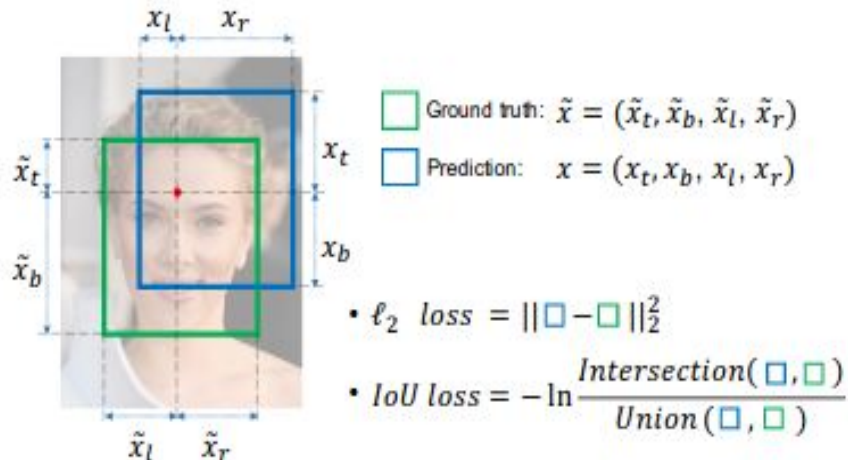
# Hand-Crafted Surrogate Losses

**Hand-crafting surrogate losses requires domain expertise.**

Example: Axis-aligned IoU loss in object detection.



# Hand-Crafted Surrogate for Object Detection



---

## Algorithm 1: $IoU$ loss Forward

---

**Input:**  $\tilde{x}$  as bounding box ground truth

**Input:**  $x$  as bounding box prediction

**Output:**  $\mathcal{L}$  as localization error

**for** each pixel  $(i, j)$  **do**

**if**  $\tilde{x} \neq 0$  **then**

$X = (x_t + x_b) * (x_l + x_r)$

$\tilde{X} = (\tilde{x}_t + \tilde{x}_b) * (\tilde{x}_l + \tilde{x}_r)$

$I_h = \min(x_t, \tilde{x}_t) + \min(x_b, \tilde{x}_b)$

$I_w = \min(x_l, \tilde{x}_l) + \min(x_r, \tilde{x}_r)$

$I = I_h * I_w$

$U = X + \tilde{X} - I$

$IoU = \frac{I}{U}$

$\mathcal{L} = -\ln(IoU)$

**else**

$\mathcal{L} = 0$

**end**

**end**

---

*Yu et al, UnitBox: An Advanced Object Detection Network, ACM-MM 2016.*

The hand-crafted IoU loss has shown improvements compared to using proxy losses.

The hand-crafted IoU loss assumes that the bounding boxes are axis aligned.

# Hand-Crafted Surrogate for Object Detection



Scene text detection

*Karatzas et al., ICDAR 2015 competition on robust reading, ICDAR'15.*



Object detection in Aerial Images

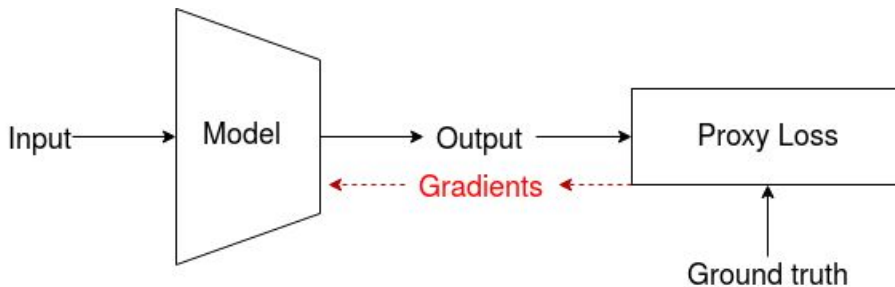
*Xia et al, DOTA: A Large-scale Dataset for Object Detection in Aerial Images, CVPR'18.*

The hand-crafted IoU loss does not generalize to the rotated bounding boxes.

# Proxy and Hand-Crafted Surrogate Losses

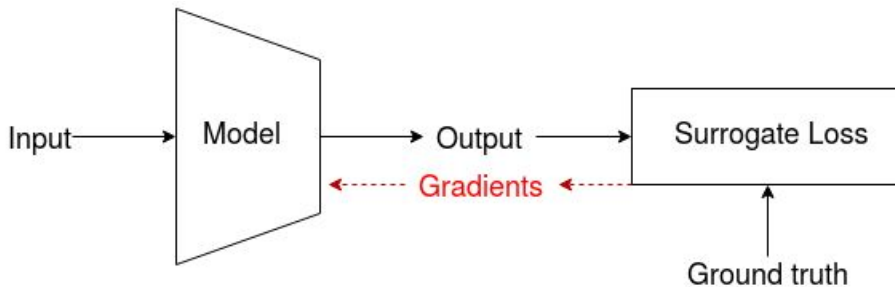
**Proxy losses may not always align with the evaluation metric.**

Examples: Smooth-L1, L2 loss in object detection; cross entropy loss in scene text recognition.



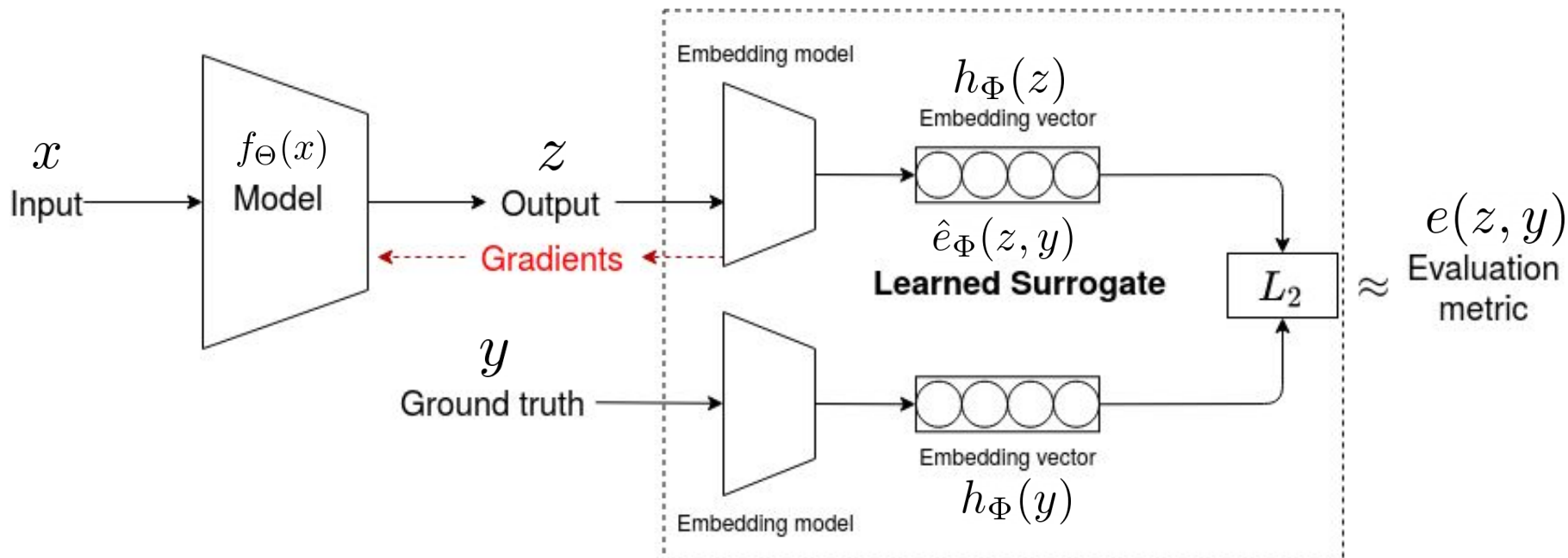
**Hand-crafting surrogate losses require domain expertise.**

Example: Axis-aligned IoU loss in object detection.



