

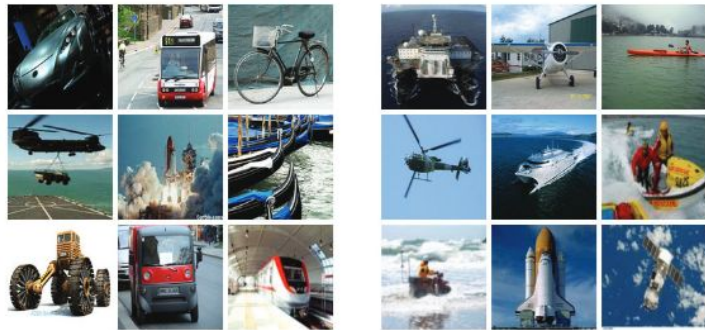
# Representation Learning via Non-Contrastive Mutual Information

Zhaohan Daniel Guo  
Bernardo Avila Pires  
Khimya Khetarpal  
Dale Schuurmans  
Bo Dai

**Poster #169**

# Self-Supervised Representation Learning

Lots of unlabeled data?

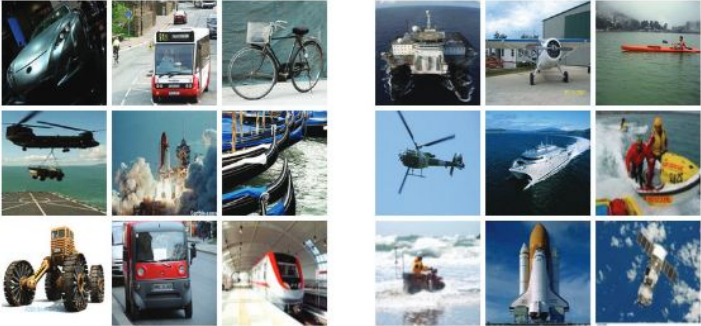


Labeling too expensive?



# Self-Supervised Representation Learning

Lots of unlabeled data?



Labeling too expensive?



Self-Supervised Representation Learning!

# Contrastive vs. Non-Contrastive

## Contrastive

eg) SimCLR, Spectral Contrastive, f-MICL

### PROS

#### Theoretically Sound

Maximize Mutual Information

### CONS

#### High Variance Loss

Pairwise comparison of all data points

## Non-Contrastive

eg) BYOL, SimSiam

### PROS

#### Batch Size Efficient

### CONS

#### Potential Representation Collapse

Theory-Practice Gap

# Contrastive vs. Non-Contrastive

## Contrastive

eg) SimCLR, Spectral Contrastive, f-MICL

### PROS

#### Theoretically Sound

Maximize Mutual Information

### CONS

#### High Variance Loss

Pairwise comparison of all data points

## Non-Contrastive

eg) BYOL, SimSiam

### PROS

#### Batch Size Efficient

### CONS

#### Potential Representation Collapse

Theory-Practice Gap

**MINC Loss - Best of Both Worlds!**

# MINC Loss - Step 1

## Start with Spectral Contrastive Loss

Theoretically maximizes a certain mutual information.

$$\max_{\phi} \underbrace{2\mathbb{E}_{p(x,x')} [\phi(x)^\top \phi(x')]}_{\text{Align embeddings for similar data points.}} - \underbrace{\mathbb{E}_{p(x)p(x')} [(\phi(x)^\top \phi(x'))^2]}_{\text{Make dissimilar embeddings orthogonal.}} \tag{3}$$

Align embeddings for similar data points.

Make dissimilar embeddings orthogonal

# MINC Loss - Step 2

## Refactor with Auxiliary

Avoids double expectation.

$$\max_{\phi} 2\mathbb{E}_{p(x,x')} [\phi(x)^\top \phi(x')] - \mathbb{E}_{p(x)p(x')} [(\phi(x)^\top \phi(x'))^2] \quad (3)$$



$$\max_{\phi} 2\mathbb{E}_{p(x,x')} [\phi(x)^\top \phi(x')] - \mathbb{E}_{p(x')} [\phi(x')^\top \Lambda \phi(x')] \\ \text{subject to } \Lambda = \mathbb{E}_{p(x)} [\phi(x) \phi^\top(x)], \quad (5)$$

Given Lambda (d x d matrix), this becomes a non-contrastive loss!

