

# CLOOB: Modern Hopfield Networks with InfoLOOB Outperform CLIP



Seminar in AI  
Pascal Pilz  
Institute for Machine Learning

# Motivation and Setup

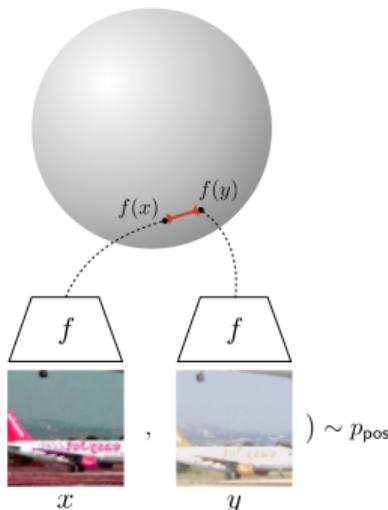


# Contrastive Learning and Zero-shot Transfer Learning

# Contrastive Learning and Zero-shot Transfer Learning

## Contrastive Learning:

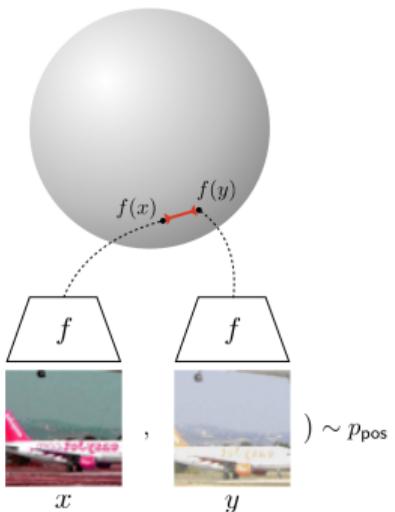
- self-supervised technique



# Contrastive Learning and Zero-shot Transfer Learning

## Contrastive Learning:

- self-supervised technique



## Zero-shot Transfer Learning:

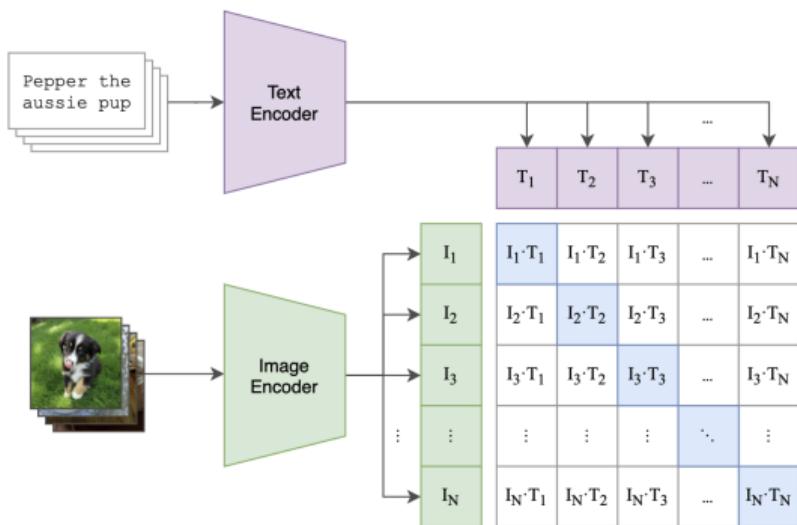
- model needs to generalize well



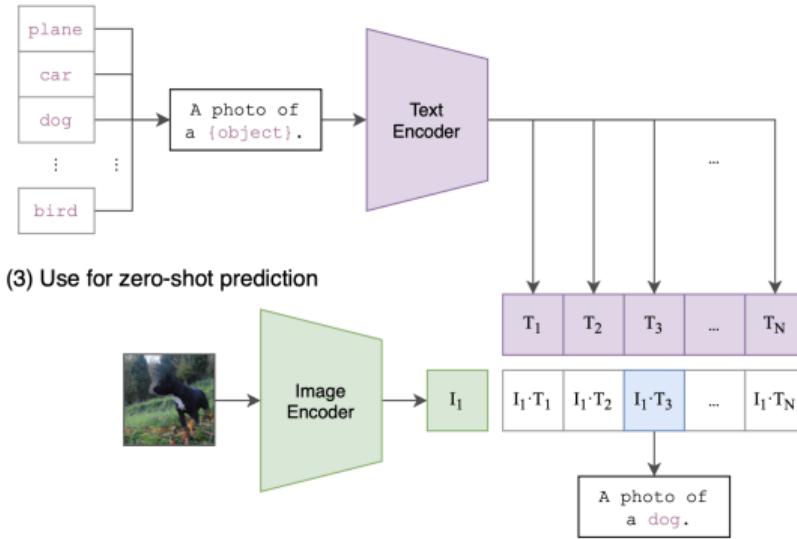
# CLIP (Contrastive Language-Image Pretraining) by OpenAI