# Learning Surrogates via Deep Embedding

The surrogate is learned via a deep embedding where the Euclidean distance between the prediction and the ground truth corresponds to the value of the evaluation metric.



# Objectives for Learning the Surrogate

The learned surrogate corresponds to the value of the evaluation metric:

$$e(z, y) \approx \hat{e}_{\Phi}(z, y)$$

# Objectives for Learning the Surrogate

The learned surrogate corresponds to the value of the evaluation metric:

$$e(z, y) \approx \hat{e}_{\Phi}(z, y)$$

The first order derivative of the learned surrogate with respect to the prediction is close to 1:

$$\left\| \frac{\partial \hat{e}_{\Phi}(z, y)}{\partial z} \right\|_2 \approx 1$$



# Objectives for Learning the Surrogate

The learned surrogate corresponds to the value of the evaluation metric:

$$e(z, y) \approx \hat{e}_{\Phi}(z, y)$$

The first order derivative of the learned surrogate with respect to the prediction is close to 1:

$$\left\| \frac{\partial \hat{e}_{\Phi}(z, y)}{\partial z} \right\|_2 \approx 1$$

Overall loss for learning the surrogate:

$$\| \hat{e}_{\Phi}(z, y) - e(z, y) \|_2^2 + \lambda \left( \left\| \frac{\partial \hat{e}_{\Phi}(z, y)}{\partial z} \right\|_2 - 1 \right)^2$$

# Post-Tuning with the Learned Surrogate

---

**Algorithm 1** Training with LS (*local-global approximation*)

---

**Inputs:** Supervised data  $D$ , random data generator  $R$ , evaluation metric  $e$ .

**Hyper-parameters:** Number of update steps  $I_a$  and  $I_b$ , learning rates  $\eta_a$  and  $\eta_b$ , number of epochs  $E$ .

**Objective:** Train the model for a given task that is  $f_\Theta(x)$  and the surrogate, *i.e.*,  $e_\Phi$ .

1: *Initialize*  $\Theta \leftarrow$  pre-trained weights,  $\Phi \leftarrow$  random weights.

2: **for** epoch = 1,...,E **do**

3:   **for**  $i = 1, \dots, I_a$  **do**

4:     sample  $(x, y) \sim P_D$ , sample  $(z_r, y_r) \sim P_R$

5:     inference  $z = f_{\Theta^{epoch-1}}(x)$

6:     compute loss  $l_{\hat{e}} = \text{loss}(z, y) + \text{loss}(z_r, y_r)$  (Equation 6)

7:      $\Phi^i \leftarrow \Phi^{i-1} - \eta_a \frac{\partial l_{\hat{e}}}{\partial \Phi^{i-1}}$

8:   **end for**

9:    $\Phi \leftarrow \Phi^{I_a}$

10:   **for**  $i = 1, \dots, I_b$  **do**

11:     sample  $(x, y) \sim P_D$

12:     inference  $z = f_{\Theta^{i-1}}(x)$

13:     compute loss  $l_f = \hat{e}_{\Phi^{epoch}}(z, y)$  (Equation 3)

14:      $\Theta^i \leftarrow \Theta^{i-1} - \eta_b \frac{\partial (l_f)}{\partial \Theta^{i-1}}$

15:   **end for**

16:    $\Theta \leftarrow \Theta^{I_b}$

17: **end for**

---

Learning the  
surrogate

Post-tuning with  
the surrogate

# Results on Scene Text Recognition

Tuning a scene text recognition model on the learned surrogate of edit distance (LS-ED) yields up to **39% improvement** on total edit distance.

Test Data	Loss Function	↑ Acc.	↑ NED	↓ TED
IIIT-5K	Cross-Entropy	84.300	0.954	945
IIIT-5K	LS-ED	86.300 +2.37%	0.953 -0.10%	837 +11.42%
SVT	Cross-Entropy	84.699	0.940	229
SVT	LS-ED	86.399 +2.00%	0.947 +0.74%	196 +14.41%
ICDAR'03	Cross-Entropy	92.558	0.972	151
ICDAR'03	LS-ED	94.070 +1.63%	0.977 +0.51%	119 +26.89%
ICDAR'13	Cross-Entropy	89.754	0.949	260
ICDAR'13	LS-ED	91.133 +1.53%	0.960 +1.15%	157 +39.61%
ICDAR'15	Cross-Entropy	71.452	0.889	1135
ICDAR'15	LS-ED	74.655 +4.48%	0.899 +1.12%	1013 +10.74%
SVTP	Cross-Entropy	74.109	0.891	424
SVTP	LS-ED	77.519 +4.60%	0.901 +1.22%	381 +10.14%
CUTE	Cross-Entropy	68.293	0.838	285
CUTE	LS-ED	71.777 +5.10%	0.868 +3.57%	234 +17.89%

# Results on Scene Text Detection

Tuning a scene text detection model on the learned surrogate of IoU for rotated bounding boxes yields a **4.25% improvement** on the F1 score.

Loss Function	↑ Recall	↑ Precision	↑ $F_1$ score
$Smooth-L_1$	71.21%	84.71%	77.37%
LS-IoU	76.79% <b>+7.83%</b>	84.93% <b>+0.25%</b>	80.66% <b>+4.25%</b>

*Ma et al, Arbitrary-Oriented Scene Text Detection via Rotation Proposals, IEEE Transactions on Multimedia 2018.*

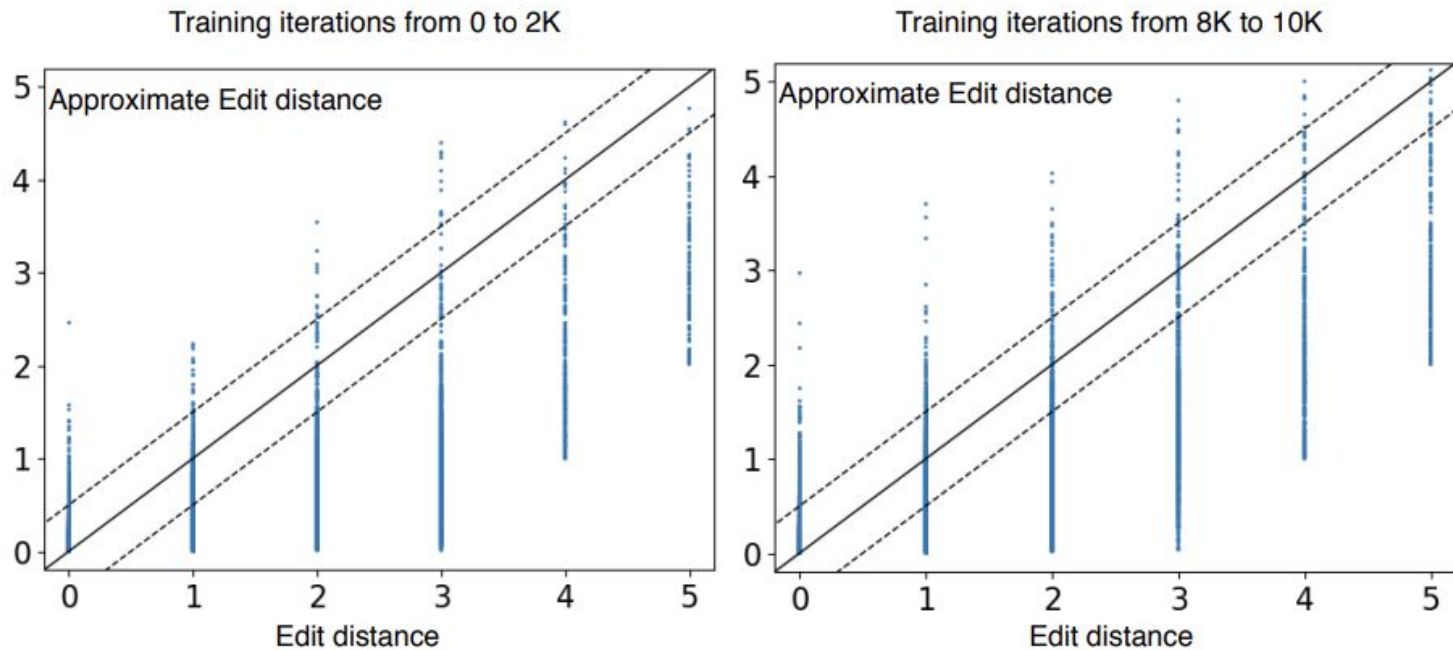
# FEDS – Filtered Edit Distance Surrogate

