

# Classification is a Strong Baseline for Deep Metric Learning

Andrew Zhai, Hao-Yu Wu

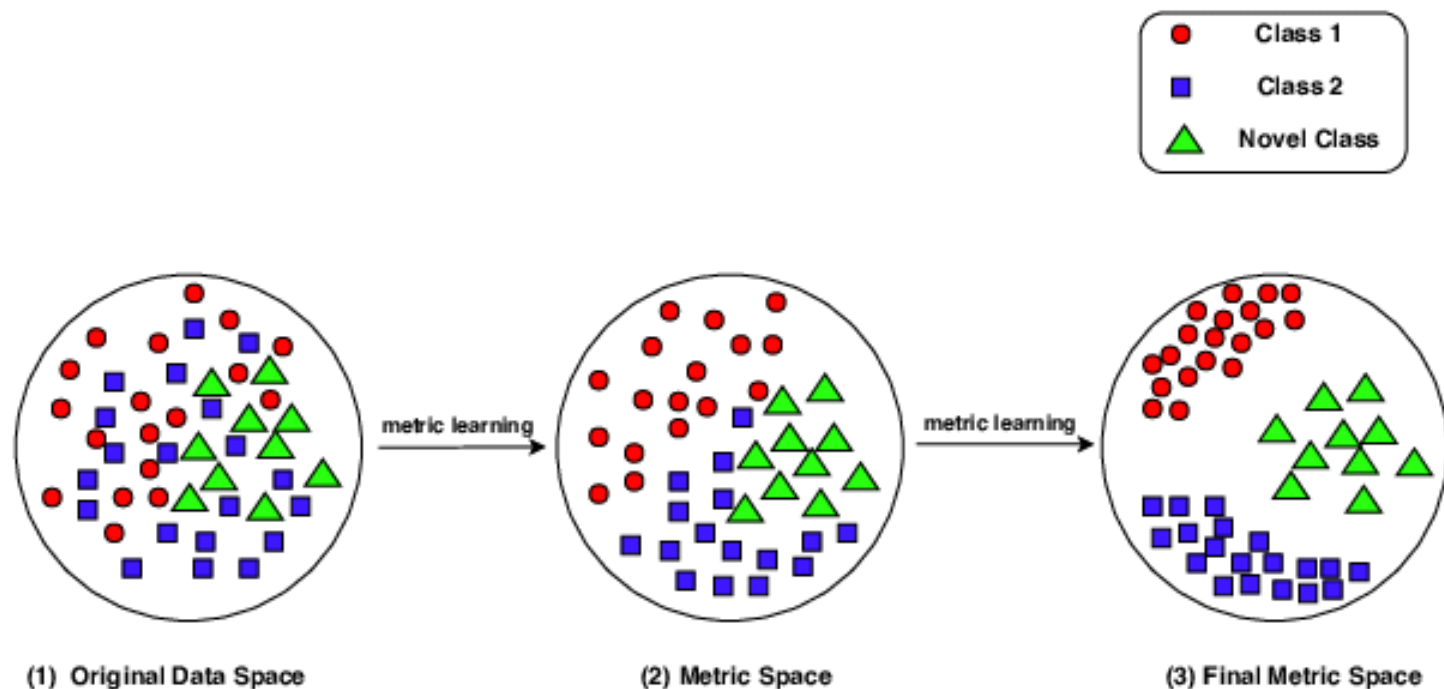
Presented by M.Rodin

November 17, 2020

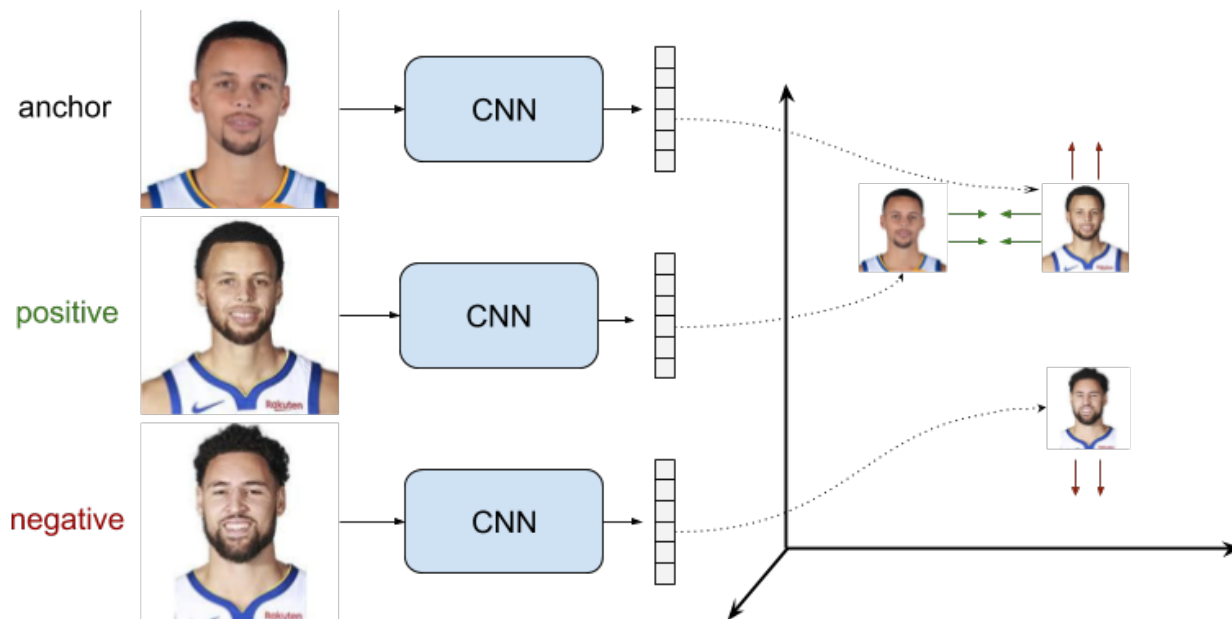
# Outline

- 1 Introduction
  - Problem formulation
  - Triplet loss
- 2 NormSoftmax
  - Architecture
  - LayerNorm
  - Binarization
  - Batch construction
  - Results
- 3 Conclusion
- 4 Practice