# CLIP (Contrastive Language-Image Pretraining) by OpenAI

(1) Contrastive pre-training

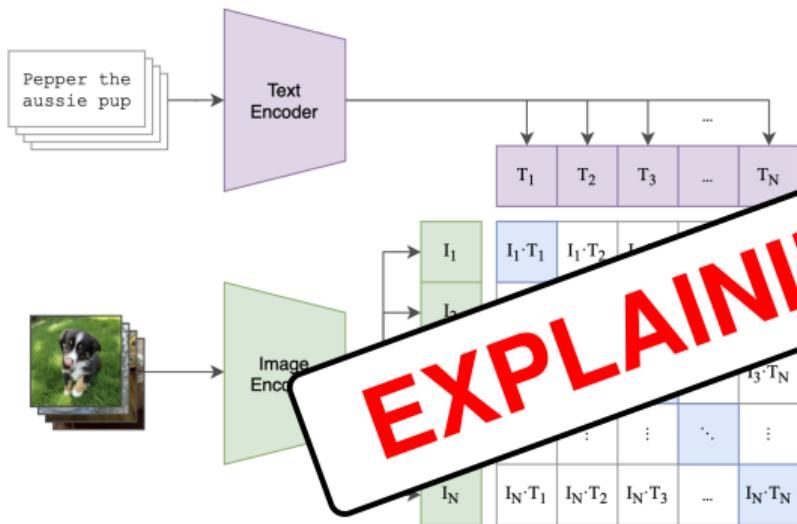


(2) Create dataset classifier from label text

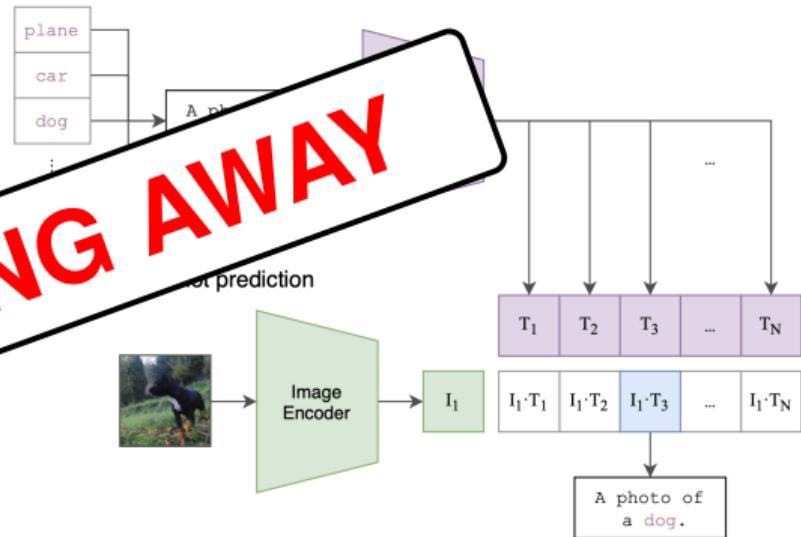


# CLIP (Contrastive Language-Image Pretraining) by OpenAI

(1) Contrastive pre-training



(2) Create dataset classifier from label text



# Explaining Away — Co-occurrences and Covariance structure

# Explaining Away — Co-occurrences and Covariance structure



# Explaining Away — Co-occurrences and Covariance structure



# Explaining Away — Co-occurrences and Covariance structure



# Modern Hopfield Networks



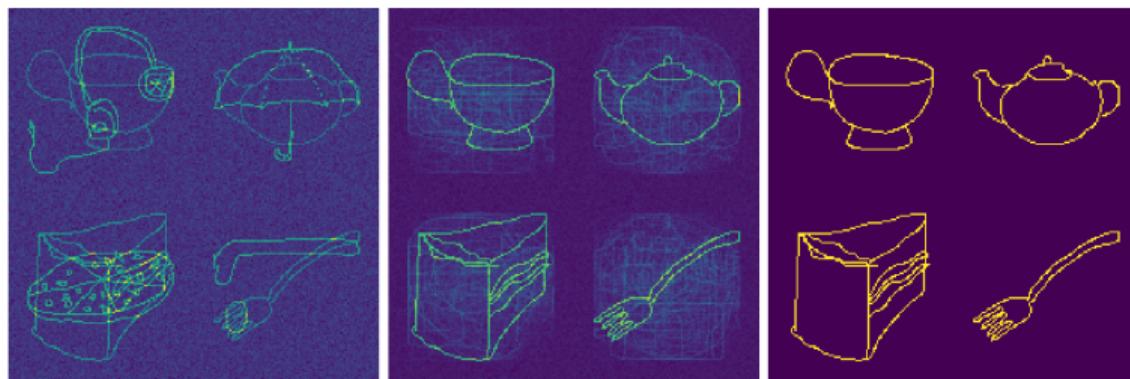
# Excursus: Modern Hopfield Networks

# Excursus: Modern Hopfield Networks

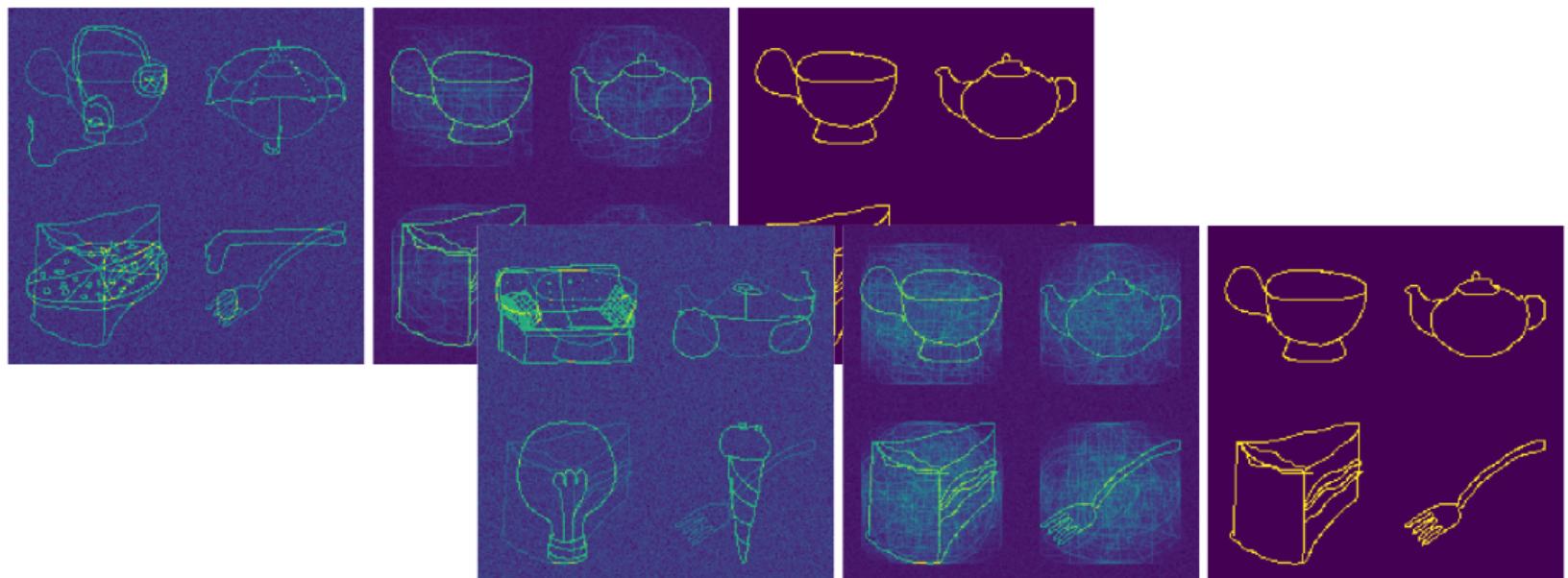