Yash Patel    Jiri Matas

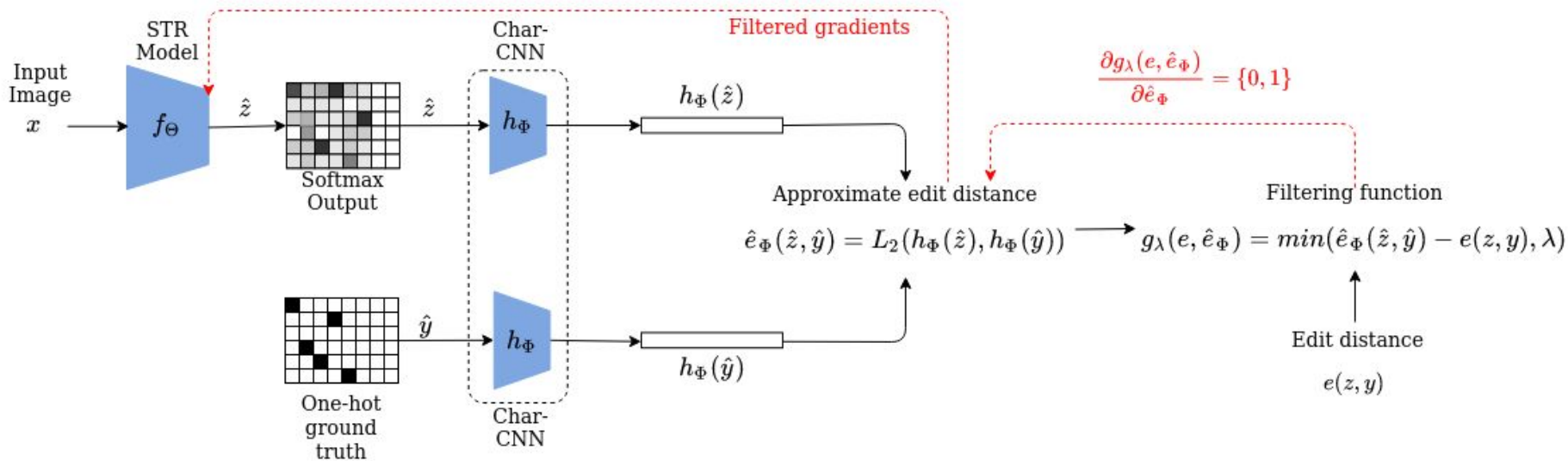
Visual Recognition Group, Czech Technical University in Prague

International Conference on Document Analysis and Recognition (ICDAR), 2021

# Quality of Approximation



# FEDS -- Filtered Edit Distance Surrogate

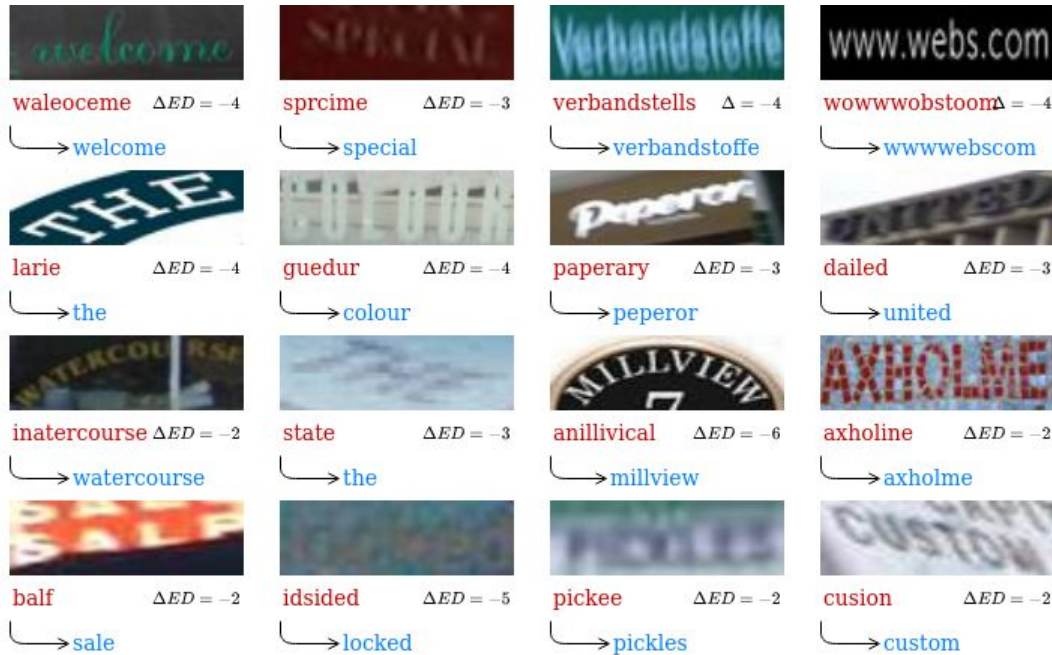


# FEDS -- Quantitative Results

Test Data	Loss Function	↑ Acc.	↑ NED	↓ TED
Synthetic Training data	Cross-Entropy [1]	85.6	0.941	4079
	LS-ED [29]	86.1 +0.61%	0.942 +0.18%	3832 +6.05%
	FEDS	86.5 +0.98%	0.946 +0.48%	3623 +11.2%
Additional weakly labelled data	Cross-Entropy [1]	88.7	0.953	3050
	LS-ED [29]	89.0 +0.41%	0.954 +0.62%	2961 +2.91%
	FEDS	89.6 +1.01%	0.956 +0.35%	2809 +7.90%



# FEDS -- Qualitative Results



# Recall@k Surrogate Loss with Large Batches and Similarity Mixup

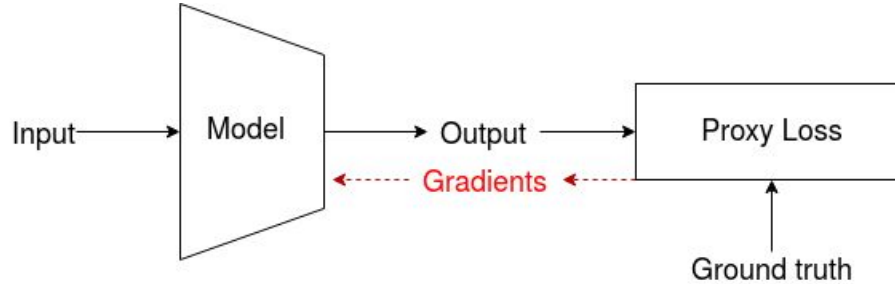
Yash Patel   Giorgos Tolias   Jiri Matas

Visual Recognition Group, Czech Technical University in Prague

IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022

# Proxy Losses

Proxy losses may not always align with the evaluation metric.



