

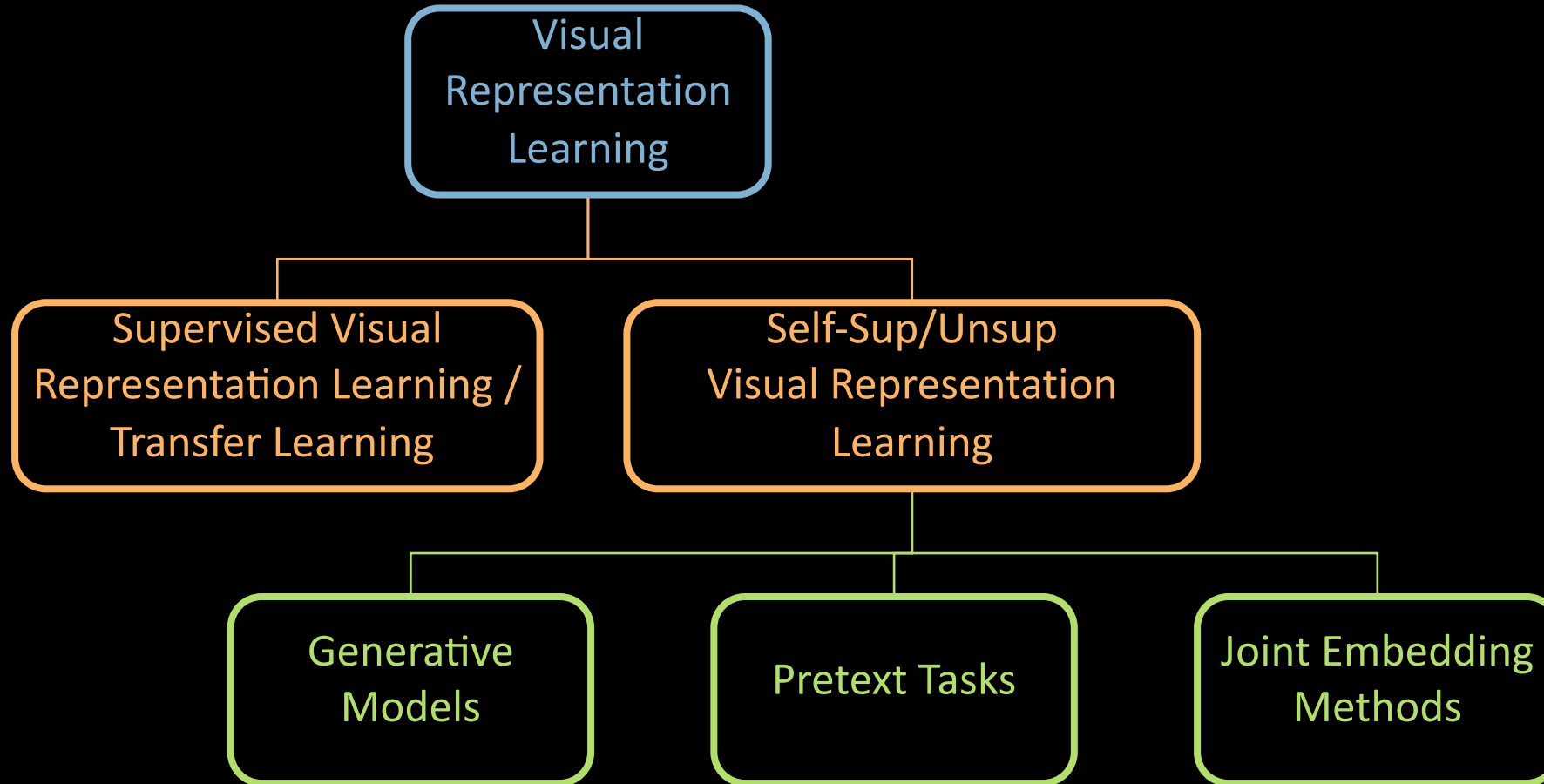
# Joint Embedding Methods

Self-Supervised Visual Representation Learning

# Visual representation learning

Overview

# Visual Representation Learning



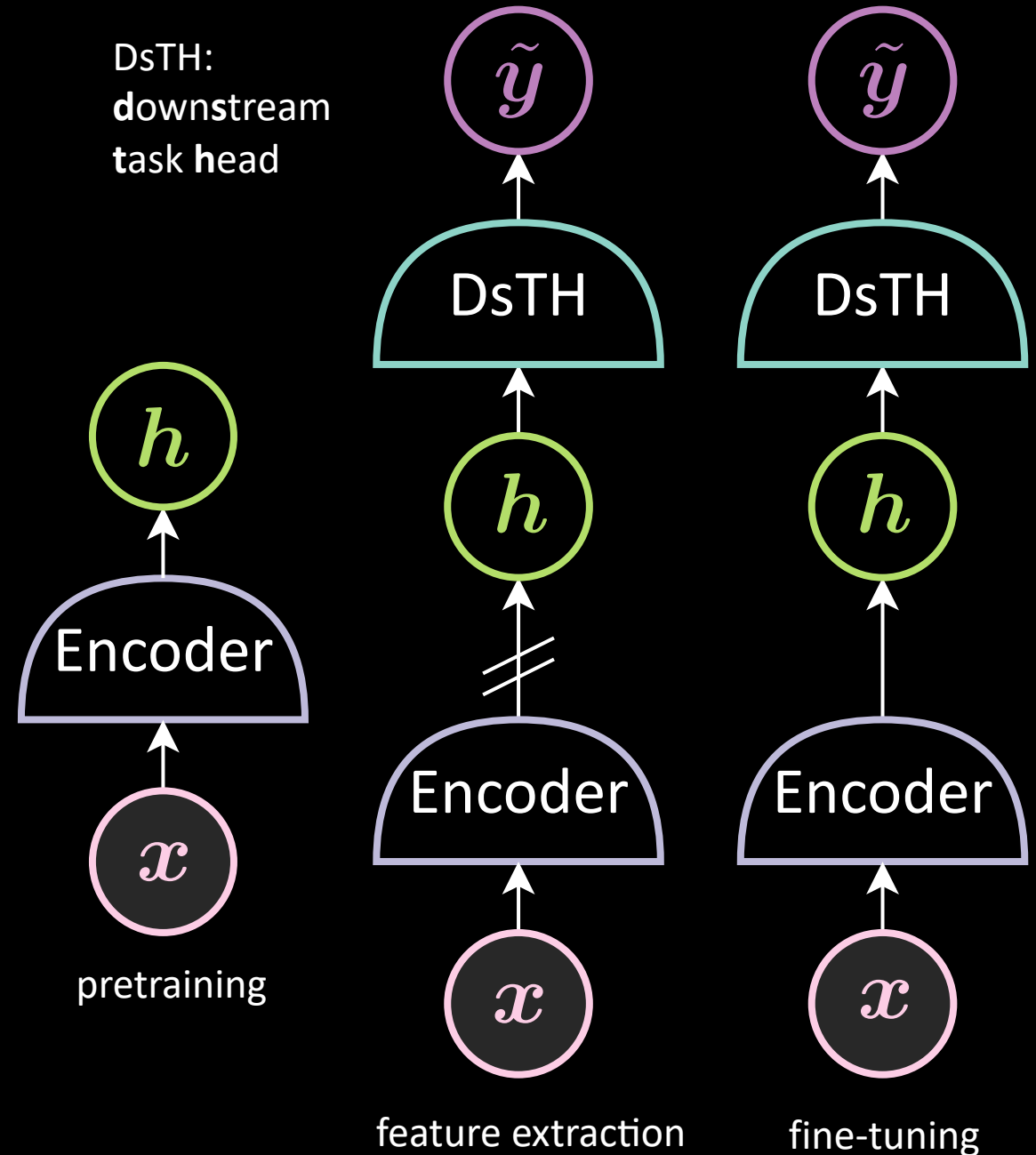
# Self-supervised visual representation learning