# Amplifying Co-occurrences and Covariance Structures

# Amplifying Co-occurrences and Covariance Structures

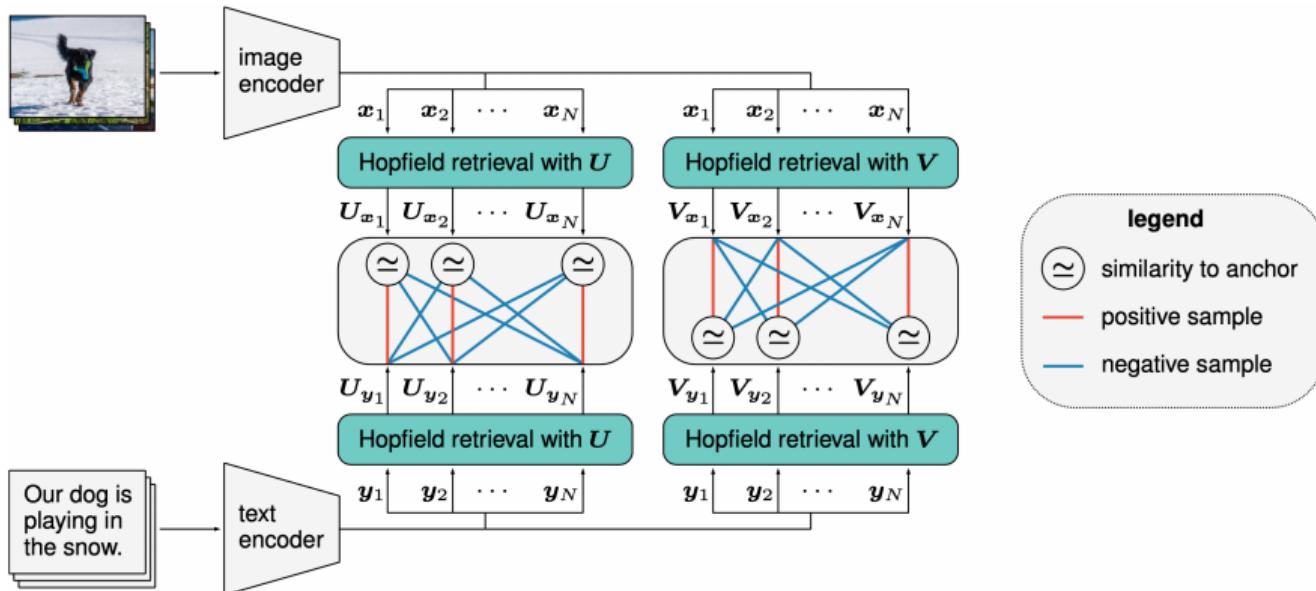


# Amplifying Co-occurrences and Covariance Structures

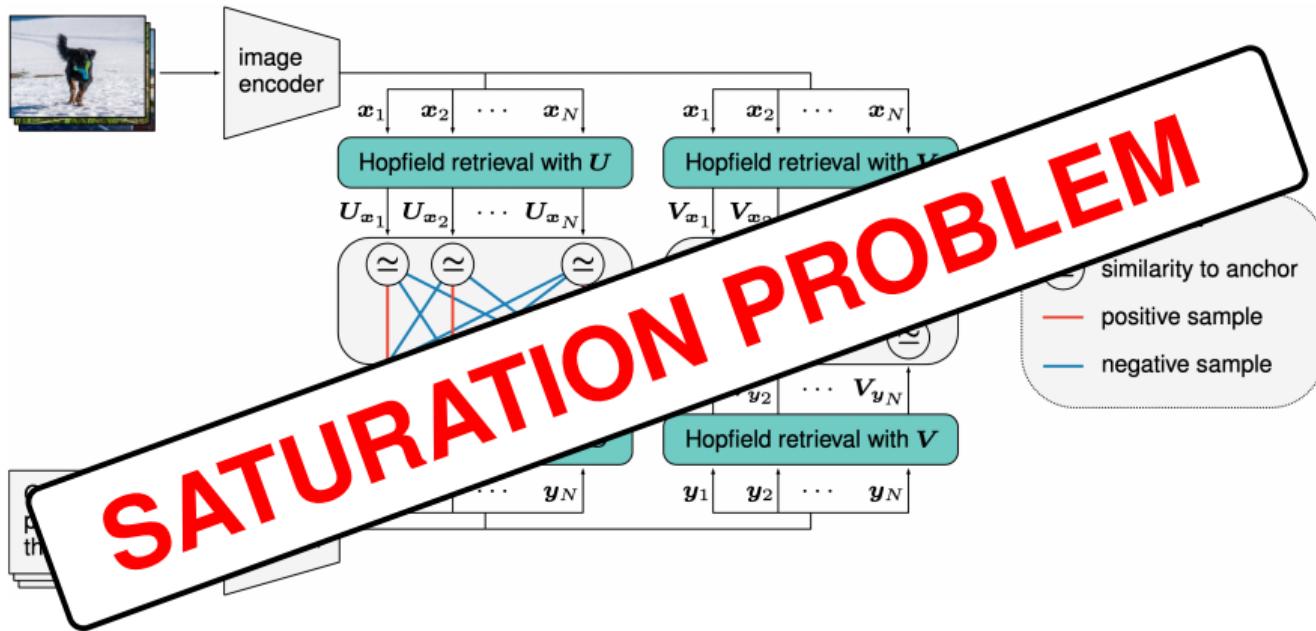


# CLOOB Architecture

# CLOOB Architecture



# CLOOB Architecture



# InfoNCE, InfoLOOB, and CLOOB



# InfoNCE (Noise Contrastive Estimation) — CLIP's Objective

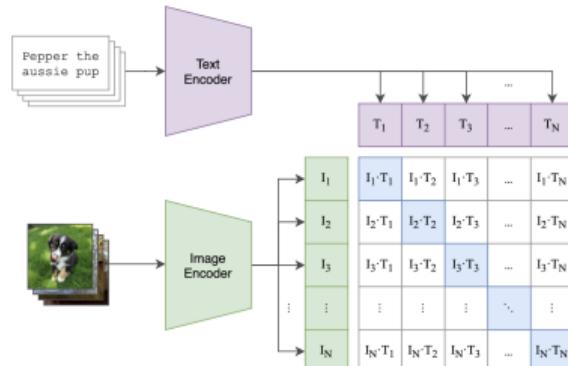
# InfoNCE (Noise Contrastive Estimation) — CLIP's Objective

$$L_{\text{InfoNCE}} = -\frac{1}{N} \ln \sum_{i=1}^N \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j=1}^N \exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_j)} - \frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j=1}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_i)}$$