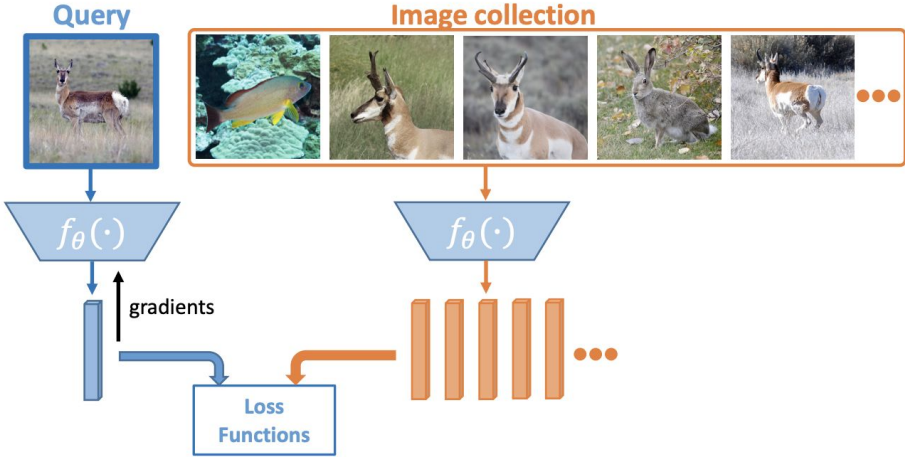
Evaluation metric in green, proxy loss functions in blue.

- **Image classification (cross entropy, softmax):** cross entropy, softmax
- **Note: These functions are non-decomposable, *i.e.*, for one fixed prediction there is no fixed target. It rather depends on comparisons within a set.**
- **Scene text recognition (edit distance):** per character cross entropy, connectionist temporal classification.
- **Image retrieval (mean average precision, recall@k):** contrastive loss, triplet loss, proxy NCA, etc.

And many more...

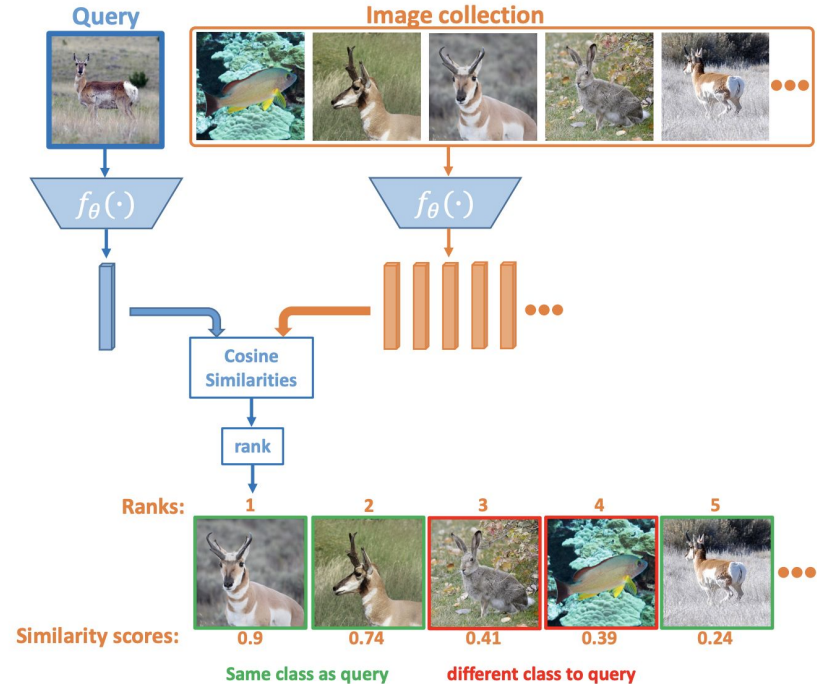
# Image Retrieval Training

- Embedding must be trained for good ranking.
- Achieved using loss functions.



# Image Retrieval Inference

- Extract embeddings from query and image collection.
- Compute similarity scores.
- Rank according to relevance to the query.



# Recall@k

Heaviside step function

Recall@k

$$R_{\Omega}^k(q) = \frac{\sum_{x \in P_q} H(k - r_{\Omega}(q, x))}{|P_q|}$$

Rank of a positive instance  $x$ , given a query  $q$

Total number of positive samples

# Recall@k

Heaviside step function

Recall@k

$$R_{\Omega}^k(q) = \frac{\sum_{x \in P_q} H(k - r_{\Omega}(q, x))}{|P_q|}$$

Rank of a positive instance  $x$ , given a query  $q$

Total number of positive samples

Rank

$$r_{\Omega}(q, x) = 1 + \sum_{z \in \Omega, z \neq x} H(s_{qz} - s_{qx})$$

Similarity score (cosine for us) between the a positive sample  $z$  and the query  $q$

Similarity score (cosine for us) between the a positive sample  $x$  and the query  $q$

# Recall@k

Heaviside step function

Recall@k

$$R_{\Omega}^k(q) = \frac{\sum_{x \in P_q} H(k - r_{\Omega}(q, x))}{|P_q|}$$

Rank of a positive instance  $x$ , given a query  $q$

Total number of positive samples

Rank

$$r_{\Omega}(q, x) = 1 + \sum_{z \in \Omega, z \neq x} H(s_{qz} - s_{qx})$$

Similarity score (cosine for us) between the a positive sample  $z$  and the query  $q$

Similarity score (cosine for us) between the a positive sample  $x$  and the query  $q$

Recall@k

$$R_{\Omega}^k(q) = \frac{\sum_{x \in P_q} H(k - 1 - \sum_{z \in \Omega, z \neq x} H(s_{qz} - s_{qx}))}{|P_q|}$$



# Recall@k

Heaviside step function

Recall@k

$$R_{\Omega}^k(q) = \frac{\sum_{x \in P_q} H(k - r_{\Omega}(q, x))}{|P_q|}$$

Rank of a positive instance  $x$ , given a query  $q$

Total number of positive samples

Rank

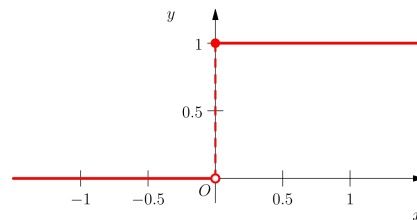
$$r_{\Omega}(q, x) = 1 + \sum_{z \in \Omega, z \neq x} H(s_{qz} - s_{qx})$$

Similarity score (cosine for us) between the a positive sample  $z$  and the query  $q$

Similarity score (cosine for us) between the a positive sample  $x$  and the query  $q$

Recall@k

$$R_{\Omega}^k(q) = \frac{\sum_{x \in P_q} H(k - 1 - \sum_{z \in \Omega, z \neq x} H(s_{qz} - s_{qx}))}{|P_q|}$$