## Step 1: pretraining

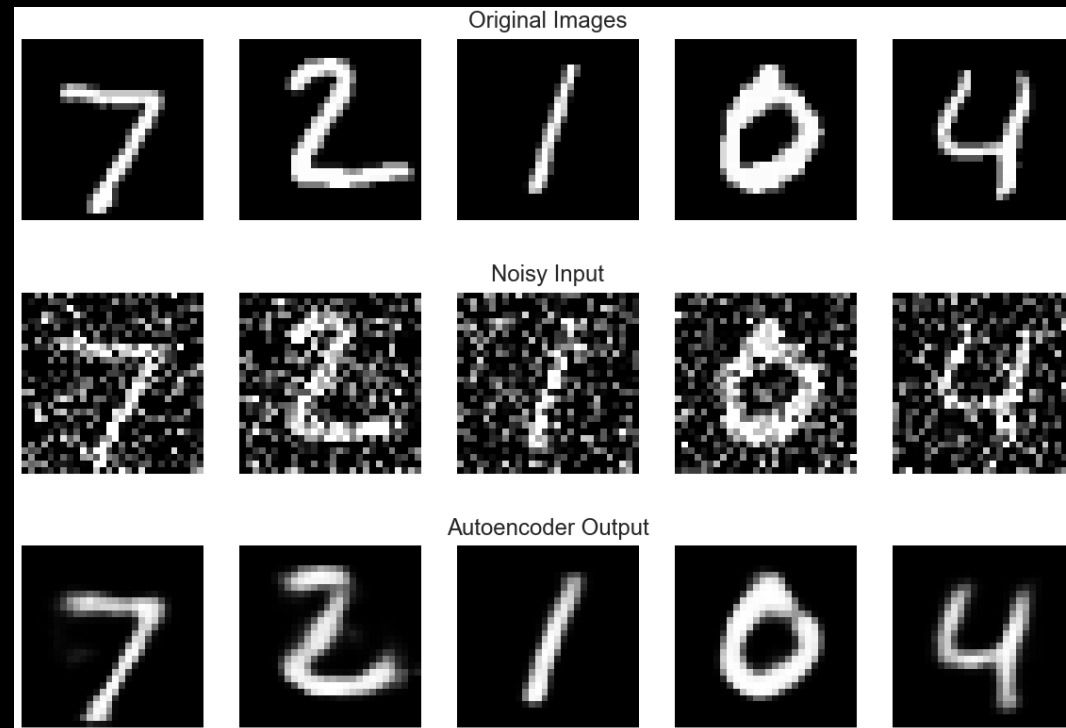
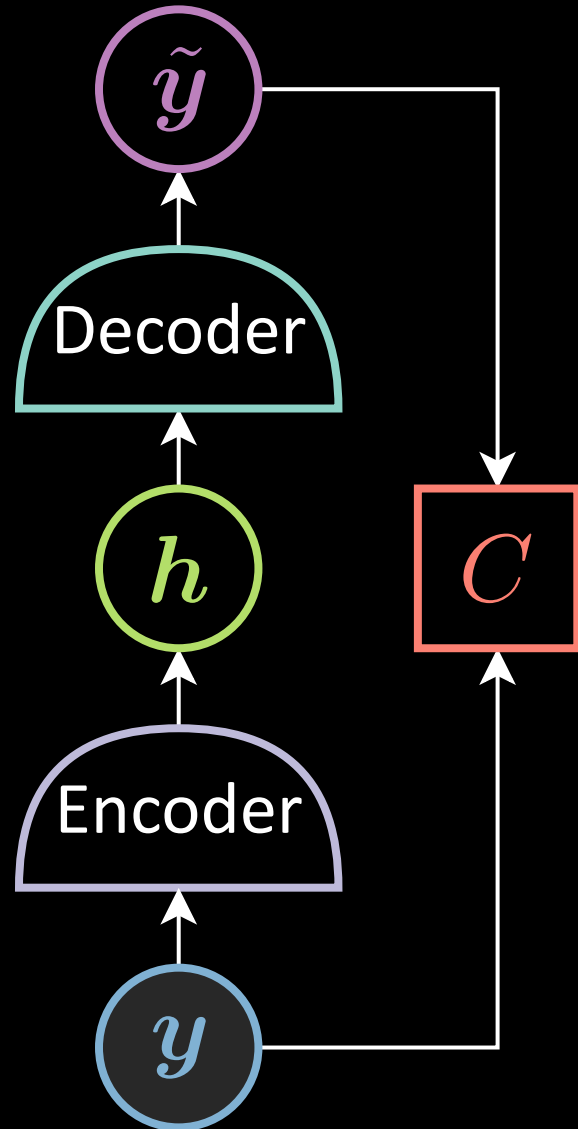
Use a large amount of unlabeled data to train a backbone network  
different methods will produce the backbone network differently