# InfoNCE (Noise Contrastive Estimation) — CLIP's Objective

$$L_{\text{InfoNCE}} = -\frac{1}{N} \ln \sum_{i=1}^N \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j=1}^N \exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_j)} - \frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j=1}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_i)}$$

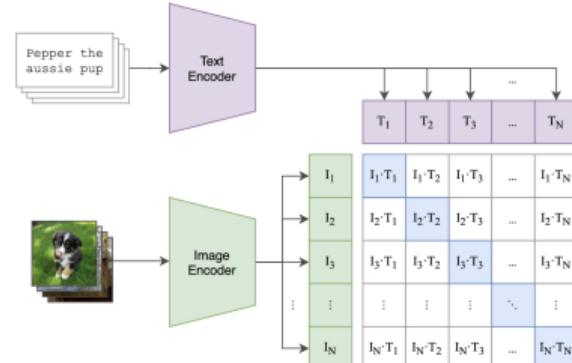
- $\mathbf{x}_i$  — image embedding
- $\mathbf{y}_i$  — text embedding



# InfoNCE (Noise Contrastive Estimation) — CLIP's Objective

$$L_{\text{InfoNCE}} = -\frac{1}{N} \ln \sum_{i=1}^N \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j=1}^N \exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_j)} - \frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j=1}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_i)}$$

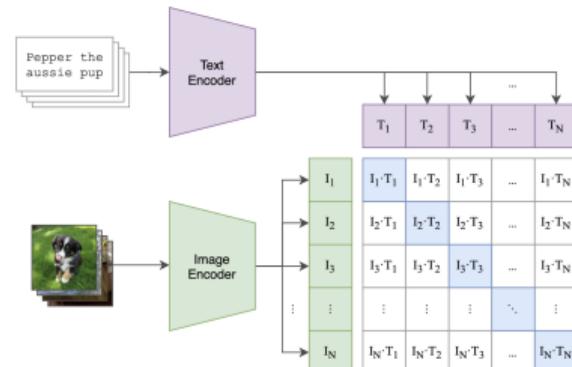
- $\mathbf{x}_i$  — image embedding
- $\mathbf{y}_i$  — text embedding
- $\tau$  — temperature, inverse entropy



# InfoNCE (Noise Contrastive Estimation) — CLIP's Objective

$$L_{\text{InfoNCE}} = -\frac{1}{N} \ln \sum_{i=1}^N \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j=1}^N \exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_j)} - \frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j=1}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_i)}$$

- $\mathbf{x}_i$  — image embedding
- $\mathbf{y}_i$  — text embedding
- $\tau$  — temperature, inverse entropy
- $\|\mathbf{x}_i\| = \|\mathbf{y}_j\| = 1 \Rightarrow \cos(\mathbf{x}_i, \mathbf{y}_j) = \mathbf{x}_i^T \mathbf{y}_j$



# InfoLOOB (Leave One Out Bound) — CLOOB's Objective

$$\mathcal{L}_{\text{InfoLOOB}} = -\frac{1}{N} \ln \sum_{i=1}^N \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j \neq i}^N \exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_j)} - \frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j \neq i}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_i)}$$

# InfoLOOB (Leave One Out Bound) — CLOOB's Objective

$$L_{\text{InfoLOOB}} = -\frac{1}{N} \ln \sum_{i=1}^N \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j \neq i}^N \exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_j)} - \frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j \neq i}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_i)}$$

$$L_{\text{InfoNCE}}(\mathbf{y}_1) = -\ln \frac{\overbrace{\exp(\tau^{-1} \mathbf{x}_1^T \mathbf{y}_1)}^a}{\underbrace{\exp(\tau^{-1} \mathbf{x}_1^T \mathbf{y}_1)}_a + \underbrace{\sum_{j=2}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_1)}_b}$$

# InfoLOOB (Leave One Out Bound) — CLOOB's Objective

$$L_{\text{InfoLOOB}} = -\frac{1}{N} \ln \sum_{i=1}^N \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j \neq i}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_j)} - \frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j \neq i}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_i)}$$

$$L_{\text{InfoNCE}}(\mathbf{y}_1) = -\ln \frac{\overbrace{\exp(\tau^{-1} \mathbf{x}_1^T \mathbf{y}_1)}^a}{\underbrace{\exp(\tau^{-1} \mathbf{x}_1^T \mathbf{y}_1) + \sum_{j=2}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_1)}_b} \quad \leftarrow \text{saturation problem}$$

# InfoLOOB (Leave One Out Bound) — CLOOB's Objective

$$L_{\text{InfoLOOB}} = -\frac{1}{N} \ln \sum_{i=1}^N \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j \neq i}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_j)} - \frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{x}_i^T \mathbf{y}_i)}{\sum_{j \neq i}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_i)}$$

$$L_{\text{InfoNCE}}(\mathbf{y}_1) = -\ln \frac{\overbrace{\exp(\tau^{-1} \mathbf{x}_1^T \mathbf{y}_1)}^a}{\underbrace{\exp(\tau^{-1} \mathbf{x}_1^T \mathbf{y}_1) + \sum_{j=2}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_1)}_b} \quad \leftarrow \text{saturation problem}$$

$$L_{\text{InfoLOOB}}(\mathbf{y}_1) = -\ln \frac{\overbrace{\exp(\tau^{-1} \mathbf{x}_1^T \mathbf{y}_1)}^a}{\underbrace{\sum_{j=2}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y}_1)}_b}$$