# Recall@k Surrogate (RS@k)

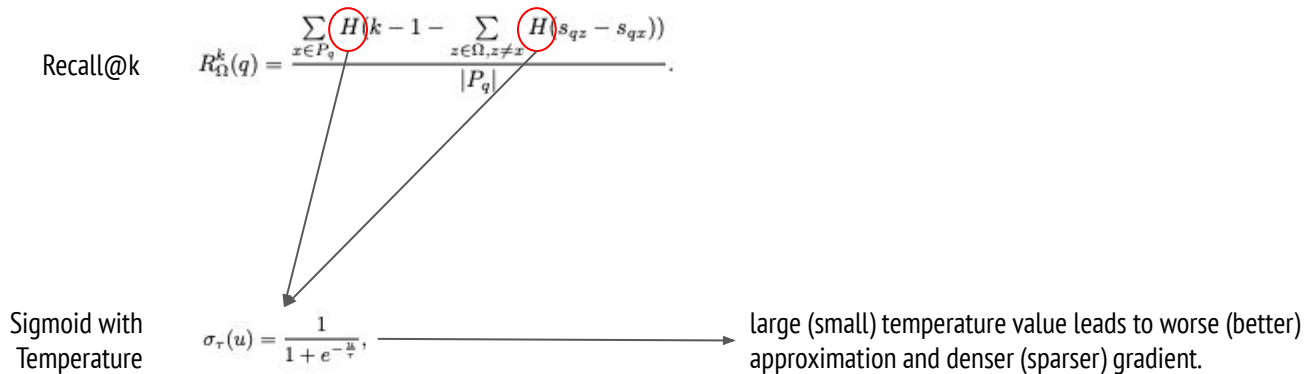
Recall@k

$$R_{\Omega}^k(q) = \frac{\sum_{x \in P_q} H(k - 1 - \sum_{z \in \Omega, z \neq x} H(s_{qz} - s_{qx}))}{|P_q|}.$$

Sigmoid with Temperature

$$\sigma_{\tau}(u) = \frac{1}{1 + e^{-\frac{u}{\tau}}},$$

large (small) temperature value leads to worse (better) approximation and denser (sparser) gradient.



# Recall@k Surrogate (RS@k)

Recall@k

$$R_{\Omega}^k(q) = \frac{\sum_{x \in P_q} H(k - 1 - \sum_{z \in \Omega, z \neq x} H(s_{qz} - s_{qx}))}{|P_q|}.$$

Sigmoid with Temperature

$$\sigma_{\tau}(u) = \frac{1}{1 + e^{-\frac{u}{\tau}}},$$

large (small) temperature value leads to worse (better) approximation and denser (sparser) gradient.

Recall@k Surrogate

$$\tilde{R}_{\Omega}^k(q) = \frac{\sum_{x \in P_q} \sigma_{\tau_1}(k - 1 - \sum_{\substack{z \in \Omega \\ z \neq x}} \sigma_{\tau_2}(s_{qz} - s_{qx}))}{|P_q|}$$

# Recall@k Surrogate (RS@k)

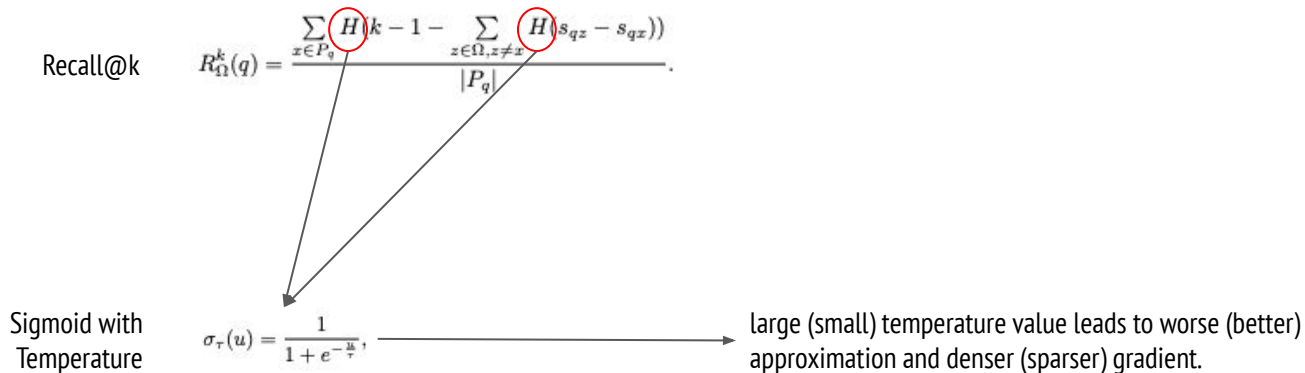
Recall@k

$$R_{\Omega}^k(q) = \frac{\sum_{x \in P_q} H(k-1 - \sum_{z \in \Omega, z \neq x} H(s_{qz} - s_{qx}))}{|P_q|}$$

Sigmoid with Temperature

$$\sigma_{\tau}(u) = \frac{1}{1 + e^{-\frac{u}{\tau}}}$$

large (small) temperature value leads to worse (better) approximation and denser (sparser) gradient.



Recall@k Surrogate

$$\tilde{R}_{\Omega}^k(q) = \frac{\sum_{x \in P_q} \sigma_{\tau_1}(k-1 - \sum_{\substack{z \in \Omega \\ z \neq x}} \sigma_{\tau_2}(s_{qz} - s_{qx}))}{|P_q|}$$

Loss

$$L^k(q) = 1 - \tilde{R}_{B \setminus q}^k(q)$$

Loss over multiple values of k

$$L^K(q) = \frac{1}{|K|} \sum_{k \in K} L^k(q)$$



# Recall@k Surrogate (RS@k)

Recall@k  
Surrogate

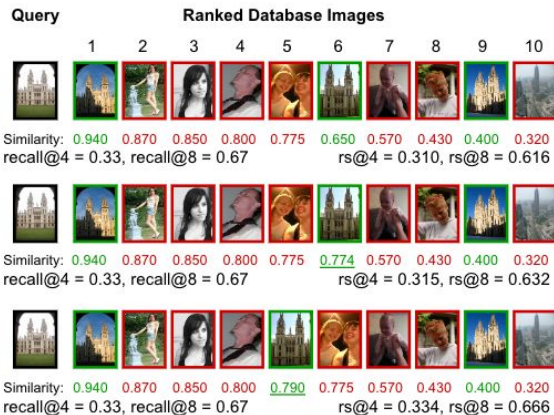
$$\tilde{R}_{\Omega}^k(q) = \frac{\sum_{x \in P_q} \sigma_{\tau_1}(k-1 - \sum_{\substack{z \in \Omega \\ z \neq x}} \sigma_{\tau_2}(s_{qz} - s_{qx}))}{|P_q|}$$



# Recall@k Surrogate (RS@k)

Recall@k  
Surrogate

$$\tilde{R}_{\Omega}^k(q) = \frac{\sum_{x \in P_q} \sigma_{\tau_1}(k-1 - \sum_{\substack{z \in \Omega \\ z \neq x}} \sigma_{\tau_2}(s_{qz} - s_{qx}))}{|P_q|}$$



# Recall@k Surrogate (RS@k)

Recall@k  
Surrogate

$$\tilde{R}_{\Omega}^k(q) = \frac{\sum_{x \in P_q} \sigma_{\tau_1}(k-1 - \sum_{\substack{z \in \Omega \\ z \neq x}} \sigma_{\tau_2}(s_{qz} - s_{qx}))}{|P_q|}$$



Figure 1. A comparison between recall@k and rs@k, the proposed differentiable recall@k surrogate. Examples show a query, the ranked database images sorted according to the similarity and the corresponding values for recall@k and rs@k and their dependence on similarity score change. Note that the values of recall@k and rs@k are close. Changes to similarity and ranking in some cases may not affect the original recall@k but can affect the surrogate, with the latter having a more significant impact than the former. Similarity values of all negatives are fixed for ease of understanding. The similarity values of the positives that were changed in rows 2, 3 and 4 are underlined.