## Step 2: evaluation

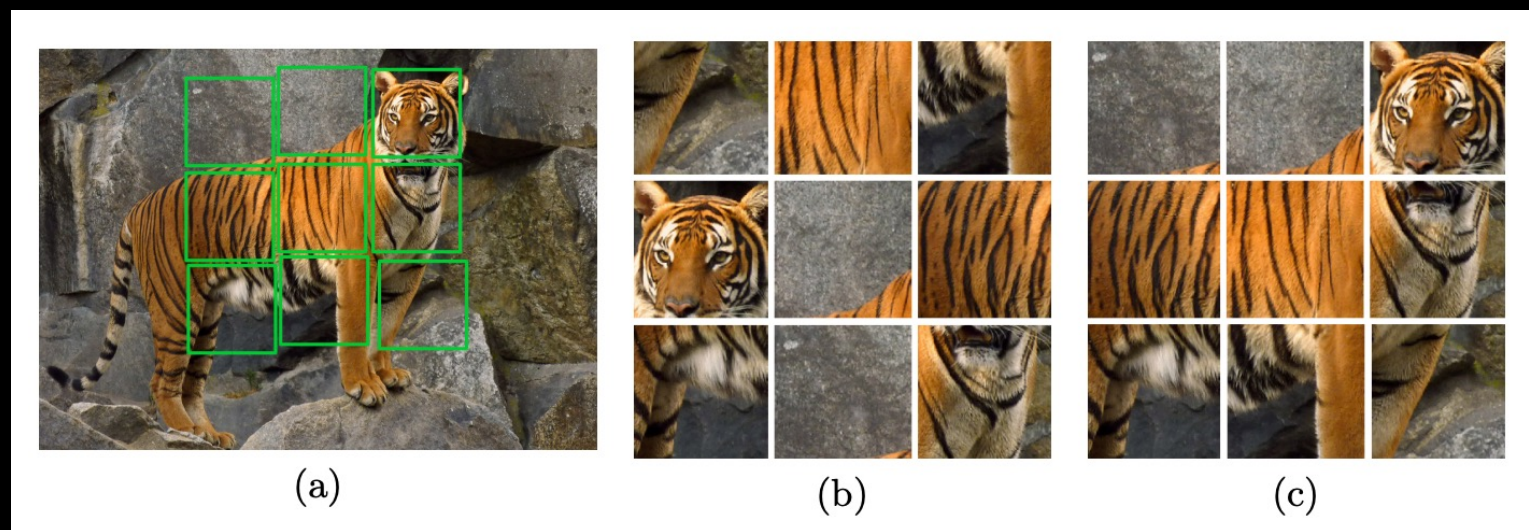
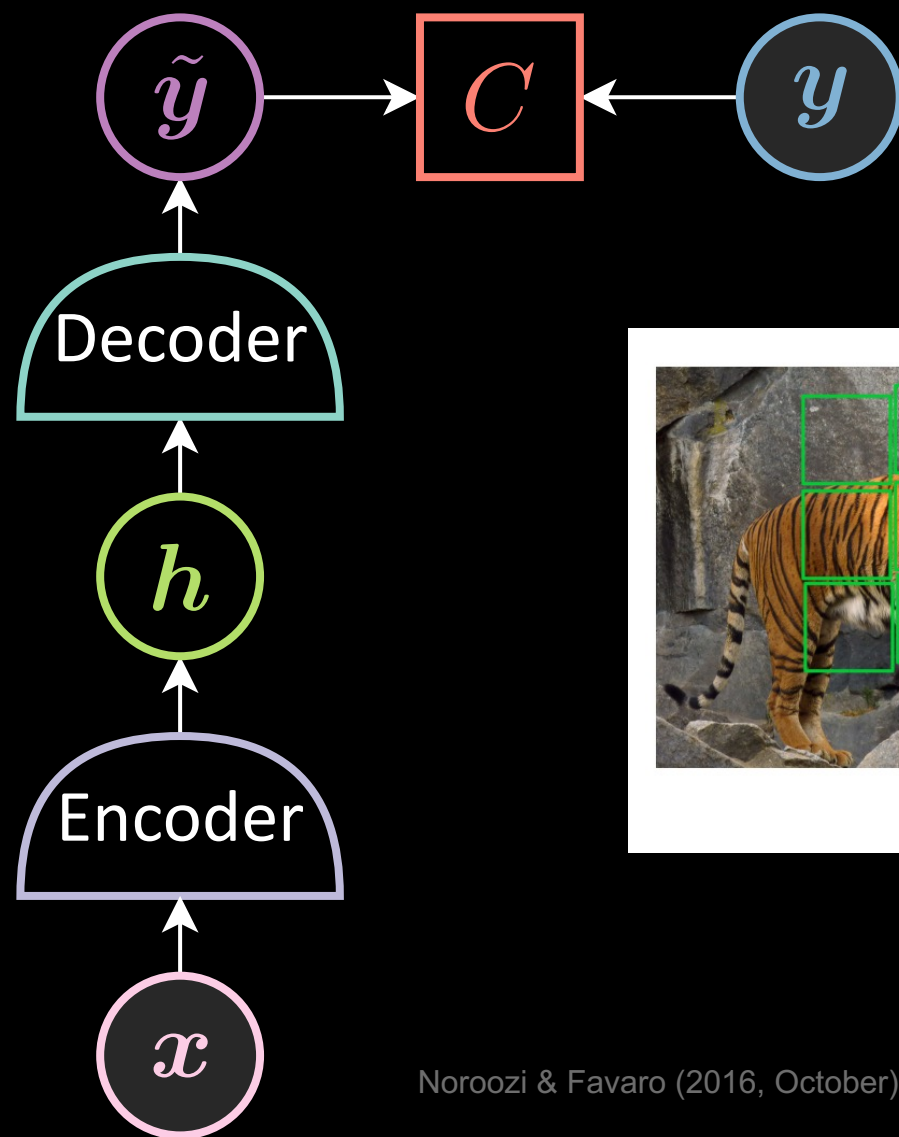
Use a small amount of labeled data to train a downstream task head network