# CLOOB (Contrastive Leave One Out Boost)

1	20
2	21
3	22
4	23
5	24
6	25
7	26
8	27
9	28
10	29
11	30
12	31
13	32
14	33
15	34
16	35
17	36
18	37
19	

# CLOOB (Contrastive Leave One Out Boost)

1	# <i>image_encoder – ResNet</i>	20
2	# <i>text_encoder – Text Transformer</i>	21
3		22
4	# $I[n, h, w, c]$ – minibatch of images	23
5	# $T[n, l]$ – minibatch of texts	24
6		25
7	# $W_i[d_i, d_e]$ – image projection	26
8	# $W_t[d_t, d_e]$ – text projection	27
9		28
10	# $\beta$ – inverse temperature Hopfield retrieval	29
11	# $\tau$ – temperature InfoLOOB	30
12		31
13		32
14		33
15		34
16		35
17		36
18		37
19		

# CLOOB (Contrastive Leave One Out Boost)

```
1 # image_encoder – ResNet                                20
2 # text_encoder – Text Transformer                      21
3                                         22
4 # I[n, h, w, c] – minibatch of images                23
5 # T[n, l] – minibatch of texts                         24
6                                         25
7 # W_i[d_i, d_e] – image projection                  26
8 # W_t[d_t, d_e] – text projection                   27
9                                         28
10 # beta – inverse temperature Hopfield retrieval    29
11 # tau – temperature InfoLOOB                        30
12                                         31
13 # extract feature representations                 32
14 I_f = image_encoder(I) #[n, d_i]                  33
15 T_f = text_encoder(T) #[n, d_t]                   34
16                                         35
17                                         36
18                                         37
19
```

# CLOOB (Contrastive Leave One Out Boost)

```
1 # image_encoder – ResNet                                20
2 # text_encoder – Text Transformer                      21
3                                         22
4 # I[n, h, w, c] – minibatch of images                23
5 # T[n, l] – minibatch of texts                         24
6                                         25
7 # W_i[d_i, d_e] – image projection                  26
8 # W_t[d_t, d_e] – text projection                   27
9                                         28
10 # beta – inverse temperature Hopfield retrieval     29
11 # tau – temperature InfoLOOB                        30
12                                         31
13 # extract feature representations                 32
14 I_f = image_encoder(I) #[n, d_i]                   33
15 T_f = text_encoder(T) #[n, d_t]                     34
16                                         35
17 # joint multimodal embedding                       36
18 x = l2_normalize(I_f @ W_i) #[n, d_e]             37
19 y = l2_normalize(T_f @ W_t) #[n, d_e]
```

# CLOOB (Contrastive Leave One Out Boost)

```
1 # image_encoder – ResNet
2 # text_encoder – Text Transformer
3
4 # I[n, h, w, c] – minibatch of images
5 # T[n, l] – minibatch of texts
6
7 # W_i[d_i, d_e] – image projection
8 # W_t[d_t, d_e] – text projection
9
10 # beta – inverse temperature Hopfield retrieval
11 # tau – temperature InfoLOOB
12
13 # extract feature representations
14 I_f = image_encoder(I) #[n, d_i]
15 T_f = text_encoder(T) #[n, d_t]
16
17 # joint multimodal embedding
18 x = l2_normalize(I_f @ W_i) #[n, d_e]
19 y = l2_normalize(T_f @ W_t) #[n, d_e]
20
21 # Hopfield retrieval H with batch stored
22 # H(beta, A, B) = B.T @ softmax(beta * A @ B.T)
23 U_x = H(beta, x, x).T #[n, d_e]
24 U_y = H(beta, y, x).T #[n, d_e]
25 V_x = H(beta, x, y).T #[n, d_e]
26
27
28
29
30
31
32
33
34
35
36
37
```

# CLOOB (Contrastive Leave One Out Boost)

```
1  # image_encoder – ResNet
2  # text_encoder – Text Transformer
3
4  # I[n, h, w, c] – minibatch of images
5  # T[n, l] – minibatch of texts
6
7  # W_i[d_i, d_e] – image projection
8  # W_t[d_t, d_e] – text projection
9
10 # beta – inverse temperature Hopfield retrieval
11 # tau – temperature InfoLOOB
12
13 # extract feature representations
14 I_f = image_encoder(I) #[n, d_i]
15 T_f = text_encoder(T) #[n, d_t]
16
17 # joint multimodal embedding
18 x = l2_normalize(I_f @ W_i) #[n, d_e]
19 y = l2_normalize(T_f @ W_t) #[n, d_e]
20
21 # Hopfield retrieval H with batch stored
22 U_x = H(beta, x, x).T #[n, d_e]
23 U_y = H(beta, y, x).T #[n, d_e]
24 V_x = H(beta, x, y).T #[n, d_e]
25 V_y = H(beta, y, y).T #[n, d_e]
26
27 # normalize retrievals
28 U_x = l2_normalize(U_x) #[n, d_e]
29 U_y = l2_normalize(U_y) #[n, d_e]
30 V_x = l2_normalize(V_x) #[n, d_e]
31 V_y = l2_normalize(V_y) #[n, d_e]
32
33
34
35
36
37
```

# CLOOB (Contrastive Leave One Out Boost)

```
1  # image_encoder – ResNet
2  # text_encoder – Text Transformer
3
4  # I[n, h, w, c] – minibatch of images
5  # T[n, l] – minibatch of texts
6
7  # W_i[d_i, d_e] – image projection
8  # W_t[d_t, d_e] – text projection
9
10 # beta – inverse temperature Hopfield retrieval
11 # tau – temperature InfoLOOB
12
13 # extract feature representations
14 I_f = image_encoder(I) #[n, d_i]
15 T_f = text_encoder(T) #[n, d_t]
16
17 # joint multimodal embedding
18 x = l2_normalize(I_f @ W_i) #[n, d_e]
19 y = l2_normalize(T_f @ W_t) #[n, d_e]
20
21 # Hopfield retrieval H with batch stored
22 # H(beta, A, B) = B.T @ softmax(beta * A @ B.T)
23 U_x = H(beta, x, x).T #[n, d_e]
24 U_y = H(beta, y, x).T #[n, d_e]
25 V_x = H(beta, x, y).T #[n, d_e]
26 V_y = H(beta, y, y).T #[n, d_e]
27
28 # normalize retrievals
29 U_x = l2_normalize(U_x) #[n, d_e]
30 U_y = l2_normalize(U_y) #[n, d_e]
31 V_x = l2_normalize(V_x) #[n, d_e]
32 V_y = l2_normalize(V_y) #[n, d_e]
33
34 # loss: info_loob(tau, anchors, samples)
35 # samples contain pos. and neg. embeddings
36 loss_i = info_loob(tau, U_x, U_y)
37 loss_t = info_loob(tau, V_y, V_x)
38 loss = (loss_i + loss_t) * tau
```