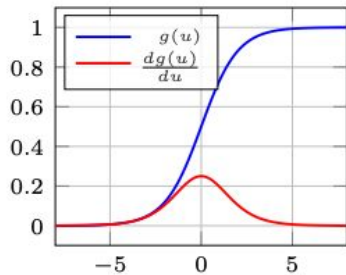


# Visualization for Temperatures

Recall@k  
Surrogate

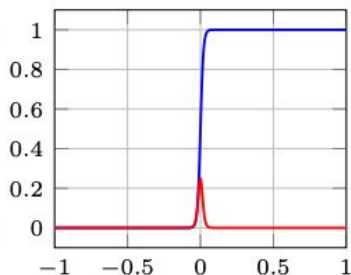
$$\tilde{R}_{\Omega}^k(q) = \frac{\sum_{x \in P_q} \sigma_{\tau_1}(k-1 - \sum_{\substack{z \in \Omega \\ z \neq x}} \sigma_{\tau_2}(s_{qz} - s_{qx}))}{|P_q|}$$

$$g(u) = \sigma_{\tau_1}(u), \tau_1 = 1$$



$$u = k - 1 - r_{\Omega}(q, x)$$

$$g(u) = \sigma_{\tau_2}(u), \tau_2 = 0.01$$



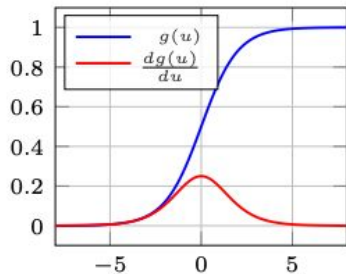
$$u = s_{qz} - s_{qx}$$

# Visualization for Temperatures

Recall@k  
Surrogate

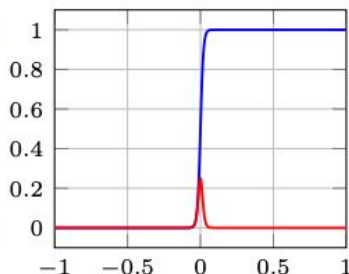
$$\tilde{R}_{\Omega}^k(q) = \frac{\sum_{x \in P_q} \sigma_{\tau_1}(k-1 - \sum_{\substack{z \in \Omega \\ z \neq x}} \sigma_{\tau_2}(s_{qz} - s_{qx}))}{|P_q|}$$

$$g(u) = \sigma_{\tau_1}(u), \tau_1 = 1$$

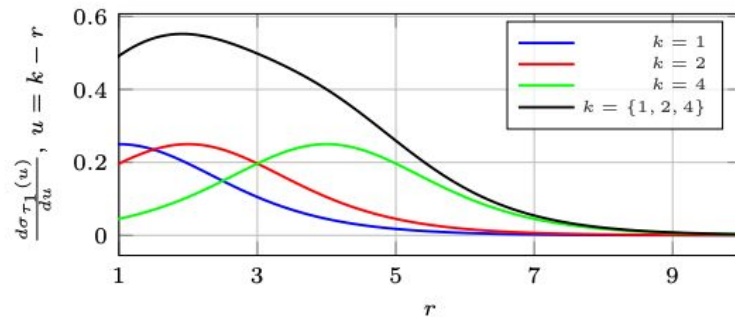


$$u = k - 1 - r_{\Omega}(q, x)$$

$$g(u) = \sigma_{\tau_2}(u), \tau_2 = 0.01$$



$$u = s_{qz} - s_{qx}$$



# Similarity Mixup (SiMix)

Synthetic sample by Mixup  $\mathbf{v}_{xz\alpha} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{z} \mid \alpha \sim U(0, 1)$ ,  $\longrightarrow$  The embedding for the virtual sample is not required in loss computation

# Similarity Mixup (SiMix)

Synthetic sample by Mixup  $\mathbf{v}_{xz\alpha} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{z} \mid \alpha \sim U(0, 1)$ ,  $\longrightarrow$  The embedding for the virtual sample is not required in loss computation

Virtual sample by SiMix  $s(w, xz\alpha) = \mathbf{w}^\top \mathbf{v}_{xz\alpha} = \alpha s_{wx} + (1 - \alpha) s_{wz}$ ,  $\longrightarrow$  Similarity scores between a real and a virtual sample can be directly computed without the embedding

# Similarity Mixup (SiMix)

Synthetic sample by Mixup  $\mathbf{v}_{xz\alpha} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{z} \quad | \quad \alpha \sim U(0, 1),$   $\longrightarrow$  The embedding for the virtual sample is not required in loss computation

Virtual sample by SiMix  $s(w, xz\alpha) = \mathbf{w}^\top \mathbf{v}_{xz\alpha} = \alpha s_{wx} + (1 - \alpha)s_{wz},$   $\longrightarrow$  Similarity scores between a real and a virtual sample can be directly computed without the embedding

$$\begin{aligned} s(xz\alpha_1, yw\alpha_2) &= \mathbf{v}_{xz\alpha_1}^\top \mathbf{v}_{yw\alpha_2} \\ &= \alpha_1\alpha_2 s_{xy} + (1 - \alpha_1)(1 - \alpha_2)s_{zw} \\ &\quad + \alpha_1(1 - \alpha_2)s_{xw} + (1 - \alpha_1)\alpha_2 s_{zy}. \end{aligned}$$

$\longrightarrow$  Similarity scores between two virtual samples can also be directly computed without the embedding

# Similarity Mixup (SiMix)

Synthetic sample by Mixup  $\mathbf{v}_{xz\alpha} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{z} \quad | \quad \alpha \sim U(0, 1)$ ,  $\longrightarrow$  The embedding for the virtual sample is not required in loss computation

Virtual sample by SiMix  $s(w, xz\alpha) = \mathbf{w}^\top \mathbf{v}_{xz\alpha} = \alpha s_{wx} + (1 - \alpha)s_{wz}$ ,  $\longrightarrow$  Similarity scores between a real and a virtual sample can be directly computed without the embedding

$$\begin{aligned} s(xz\alpha_1, yw\alpha_2) &= \mathbf{v}_{xz\alpha_1}^\top \mathbf{v}_{yw\alpha_2} \\ &= \alpha_1\alpha_2 s_{xy} + (1 - \alpha_1)(1 - \alpha_2)s_{zw} \\ &\quad + \alpha_1(1 - \alpha_2)s_{xw} + (1 - \alpha_1)\alpha_2 s_{zy}. \end{aligned}$$

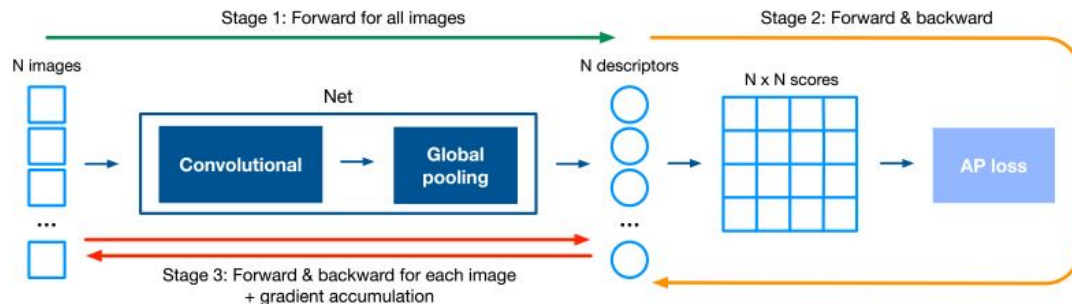
$\longrightarrow$  Similarity scores between two virtual samples can also be directly computed without the embedding

Unlike other mixup techniques:

- The embedding of “virtual” sample is never computed. This makes SiMix more computationally and memory efficient.
- The “virtual” sample is used as a positive, negative and a query.