# Problem formulation

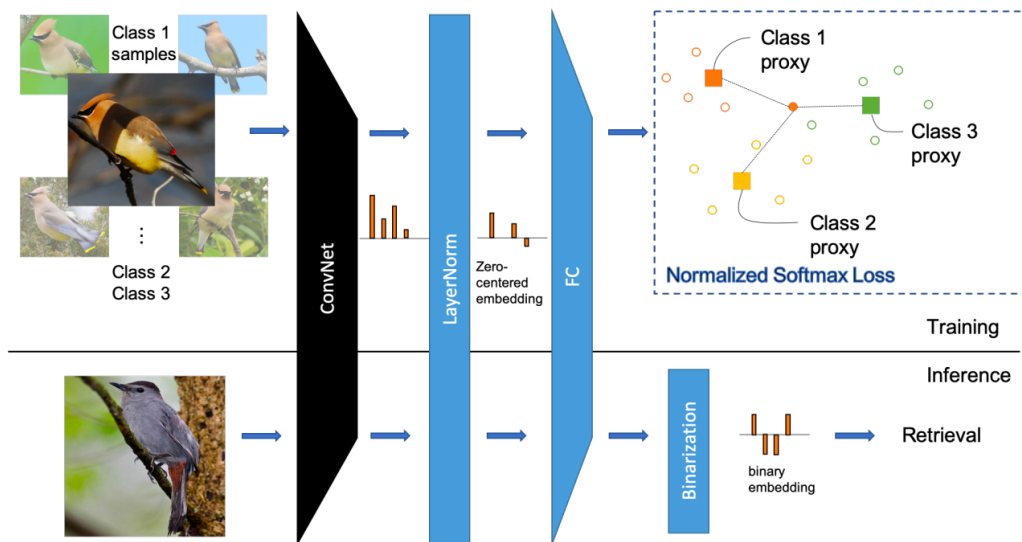


# Triplet loss



$$L(r_a, r_p, r_n) = \max(0, m + d(r_a, r_p) - d(r_a, r_n)) \quad (1)$$

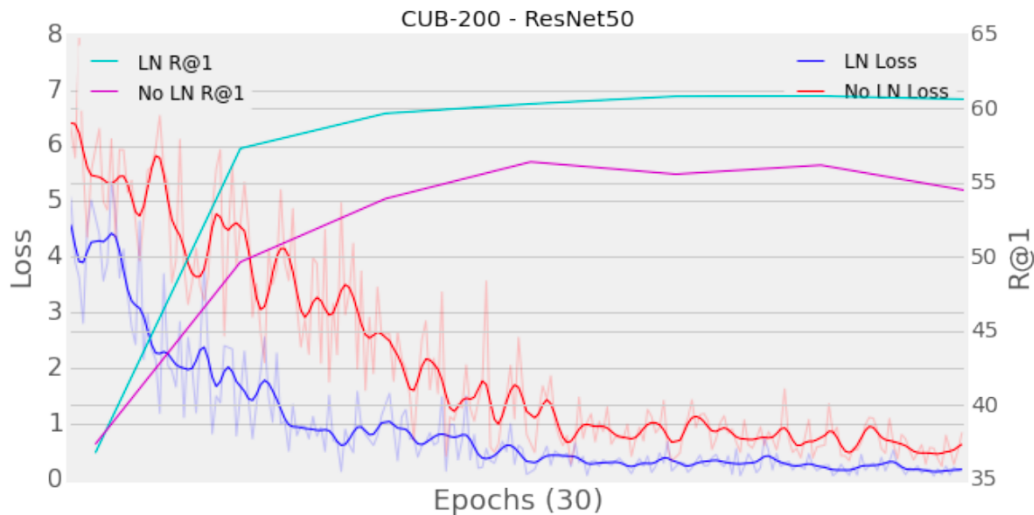
# NormSoftmax: Architecture



$$L_{\text{norm}} = -\log \left( \frac{\exp(x^T p_y / \sigma)}{\sum_{z \in Z} \exp(x^T p_z / \sigma)} \right) \quad (2)$$

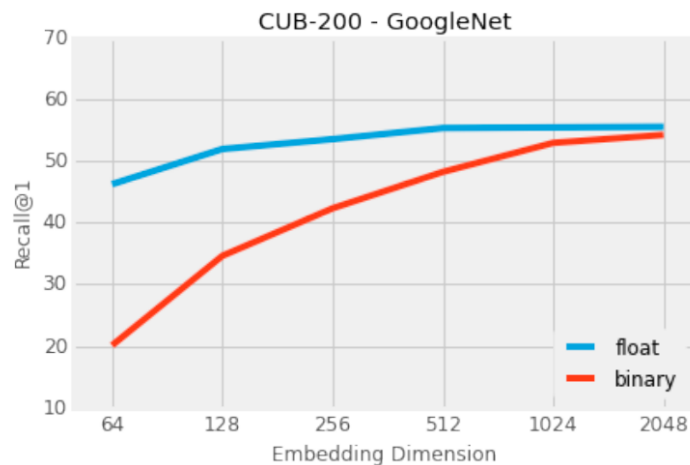
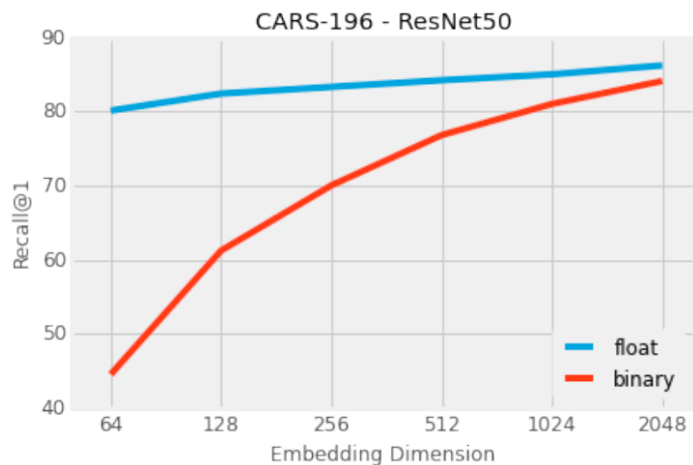
# NormSoftmax: LayerNorm

- Allows us to easily binarize embeddings via thresholding at zero
- Helps the network better initialize new parameters and reach better optima



# NormSoftmax: Binarization

- thresholding at zero
- 2048 binary vector = 64 float vector (256 bytes in memory)



# NormSoftmax: Batch construction

<b>S</b>	-	1	3	12	25	37	75
<b>C</b>	-	75	25	6	3	2	1
<b>R@1</b>	59.5	59.6	60.0	60.8	61.3	59.6	40.9

Table 4: ResNet50 Recall@1 on CUB-200-2011 dataset across varying samples per class for batch size of 75. **(S)** Samples per class in batch. **(C)** Distinct classes in batch. First column shows no class balancing in batch