# Experiments



# Experiments

# Experiments

- Conceptual Captions (CC) Pretraining
  - rich textual description
  - only 2.9 million images

# Experiments

- Conceptual Captions (CC) Pretraining
  - rich textual description
  - only 2.9 million images
- Yahoo Flickr Create Commons (YFCC) Pretraining
  - less rich textual description
  - 15 million images

# Experiments

- Conceptual Captions (CC) Pretraining
  - rich textual description
  - only 2.9 million images
- Yahoo Flickr Create Commons (YFCC) Pretraining
  - less rich textual description
  - 15 million images
- Both tested on seven image classification datasets

# Experiments

- Conceptual Captions (CC) Pretraining
  - rich textual description
  - only 2.9 million images
- Yahoo Flickr Create Commons (YFCC) Pretraining
  - less rich textual description
  - 15 million images
- Both tested on seven image classification datasets
- Ablation studies
  - Modern Hopfield networks
  - InfoLOOB

# CLIP and CLOOB: Results

CC — mean accuracy over 5 runs

Dataset	CLIP RN-50	CLOOB RN-50	CLIP* RN-50	CLOOB* RN-50
Birdsnap	$2.26 \pm 0.20$	<b><math>3.06 \pm 0.30</math></b>	$2.8 \pm 0.16$	<b><math>3.24 \pm 0.31</math></b>
Country211	$0.67 \pm 0.11$	$0.67 \pm 0.05$	$0.7 \pm 0.04$	$0.73 \pm 0.05$
Flowers102	$12.56 \pm 0.38$	$13.45 \pm 1.19$	$13.32 \pm 0.43$	$14.36 \pm 1.17$
GTSRB	$7.66 \pm 1.07$	$6.38 \pm 2.11$	$8.96 \pm 1.70$	$7.03 \pm 1.22$
UCF101	$20.98 \pm 1.55$	$22.26 \pm 0.72$	$21.63 \pm 0.65$	<b><math>23.03 \pm 0.85</math></b>
Stanford Cars	$0.91 \pm 0.10$	<b><math>1.23 \pm 0.10</math></b>	$0.99 \pm 0.16$	<b><math>1.41 \pm 0.32</math></b>
ImageNet	$20.33 \pm 0.28$	<b><math>23.97 \pm 0.15</math></b>	$21.3 \pm 0.42$	<b><math>25.67 \pm 0.22</math></b>
ImageNet V2	$20.24 \pm 0.50$	<b><math>23.59 \pm 0.15</math></b>	$21.24 \pm 0.22$	<b><math>25.49 \pm 0.11</math></b>

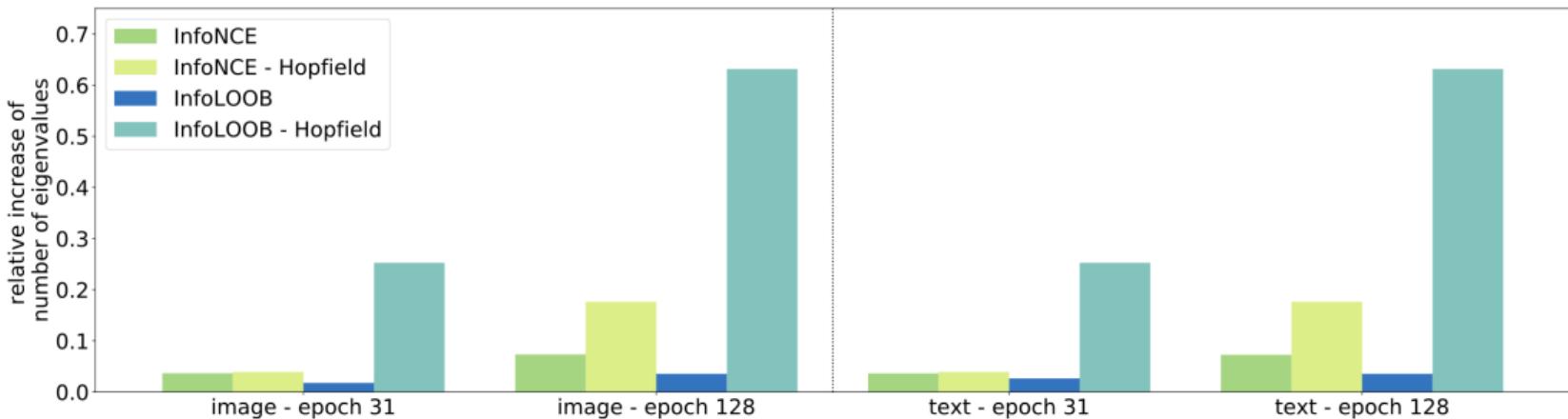
Bold: statistically significant

YFCC — one run

Dataset	RN-50		RN-101		RN-50x4	
	CLIP	CLOOB	CLIP	CLOOB	CLIP	CLOOB
Birdsnap	21.8	<b>28.9</b>	22.6	<b>30.3</b>	20.8	<b>32.0</b>
Country211	6.9	<b>7.9</b>	7.8	<b>8.5</b>	8.1	<b>9.3</b>
Flowers102	48.0	<b>55.1</b>	48.0	<b>55.3</b>	50.1	<b>54.3</b>
GTSRB	7.9	<b>8.1</b>	7.4	<b>11.6</b>	9.4	<b>11.8</b>
UCF101	<b>27.2</b>	25.3	28.6	<b>28.8</b>	31.0	<b>31.9</b>
Stanford Cars	3.7	<b>4.1</b>	3.8	<b>5.5</b>	3.5	<b>6.1</b>
ImageNet	34.6	<b>35.7</b>	35.3	<b>37.1</b>	37.7	<b>39.0</b>
ImageNet V2	33.4	<b>34.6</b>	34.1	<b>35.6</b>	35.9	<b>37.3</b>