# Training with Large Batches

- **Step-1:** Iteratively feed-forward through the large batch and only store the embeddings (discard model activations).
- **Step-2:** Iteratively compute the loss.
- **Step-3:** Compute gradients of loss with respect to the embeddings.
- **Step-4:** Iteratively feed-forward through the large batch and backpropagation through the model.
- **Step-5:** Update model weights.



Revaud et al., *Learning with Average Precision: Training Image Retrieval with a Listwise Loss*, ICCV 2019.

Note: A public implementation of this training procedure does not exist. We will release it with Camera Ready.

Patel et al., *Recall@k Surrogate Loss with Large Batches and Similarity Mixup*, CVPR 2022.

# Recall@k Surrogate Loss with Large Batches and Similarity Mixup

---

## Algorithm 1 Training with RS@k and SiMix.

---

```
1: procedure TRAIN-RS@K( $X, Y, M, m$ )
2:    $X$  : training images
3:    $Y$  : class labels
4:    $M$  : mini-batch size
5:    $m$  : number of images per class in mini-batch
6:
7:    $\theta \leftarrow$  initialize according to pre-training            $\triangleright$  use ImageNet
8:   for iteration  $\in [1, \dots, \text{number-of-iterations}]$  do
9:      $loss \leftarrow 0$                                           $\triangleright$  set batch loss to zero
10:     $B \leftarrow$  BATCH-SAMPLER( $X, Y, M, m$ )
11:     $\hat{B} \leftarrow$  VIRTUAL-BATCH( $B$ )                           $\triangleright$  enumerate virtual examples
12:    for  $(x, z) \in B \times B$  do compute  $s(x, z)$             $\triangleright$  use  $\mathbf{x}^\top \mathbf{z}$ 
13:    for  $(x, z) \in B \times \hat{B}$  do compute  $s(x, z)$           $\triangleright$  use (9)
14:    for  $(x, z) \in \hat{B} \times \hat{B}$  do compute  $s(x, z)$           $\triangleright$  use (10)
15:     $B \leftarrow B \cup \hat{B}$                                      $\triangleright$  expand batch with virtual examples
16:    for  $q \in B$  do                                          $\triangleright$  use each image in the batch as query
17:       $loss \leftarrow loss + L^K(q)$                            $\triangleright$  Recall@k loss (7)
18:    end for
19:     $\theta \leftarrow$  MINIMIZE( $\frac{loss}{|B|}$ )                        $\triangleright$  SGD update
20:  end for
21: end procedure
```



# Datasets

Dataset	#Images	#Classes	#Avg
iNaturalist Train [58]	325,846	5,690	57.3
iNaturalist Test [58]	136,093	2,452	55.5
VehicleID Train [28]	110,178	13,134	8.4
VehicleID Test [28]	40,365	4,800	8.4
SOP Train [37]	59,551	11,318	5.3
SOP Test [37]	60,502	11,316	5.3
Cars196 Train [26]	8,054	98	82.1
Cars196 Test [26]	8,131	98	82.9
$\mathcal{R}$ Oxford [42]	4,993	11	n/a
$\mathcal{R}$ Paris [42]	6,322	11	n/a
GLDv1 [36]	1,060,709	12,894	82.3

Datasets are diverse in the number of training examples, the number of classes, and the number of examples per class, ranging from class balanced to long-tailed.

# Results - iNaturalist, SOP, VehicleID, Cars196

Method	Arch. dim	iNaturalist [58]				SOP [37]				VehicleID [28]						Cars196 [26]			
		r@k																	
		1	4	16	32	10 <sup>0</sup>	10 <sup>1</sup>	10 <sup>2</sup>	10 <sup>3</sup>	Small		Medium		Large		1	2	4	8
ProxyNCA [33]	$I_1^{128}$	<u>61.6</u>	<u>77.4</u>	<u>87.0</u>	<u>90.6</u>	73.7	-	-	-	-	-	-	-	-	-	73.2	82.4	86.4	88.7
Margin [66]	$R_{50}^{128}$	58.1	75.5	86.8	<u>90.7</u>	72.7	86.2	93.8	98.0	-	-	-	-	-	-	<u>79.6</u>	<u>86.5</u>	<u>91.9</u>	<u>95.1</u>
Divide [50]	$R_{50}^{128}$	-	-	-	-	75.9	88.4	94.9	98.1	87.7	92.9	85.7	90.4	82.9	90.2	-	-	-	-
MIC [47]	$R_{50}^{128}$	-	-	-	-	77.2	89.4	95.6	-	86.9	93.4	-	-	82.0	91.0	-	-	-	-
Cont. w/M [64]	$R_{50}^{128}$	-	-	-	-	<u>80.6</u>	<u>91.6</u>	<u>96.2</u>	<u>98.7</u>	<u>94.7</u>	<u>96.8</u>	<u>93.7</u>	<u>95.8</u>	<u>93.0</u>	<u>95.8</u>	-	-	-	-
RS@k <sup>†</sup>	$R_{50}^{128}$	69.3	82.9	90.6	93.1	80.6	91.6	96.4	<b>98.8</b>	<b>95.6</b>	<b>97.8</b>	<b>94.4</b>	<b>96.8</b>	<b>93.5</b>	<b>96.6</b>	78.1	85.8	91.1	94.5
RS@k <sup>†</sup> +SiMix	$R_{50}^{128}$	<b>69.6</b>	<b>83.3</b>	<b>91.2</b>	<b>93.8</b>	<b>80.9</b>	<b>91.7</b>	<b>96.5</b>	<b>98.8</b>	95.4	97.5	93.8	96.6	93.0	96.2	<b>84.7</b>	<b>90.9</b>	<b>94.7</b>	<b>96.9</b>
		+21%	+26%	+32%	+33%	+1.5%	+1.2%	+7.9%	+7.7%	+17%	+31%	+11%	+24%	+7.1%	+19%	+25%	+33%	+35%	+37%
FastAP [7]	$R_{50}^{512}$	60.6	77.0	87.2	90.6	76.4	89.0	95.1	98.2	91.9	96.8	90.6	95.9	87.5	95.1	-	-	-	-
MS [63]	$I_3^{512}$	-	-	-	-	78.2	90.5	96.0	98.7	-	-	-	-	-	-	84.1	90.4	94.0	96.1
NormSoftMax [68]	$R_{50}^{512}$	-	-	-	-	78.2	90.6	96.2	-	-	-	-	-	-	-	84.2	90.4	94.4	96.9
Blackbox AP [46]	$R_{50}^{512}$	62.9	79.0	88.9	92.1	78.6	90.5	96.0	98.7	-	-	-	-	-	-	-	-	-	-
Cont. w/M [64]	$I_3^{512}$	-	-	-	-	79.5	90.8	96.1	98.7	94.6	96.9	<u>93.4</u>	96.0	<u>93.0</u>	96.1	-	-	-	-
HORDE [23]	$R_{50}^{512}$	-	-	-	-	80.1	91.3	96.2	-	-	-	-	-	-	-	86.2	91.9	95.1	97.2
ProxyNCA++ [55]	$R_{50}^{512}$	-	-	-	-	<u>80.7</u>	<u>92.5</u>	96.7	98.9	-	-	-	-	-	-	<u>86.5</u>	<u>92.5</u>	<u>95.7</u>	<b>97.7</b>
SAP [6]	$R_{50}^{512}$	<u>67.2</u>	<u>81.8</u>	<u>90.3</u>	<u>93.1</u>	80.1	91.5	<u>96.6</u>	<u>99.0</u>	<u>94.9</u>	<u>97.6</u>	93.3	<u>96.4</u>	91.9	<u>96.2</u>	76.1	84.3	89.8	93.8
SAP <sup>†</sup> [6] +GeM +LN	$R_{50}^{512}$	68.7	82.7	90.9	93.5	80.3	92.0	96.9	99.0	94.2	97.2	92.7	96.2	91.0	95.8	78.2	85.6	90.8	94.3
RS@k <sup>†</sup>	$R_{50}^{512}$	71.2	84.0	91.3	93.6	<b>82.8</b>	<b>92.9</b>	<b>97.0</b>	99.0	<b>95.7</b>	<b>97.9</b>	<b>94.6</b>	<b>96.9</b>	<b>93.8</b>	<b>96.6</b>	80.7	88.3	92.8	95.7
RS@k <sup>†</sup> +SiMix	$R_{50}^{512}$	<b>71.8</b>	<b>84.7</b>	<b>91.9</b>	<b>94.3</b>	82.1	92.8	<b>97.0</b>	<b>99.1</b>	95.3	97.7	94.2	96.5	93.3	96.4	<b>88.2</b>	<b>93.0</b>	<b>95.9</b>	97.4
		+14%	+16%	+16%	+17%	+11%	+5.3%	+12%	+10%	+16%	+13%	+18%	+14%	+11%	+10%	+13%	+6.7%	+4.7%	-13%
SAP <sup>†</sup> [6]	ViT-B/32 <sup>512</sup>	72.2	84.6	91.6	93.9	83.7	94.0	97.8	99.3	94.8	97.7	93.5	96.8	92.1	96.3	78.1	85.7	91.0	94.8
RS@k <sup>†</sup>	ViT-B/32 <sup>512</sup>	75.9	87.1	93.1	95.1	85.1	94.6	98.0	99.3	95.1	97.7	94.1	96.7	93.2	96.5	78.1	86.4	92.3	95.6
SAP <sup>†</sup> [6]	ViT-B/16 <sup>512</sup>	79.1	89.0	94.2	95.8	86.6	95.4	98.4	99.5	95.5	97.7	94.2	96.9	93.1	96.6	86.2	92.1	95.1	97.2
RS@k <sup>†</sup>	ViT-B/16 <sup>512</sup>	83.9	92.1	95.9	97.2	88.0	96.1	98.6	99.6	96.2	98.0	95.2	97.2	94.7	97.1	89.5	94.2	96.6	98.3