# Generative Models- Autoencoder



# Pretext Tasks



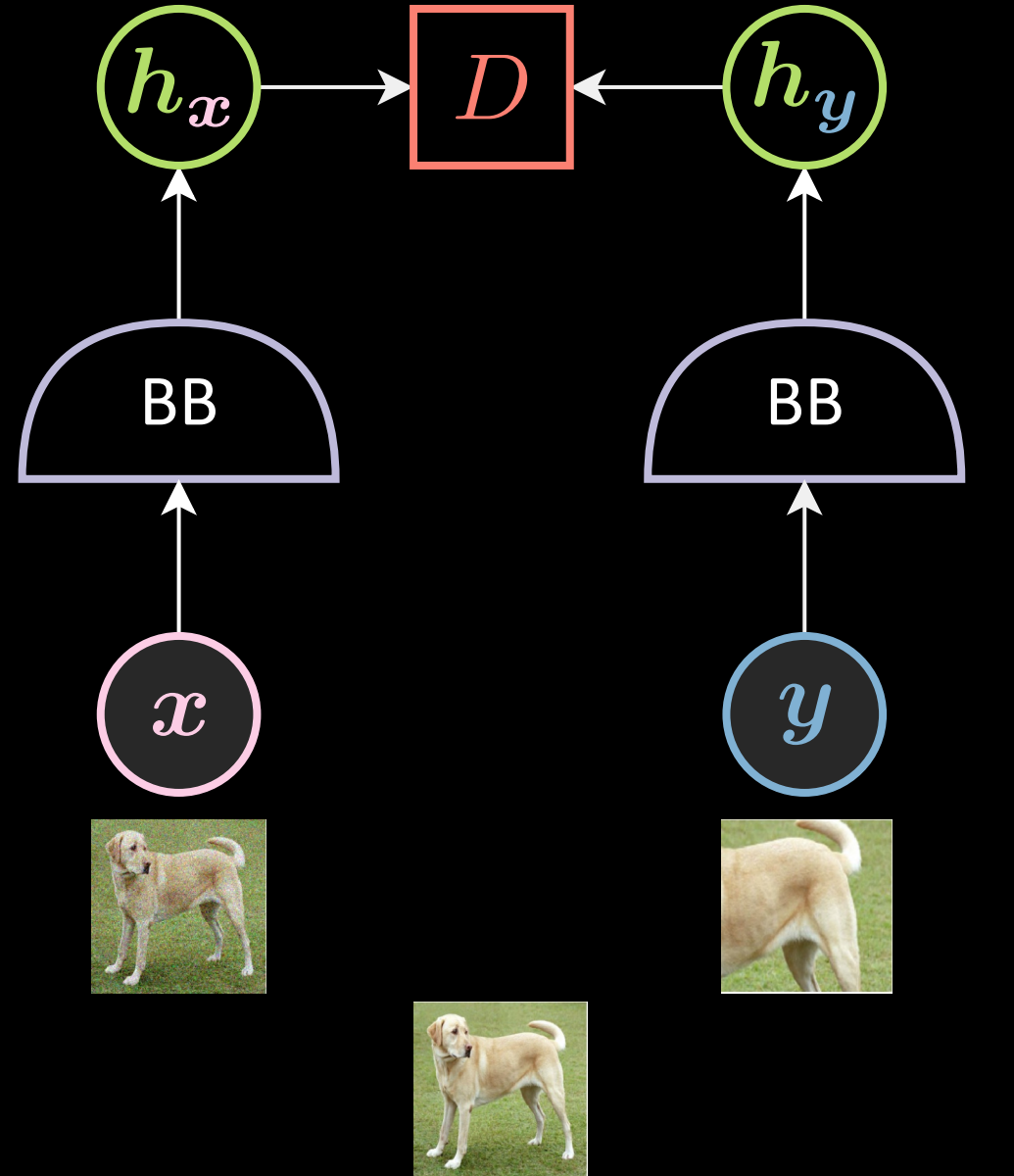
Noroozi & Favaro (2016, October). Unsupervised learning of visual representations by solving jigsaw puzzles.

# Joint embedding methods

Siamese nets & co.

# Joint Embedding Methods

Good backbone network should be robust to certain distortions  
(invariant to data augmentation)