# NormSoftmax: Results

Recall@K	Net.	CARS-196				CUB-200			
		1	2	4	8	1	2	4	8
Contrastive <sup>128</sup> [18]	G	21.7	32.3	46.1	58.9	26.4	37.7	49.8	62.3
Lifted Struct <sup>128</sup> [18]	G	49.0	60.3	72.1	81.5	47.2	58.9	70.2	80.2
Clustering <sup>64</sup> [19]	B	58.1	70.6	80.3	87.8	48.2	61.4	71.8	81.9
Npairs <sup>64</sup> [7]	G	71.1	79.7	86.5	91.6	51.0	63.3	74.3	83.2
Angular Loss <sup>512</sup> [9]	G	71.4	81.4	87.5	92.1	54.7	66.3	76.0	83.9
Proxy NCA <sup>64</sup> [7]	B	73.2	82.4	86.4	88.7	49.2	61.9	67.9	72.4
HDC <sup>384</sup> [30]	G	73.7	83.2	89.5	93.8	53.6	65.7	77.0	85.6
Margin <sup>128</sup> [7]	R50	79.6	86.5	91.9	95.1	<u>63.6</u>	74.4	83.1	90.0
HTL <sup>512</sup> [24]	B	81.4	88.0	92.7	95.7	57.1	68.8	78.7	86.5
A-BIER <sup>512</sup> [5]	G	82.0	89.0	93.2	96.1	57.5	68.7	78.3	86.2
ABE-8 <sup>512</sup> [26]	G†	85.2	90.5	94.0	96.1	60.6	71.5	79.8	87.4
DREML <sup>576</sup> [4]	R18	86.0	91.7	95.0	97.2	63.9	75.0	83.1	89.7
LMCL <sup>512</sup> [23]	R50	73.9	81.7	87.4	91.5	58.7	70.3	79.9	86.9
LMCL <sup>*2048</sup> [23]	R50	88.3	93.1	95.7	97.4	61.2	71.4	80.4	87.4
NormSoftMax <sup>1024</sup>	B	87.9	93.2	96.2	98.1	62.2	73.9	82.7	89.4
NormSoftmax <sup>128</sup>	R50	81.6	88.7	93.4	96.3	56.5	69.6	79.9	87.6
NormSoftmax <sup>512</sup>	R50	84.2	90.4	94.4	96.9	61.3	73.9	83.5	90.0
NormSoftmax <sup>2048</sup>	R50	<b>89.3</b>	<b>94.1</b>	<b>96.4</b>	<b>98.0</b>	<b>65.3</b>	<b>76.7</b>	<b>85.4</b>	<b>91.8</b>
NormSoftmax <sup>2048bits</sup>	R50	<u>88.7</u>	<u>93.7</u>	<b>96.4</b>	<b>98.0</b>	63.3	<u>75.2</u>	<u>84.3</u>	<u>91.0</u>

- Classification-based metric learning approaches can achieve state-of-the-art
- Binarization is allowing us to achieve SOTA performance with the same memory footprint as 64 dimensional float embeddings

# Practice

Tutorial in the similarity section

<https://github.com/microsoft/computervision-recipes/>

Reference:  
10.jpg



29.jpg  
rank: 1  
dist: 0.81



6.jpg  
rank: 2  
dist: 0.91



28.jpg  
rank: 3  
dist: 0.94



5.jpg  
rank: 4  
dist: 0.97



14.jpg  
rank: 5  
dist: 0.98



82.jpg  
rank: 6  
dist: 1.00



# Practice: Document stamps detection

jupyter check\_model\_distances Last Checkpoint: 10.06.2020 (autosaved)



Logout

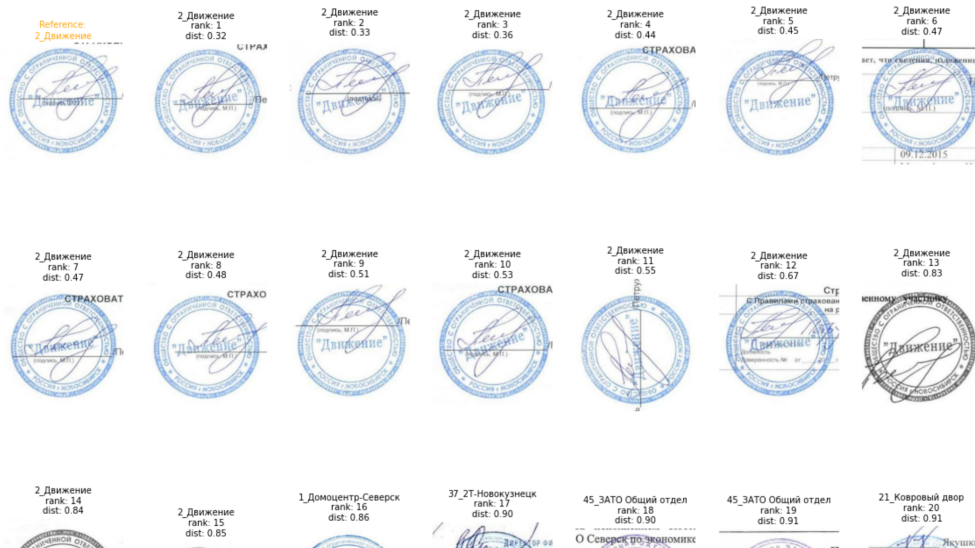
File Edit View Insert Cell Kernel Widgets Help

Trusted

catalyst

Code

```
dist = compute_distances(vectors[index], vectors, method='l2')  
plot_distances(dist, ds, num_rows=8, num_cols=7, figsize=(20,40))
```



Thank you for your attention