# MINC Loss - Step 3

Use EMA and Generalized Hebbian Algorithm (GHA)

EMA Update:  $\Lambda_{t+1} \leftarrow (1 - \alpha)\Lambda_t + \alpha(\phi(x)\phi(x)^\top)$

$$\max_{\phi} 2\mathbb{E}[\phi(x)^\top \phi(x)] + \mathbb{E}[\phi(x)^\top \underbrace{\mathbf{LT}(\Lambda)}_{\text{Lower Triangular Transformation}} \phi(x)]$$

Lower Triangular Transformation:

From Generalized Hebbian Algorithm -  
Asymptotic Orthogonality of Embeddings

# Approximate Power Iteration

## Framing MINC as Approximate Power Iteration

The MINC objective can be seen doing approximate power iteration for finding the eigendecomposition of a low-rank approximation:

Low-Rank Eigendecomposition:  $\min_{\phi} \|M - FF'^{\top}\|_{\mathbb{F}}^2,$



Iterate Fixed Point Equation:  $\int \frac{p(x, x')}{\sqrt{p(x)p(x')}} \sqrt{p(x)} \phi(x) dx = \Lambda \sqrt{p(x')} \phi(x')$

Subject to Orthogonality:  $\int \phi(x) \phi(x)^{\top} p(x) dx = \int \phi(x') \phi(x')^{\top} p(x') dx' = \Lambda,$

# Approximate Power Iteration

## Framing MINC as Approximate Power Iteration

The MINC objective can be seen doing approximate power iteration for finding the eigendecomposition of a low-rank approximation:

Low-Rank Eigendecomposition:  $\min_{\phi} \|M - FF'^{\top}\|_{\mathbb{F}}^2,$



Iterate Fixed Point Equation:  $\int \frac{p(x, x')}{\sqrt{p(x)p(x')}} \sqrt{p(x)} \phi(x) dx \equiv \Lambda y$

Subject to  
Orthogonality:

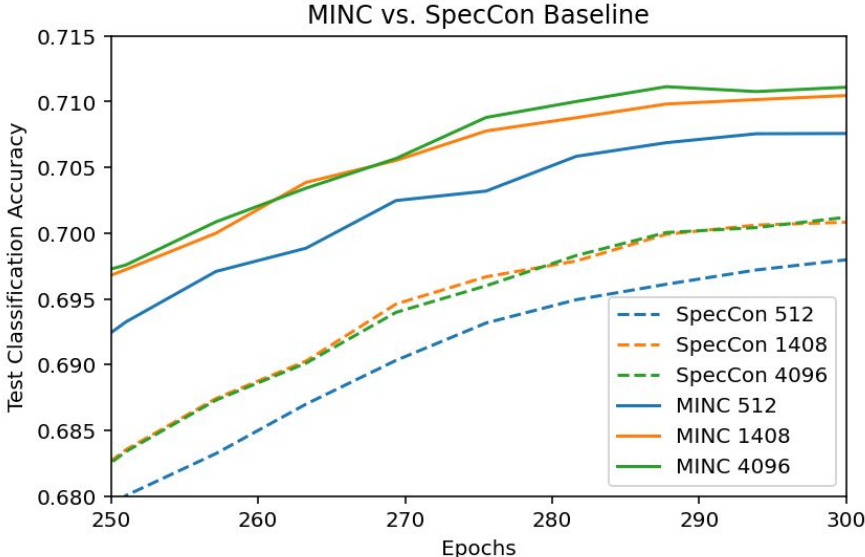
$$\int \phi(x) \phi(x)^{\top} p(x) dx = \int \phi(x') \phi(x')^{\top} p(x') dx \equiv \Lambda$$

See Poster #169  
and Paper for  
more details!

# Experimental Results

## MINC Outperforms Spectral Contrastive Loss at multiple batch sizes!

MINC at small batch size of 512 outperforms Spectral Contrastive at large batch size of 4096! More results in poster #169 and paper!





Thank you!  
See Poster #169 today  
(May 3)