Bold: higher value

# Results Ablation Study



# Conclusion



# Critical Assessment

# Critical Assessment

- Reproducibility
  - hyperparameters, experiments well-defined ✓
  - code, datasets public ✓

# Critical Assessment

- Reproducibility
  - hyperparameters, experiments well-defined ✓
  - code, datasets public ✓
- Thoroughness
  - two theorems, both proven ✓
  - all assumptions listed ✓

# Critical Assessment

- Reproducibility
  - hyperparameters, experiments well-defined ✓
  - code, datasets public ✓
- Thoroughness
  - two theorems, both proven ✓
  - all assumptions listed ✓
- Fairness of comparison
  - as faithful to original CLIP as possible ✓
  - CLIP's 400 million image dataset not public

# Critical Assessment

- Reproducibility
  - hyperparameters, experiments well-defined ✓
  - code, datasets public ✓
- Thoroughness
  - two theorems, both proven ✓
  - all assumptions listed ✓
- Fairness of comparison
  - as faithful to original CLIP as possible ✓
  - CLIP's 400 million image dataset not public

Paper fulfills all NeurIPS check marks

# Summary

→ CLIP

# Summary

→ CLIP

✗ explaining away

# Summary

- CLIP
  - ✗ explaining away
- Modern Hopfield networks

# Summary

→ CLIP

- ✗ explaining away

→ Modern Hopfield networks

- ✗ saturation problem

# Summary

→ CLIP

- ✗ explaining away

→ Modern Hopfield networks

- ✗ saturation problem

→ CLOOB

# Summary

→ CLIP

✗ explaining away

→ Modern Hopfield networks

✗ saturation problem

→ CLOOB

✓ Modern Hopfield networks

# Summary

→ CLIP

✗ explaining away

→ Modern Hopfield networks

✗ saturation problem

→ CLOOB

✓ Modern Hopfield networks

✓ InfoLOOB

# Summary

→ CLIP

- ✗ explaining away

→ Modern Hopfield networks

- ✗ saturation problem

→ CLOOB

- ✓ Modern Hopfield networks
- ✓ InfoLOOB

→ 😊

# Summary

→ CLIP

✗ explaining away

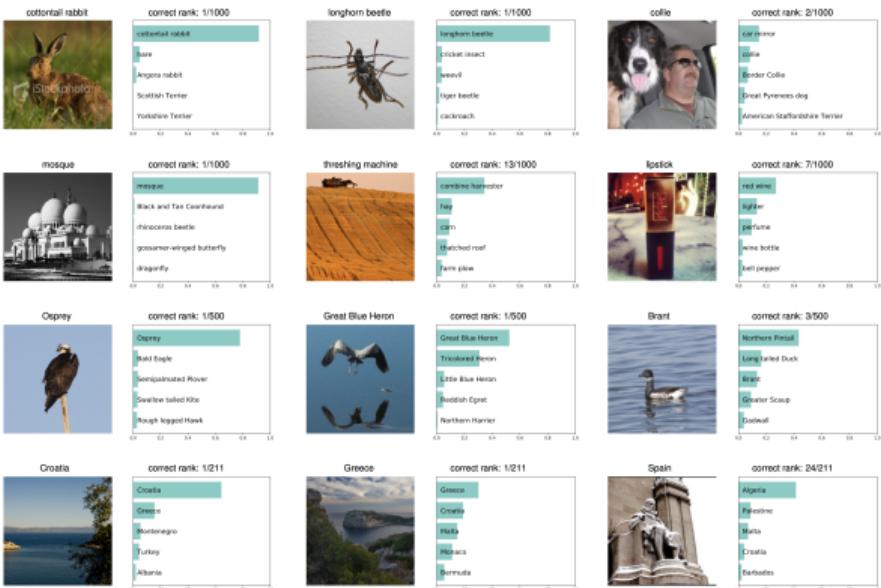
→ Modern Hopfield networks

✗ saturation problem

→ CLOOB

- ✓ Modern Hopfield networks
- ✓ InfoLOOB

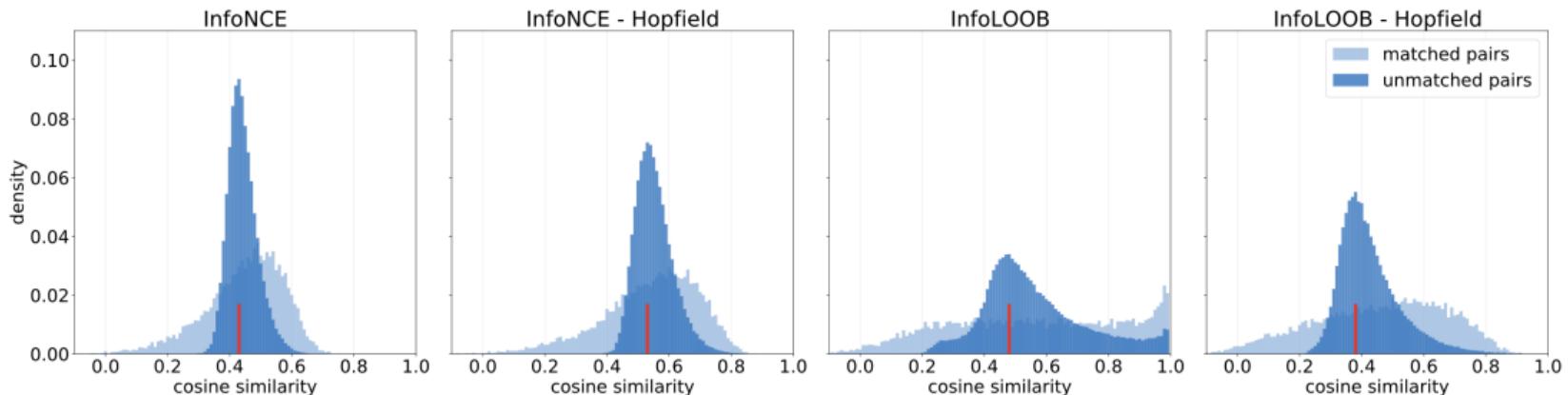
→ 😊



# Additional Material

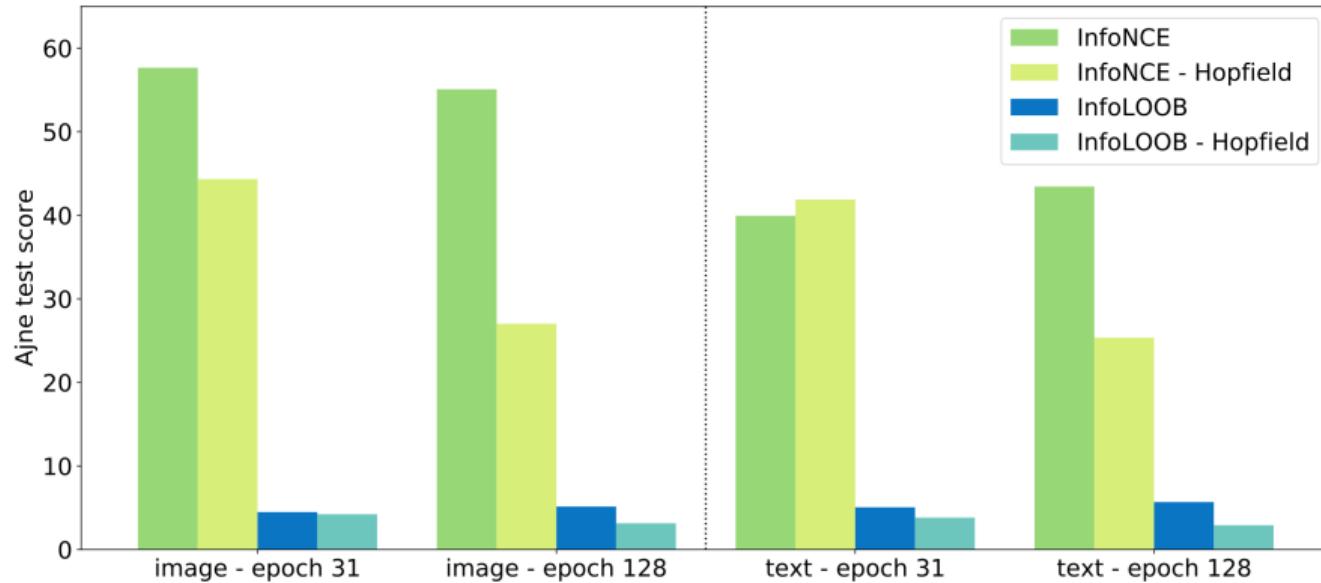