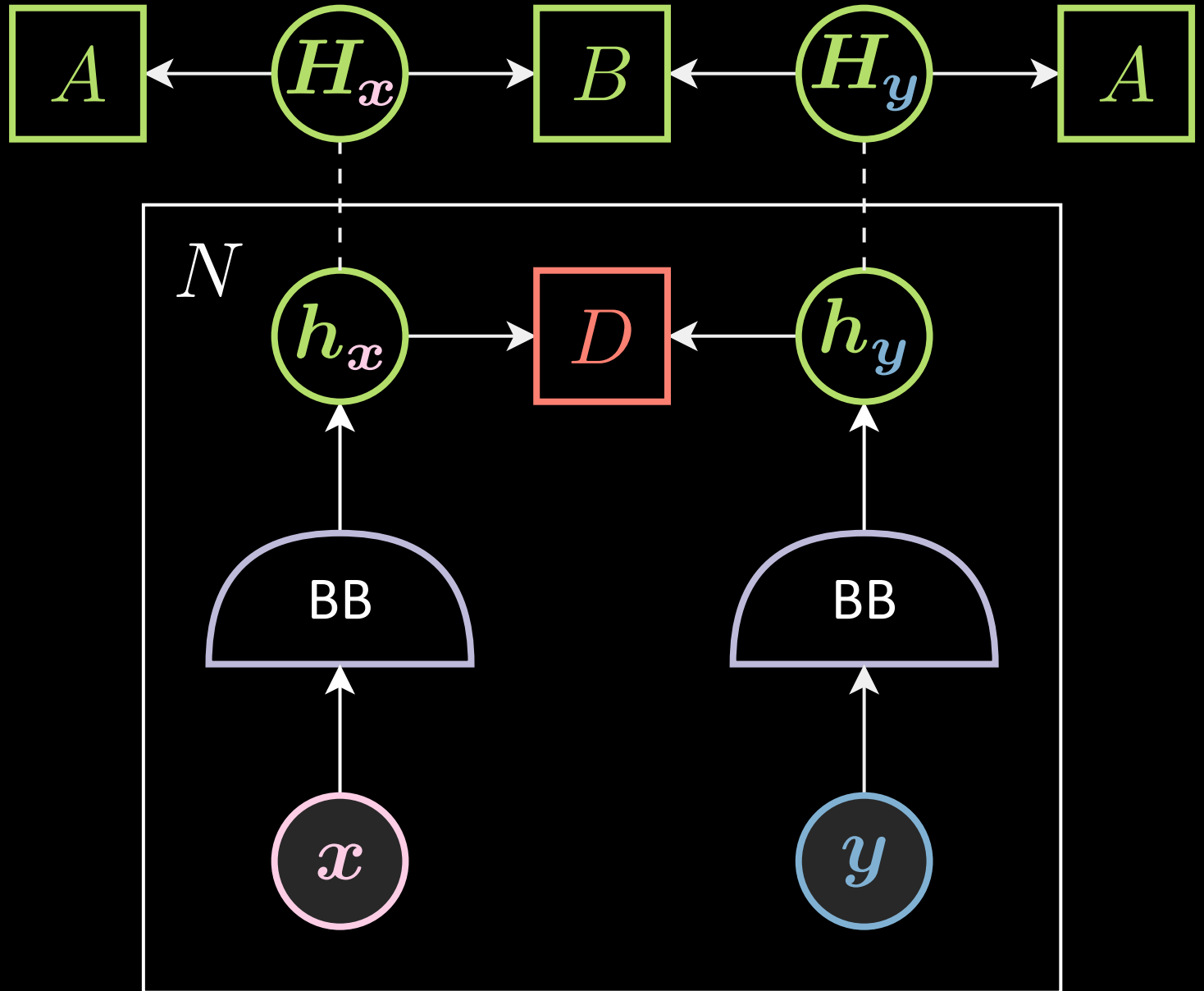


BB: backbone



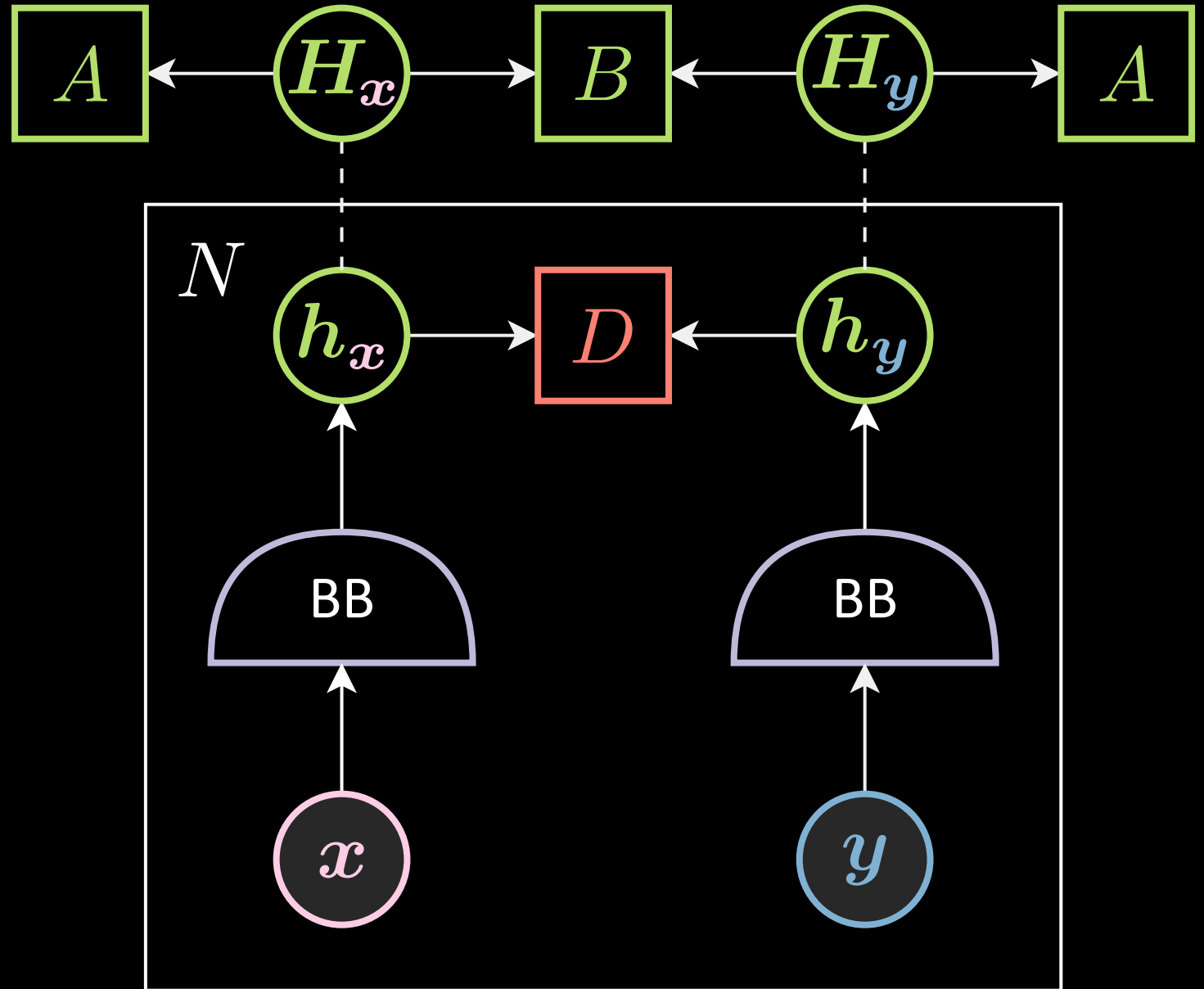
# JEM

We add extra loss to prevent the trivial solution (constant embeddings)