Recall@k(%) performances. Best results are shown with bold, previous state-of-the-art with underline and relative gains over the state-of-the-art in % of error reduction with blue. All the methods marked with † were trained using the same pipeline by us.

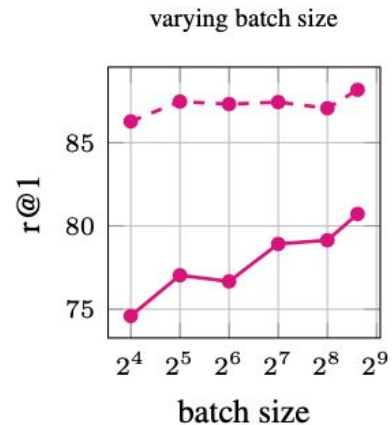
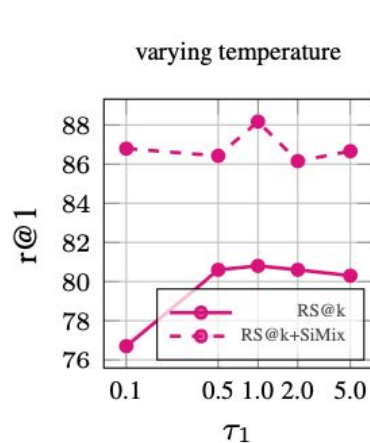
# Results - Revisited Oxford and Paris

Arch.	Loss	Train-set	Mean		$\mathcal{R}O$		$\mathcal{R}O+\mathcal{R}1M$		$\mathcal{R}Par$		$\mathcal{R}P+\mathcal{R}1M$		
			all	$\mathcal{R}1M$	med	hard	med	hard	med	hard	med	hard	
GeM*	AP [19]	Landmarks-clean [2] [14]	[44]/[56]	49.7	36.7	67.1	42.3	47.8	22.5	80.3	60.9	51.9	24.6
GeM*	AP [19]	GLDv1 [36]	[44]/github	-	-	66.3	42.5	-	-	80.2	60.8	-	-
GeM†	SAP [6]	GLDv1 [36]	[6]	52.7	40.6	67.9	46.3	49.5	25.8	81.7	63.3	57.4	29.8
GeM†	RS@k	GLDv1 [36]	ours	53.1	41.0	68.3	46.1	50.1	25.8	82.1	63.9	57.9	30.2
GeM+SiMix†	RS@k	GLDv1 [36]	ours	53.1	41.8	68.4	45.3	51.0	26.4	81.2	62.4	58.7	31.1

Performance comparison (mAP) on ROxford and RParis with 1m distractor images (R1m). Mean performance is reported across all setups or the large-scale setups only. \* denotes that the FC layer is not part of the training but is added afterward to implement whitening. Batch size is 4096 for all methods; SiMix virtually increases it to 10240. ResNet101 is used as a backbone for all methods.

# Results - Ablation

Method	r@1	r@2	r@4	r@8	r@16	Avg
RS@{1} <sup>†</sup>	81.1	87.7	92.0	95.0	96.9	90.5
RS@{1, 2} <sup>†</sup>	80.2	87.2	91.9	95.0	97.2	90.3
RS@{1, 2, 4} <sup>†</sup>	79.6	86.5	91.2	94.5	96.8	89.7
RS@{1, 2, 4, 8} <sup>†</sup>	79.3	86.3	91.0	94.5	96.9	89.6
RS@{1, 2, 4, 8, 16} <sup>†</sup>	80.8	87.6	92.2	95.0	97.1	90.5
RS@{2, 4, 8, 16} <sup>†</sup>	80.3	87.5	92.3	95.4	97.5	90.6
RS@{4, 8, 16} <sup>†</sup>	79.6	87.1	91.7	95.0	97.3	90.1
RS@{8, 16} <sup>†</sup>	79.6	87.1	91.7	95.0	97.3	90.1
RS@{16} <sup>†</sup>	75.8	83.9	89.8	93.6	96.4	87.9



# Thank You!

## For more details and applications, kindly refer to our papers:

1. Learning Surrogates via Deep Embedding, Y Patel, T Hodan, J Matas, *European Conference on Computer Vision (ECCV) 2020*.
2. Saliency Driven Perceptual Image Compression, Y Patel, S Appalaraju, R Manmatha, *Winter Applications of Computer Vision (WACV) 2021*.
3. FEDS--Filtered Edit Distance Surrogate, Y Patel, J Matas, *International Conference on Document Analysis and Recognition (ICDAR), 2021*.
4. Neural Network-based Acoustic Vehicle Counting, S Djukanović, Y Patel, J Matas, T Virtanen, *European Signal Processing Conference (EUSIPCO), 2021*.
5. Recall@k Surrogate Loss with Large Batches and Similarity Mixup, Y Patel, G Toliás, J Matas, *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022*.

My homepage: [yash0307.github.io/](https://yash0307.github.io/)