# Distribution of Cosine Similarity



Distribution of cosine similarity of matched pairs and of the 10 unmatched pairs that have the highest similarity score with anchor

# Ajne Test



# InfoLOOB – Further details

	$\mathbf{u}_i$ - stored image emb.	$\mathbf{v}_i$ - stored text emb.
$\mathbf{x}_i$ - image query emb.	$\mathbf{U}_{\mathbf{x}_i}$ - image-retr. image emb.	$\mathbf{V}_{\mathbf{x}_i}$ - image-retr. text emb.
$\mathbf{y}_i$ - text query emb.	$\mathbf{U}_{\mathbf{y}_i}$ - text-retr. image emb.	$\mathbf{V}_{\mathbf{y}_i}$ - text-retr. text emb.

$$L_{\text{InfoLOOB}} = -\frac{1}{N} \ln \sum_{i=1}^N \frac{\exp(\tau^{-1} \mathbf{U}_{\mathbf{x}_i}^T \mathbf{U}_{\mathbf{y}_i})}{\sum_{j \neq i}^N \exp(\tau^{-1} \mathbf{U}_{\mathbf{x}_i}^T \mathbf{U}_{\mathbf{y}_j})} - \frac{1}{N} \sum_{i=1}^N \ln \frac{\exp(\tau^{-1} \mathbf{V}_{\mathbf{x}_i}^T \mathbf{V}_{\mathbf{y}_i})}{\sum_{j \neq i}^N \exp(\tau^{-1} \mathbf{V}_{\mathbf{x}_j}^T \mathbf{V}_{\mathbf{y}_i})}$$

# Gradients of InfoLOOB and InfoNCE

$$\frac{\partial}{\partial \mathbf{y}} L_{\text{InfoNCE}}(\mathbf{y}) = \frac{\partial}{\partial \mathbf{y}} - \ln \frac{\exp(\tau^{-1} \mathbf{x}_1^T \mathbf{y}_1)}{\sum_{j=1}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y})} = -\tau^{-1} \mathbf{y}^T \mathbf{x}_1 + \tau^{-1} \text{lse}(\tau^{-1}, \mathbf{X}^T \mathbf{y})$$

$$\frac{\partial}{\partial \mathbf{y}} L_{\text{InfoLOOB}}(\mathbf{y}) = \frac{\partial}{\partial \mathbf{y}} - \ln \frac{\exp(\tau^{-1} \mathbf{x}_1^T \mathbf{y}_1)}{\sum_{j \neq 1}^N \exp(\tau^{-1} \mathbf{x}_j^T \mathbf{y})} = -\tau^{-1} \mathbf{y}^T \mathbf{x}_1 + \tau^{-1} \text{lse}(\tau^{-1}, \tilde{\mathbf{X}}^T \mathbf{y})$$

$$\frac{\partial}{\partial \mathbf{y}} L_{\text{InfoNCE}}(\mathbf{y}) = -\tau^{-1} (1 - p_1) (\mathbf{x}_1 - \tilde{\mathbf{X}} \text{softmax}(\tau^{-1} \tilde{\mathbf{X}}^T \mathbf{y})) = (1 - p_1) \frac{\partial}{\partial \mathbf{y}} L_{\text{InfoLOOB}}(\mathbf{y})$$

$$\text{lse}(\beta, \alpha) = \beta^{-1} \log \left( \sum_{i=1}^N \exp(\beta \alpha_i) \right)$$

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N), \quad \tilde{\mathbf{X}} = (\mathbf{x}_2, \dots, \mathbf{x}_N)$$