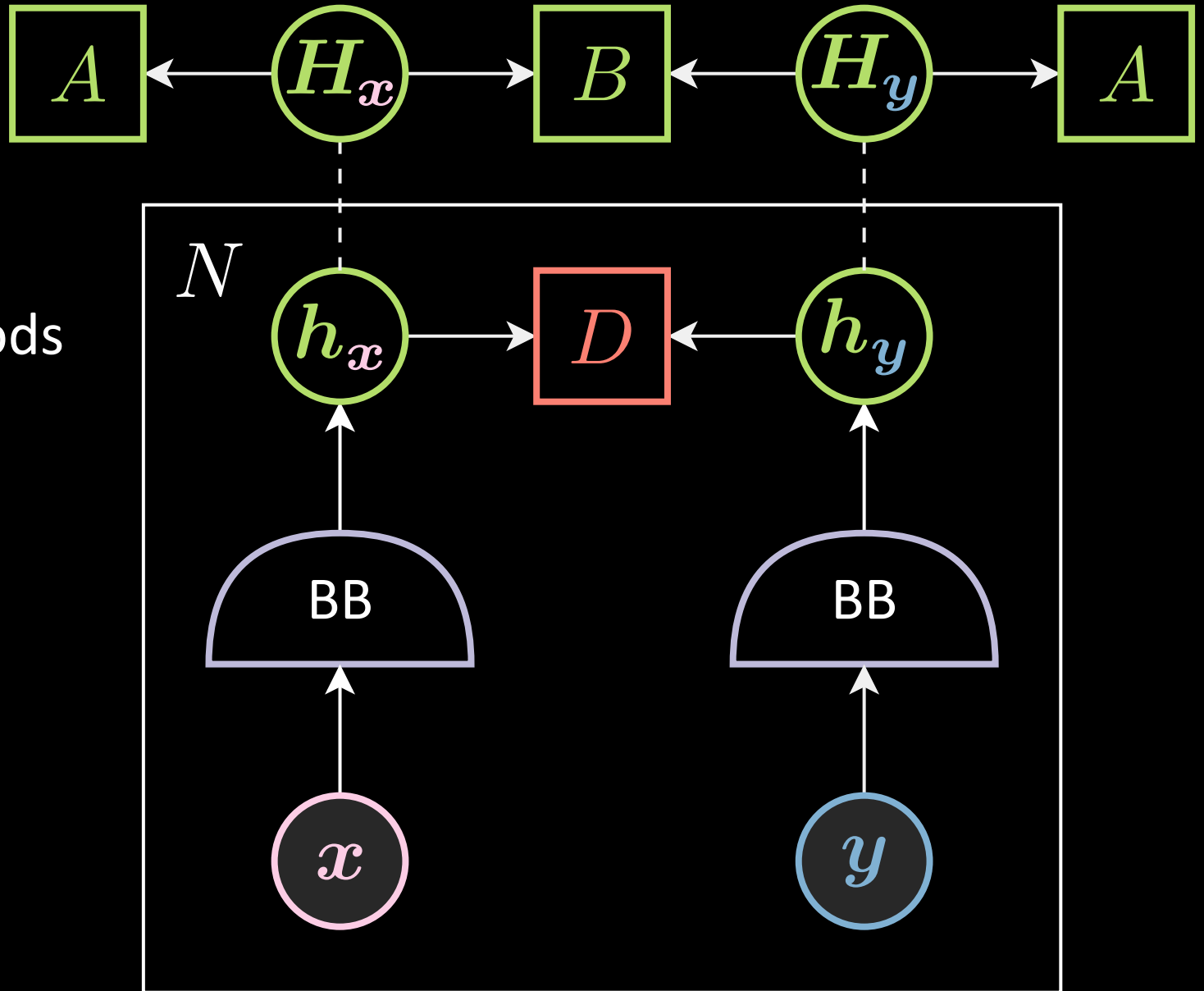
# JEM

1. Data augmentation
2. Backbone network
3. Energy function
4. Loss functional



# JEM

1. Contrastive methods
2. Non-contrastive methods
3. Clustering methods
4. Other methods



# JEM loss functions

- A term that pushes the positive pair closer
- An (implicit) term that prevents the trivial solution (constant output)

To make the training stable, people usually normalize the embeddings or put a hinge on the loss function to prevent the norm of embeddings becoming too large or too small

# Contrastive methods

Pull up on contrastive samples

# Contrastive Methods

The loss function should push

1. the positive pairs closer

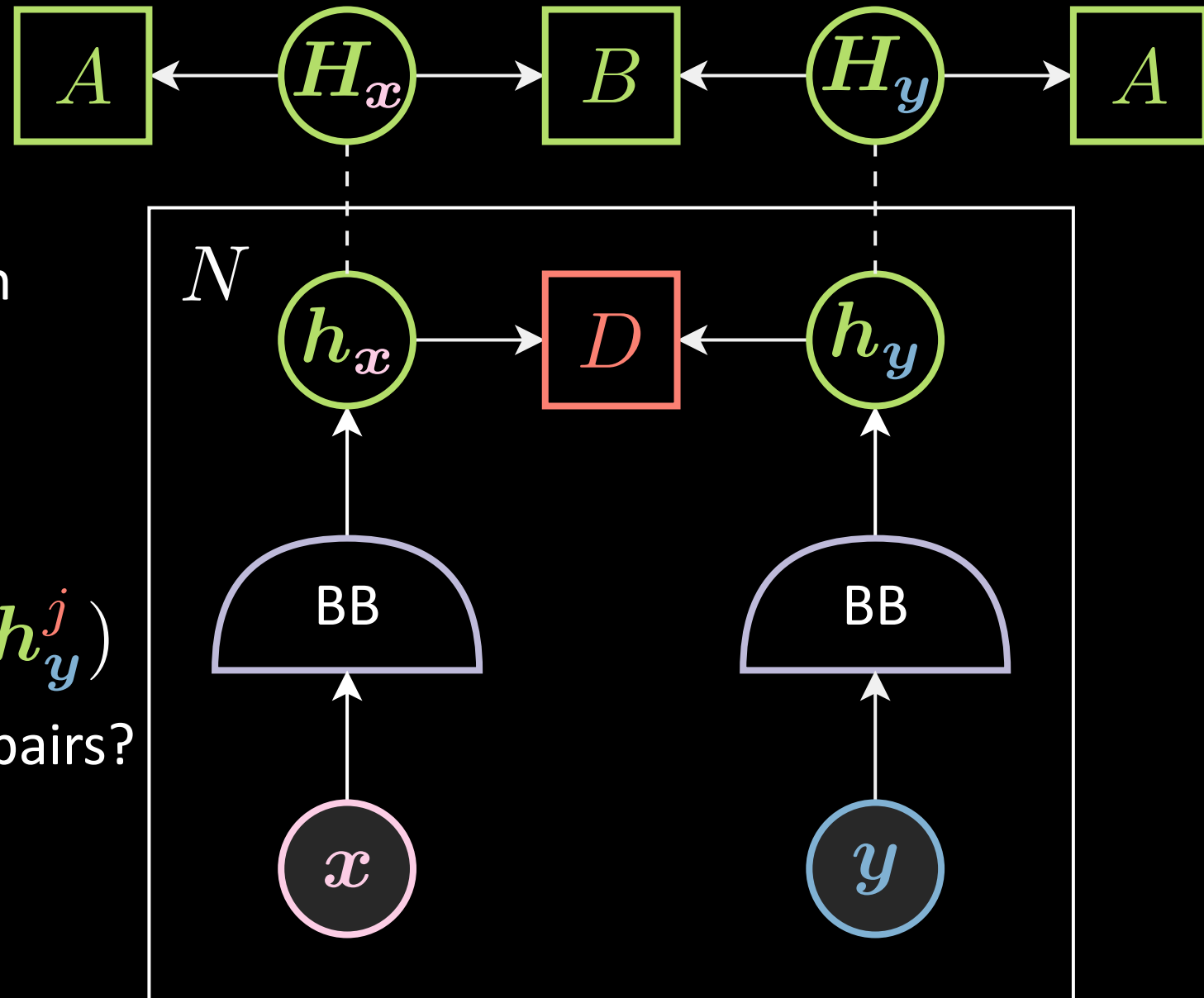
$$(h_x^i, h_y^i)$$

2. the negative pairs away

$$(h_x^i, h_x^j), (h_x^i, h_y^j), (h_y^i, h_y^j)$$

How to find a good negative pairs?

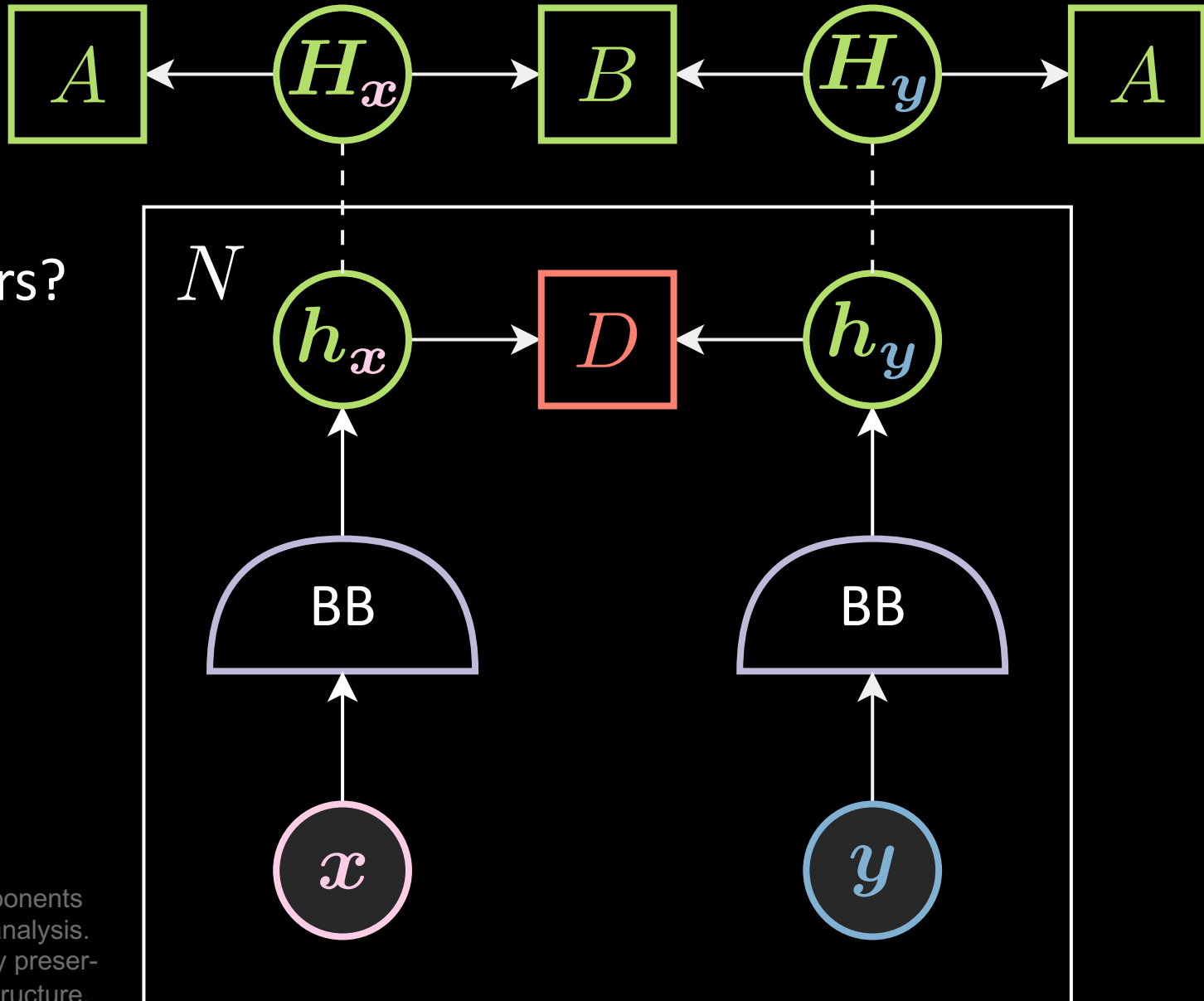
hard negative mining?



# Contrastive Methods

## SimCLR and MoCo

- How to find a good negative pairs?
- Use large batch size!
- Both SimCLR and MoCo use the InfoNCE loss function:



Goldberger, Hinton, Roweis & Salakhutdinov (2004). Neighbourhood components analysis.

Salakhutdinov & Hinton (2007, March). Learning a nonlinear embedding by preserving class neighbourhood structure.

Van den Oord, Li & Vinyals (2018). Representation learning with contrastive predictive coding.

Chen, Kornblith, Norouzi & Hinton (2020, November). A simple framework for contrastive learning of visual representations.

He, Fan, Wu, Xie & Girshick (2020). Momentum contrast for unsupervised visual representation learning.

# The InfoNCE cost function

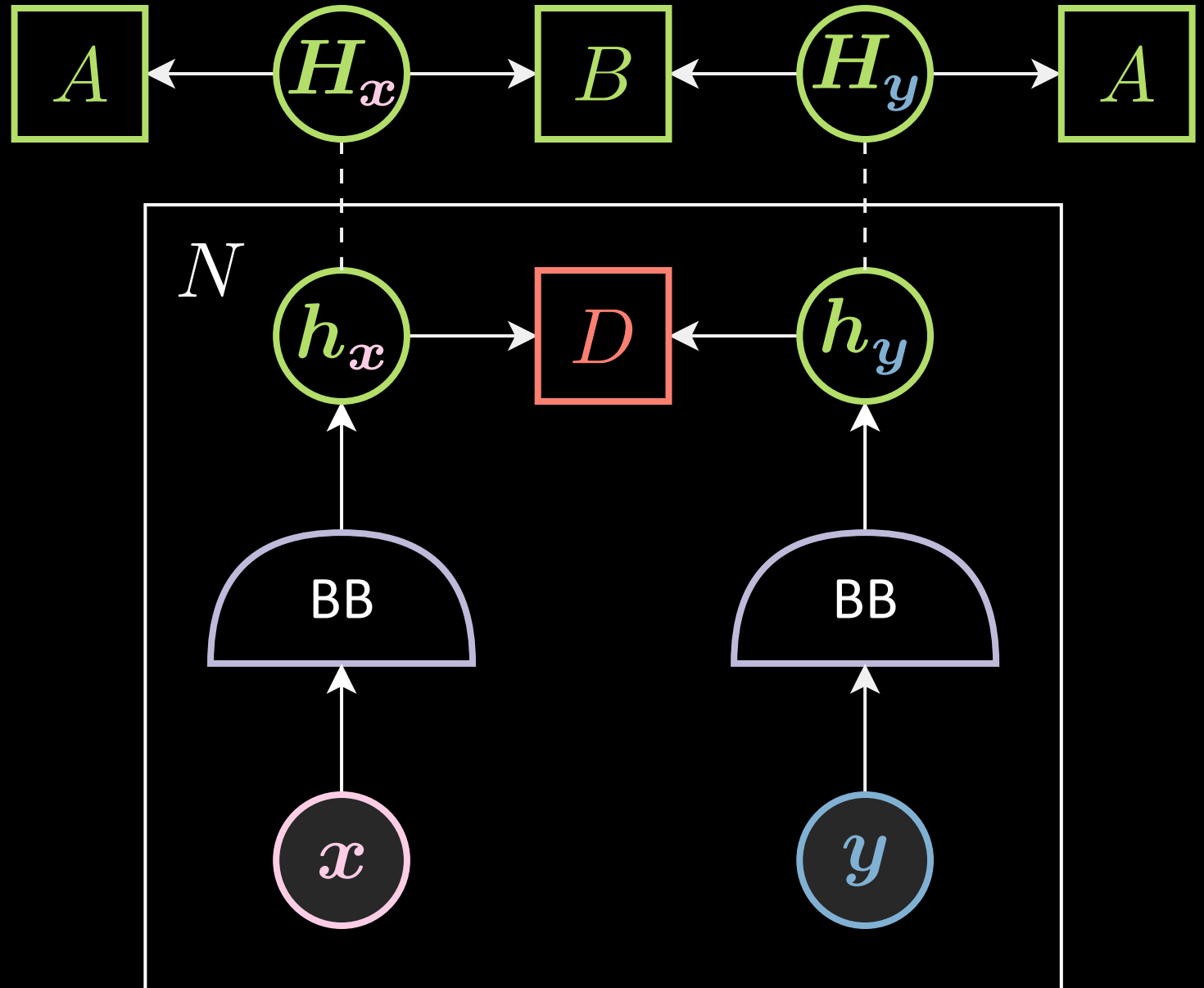
$$\begin{aligned} D(\mathbf{h}_x, \mathbf{h}_y) &= \\ &= -\log \frac{\exp(\beta \operatorname{sim}(\mathbf{h}_x, \mathbf{h}_y))}{\sum_j^N \exp(\beta \operatorname{sim}(\mathbf{h}_x, \mathbf{h}_x^j)) + \sum_j^N \exp(\beta \operatorname{sim}(\mathbf{h}_x, \mathbf{h}_y^j))} \\ &= -\beta \operatorname{sim}(\mathbf{h}_x, \mathbf{h}_y) + \log \left[ \sum_j^N \exp(\beta \operatorname{sim}(\mathbf{h}_x, \mathbf{h}_x^j)) + \sum_j^N \exp(\beta \operatorname{sim}(\mathbf{h}_x, \mathbf{h}_y^j)) \right] \\ &= -\beta \operatorname{sim}(\mathbf{h}_x, \mathbf{h}_y) + \operatorname{softmax}_{\beta, j} [\operatorname{sim}(\mathbf{h}_x, \mathbf{h}_x^j), \operatorname{sim}(\mathbf{h}_x, \mathbf{h}_y^j)] \\ \operatorname{sim}(\mathbf{h}_x, \mathbf{h}_y) &= \frac{\mathbf{h}_x^\top \mathbf{h}_y}{\|\mathbf{h}_x\| \|\mathbf{h}_y\|} \end{aligned}$$

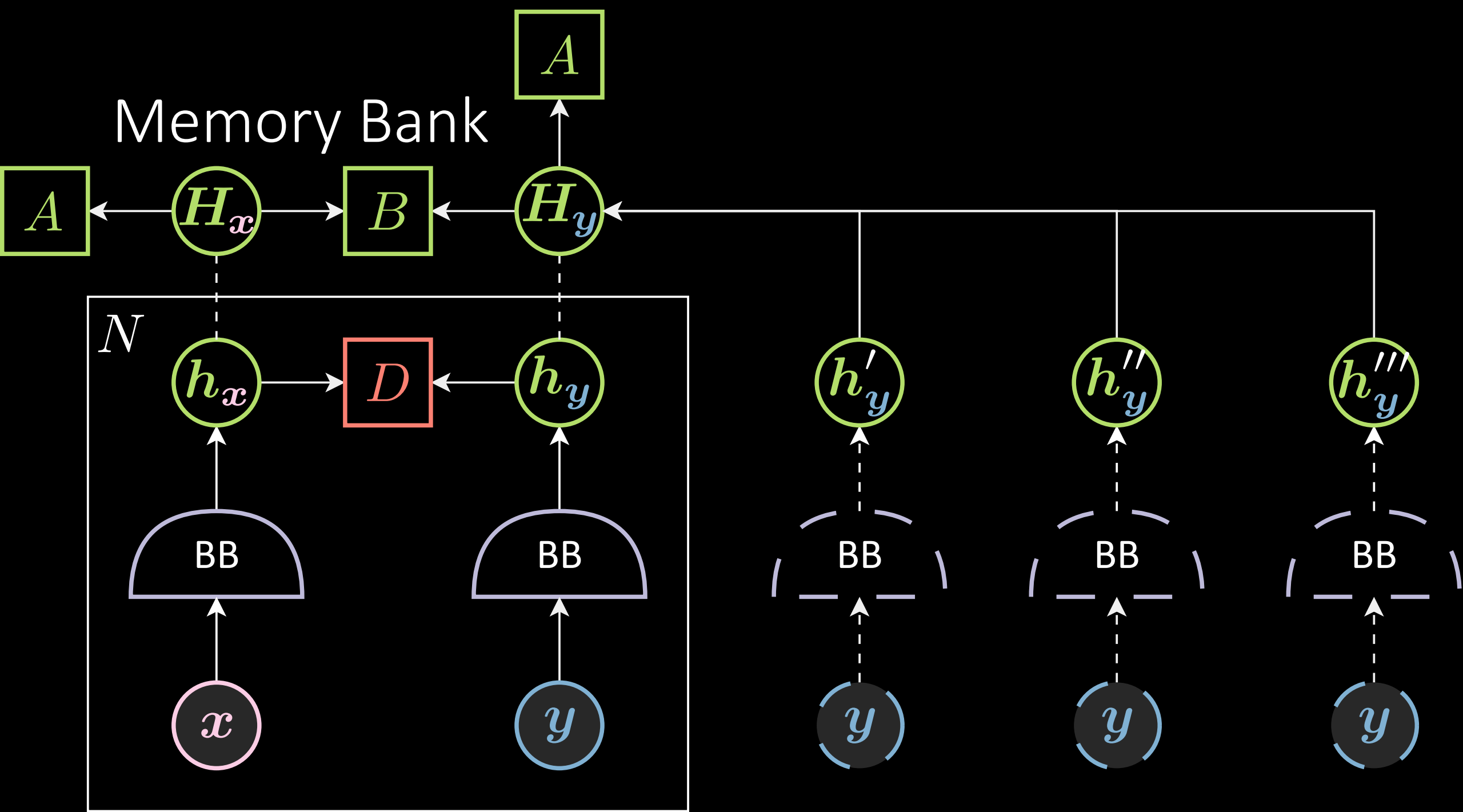


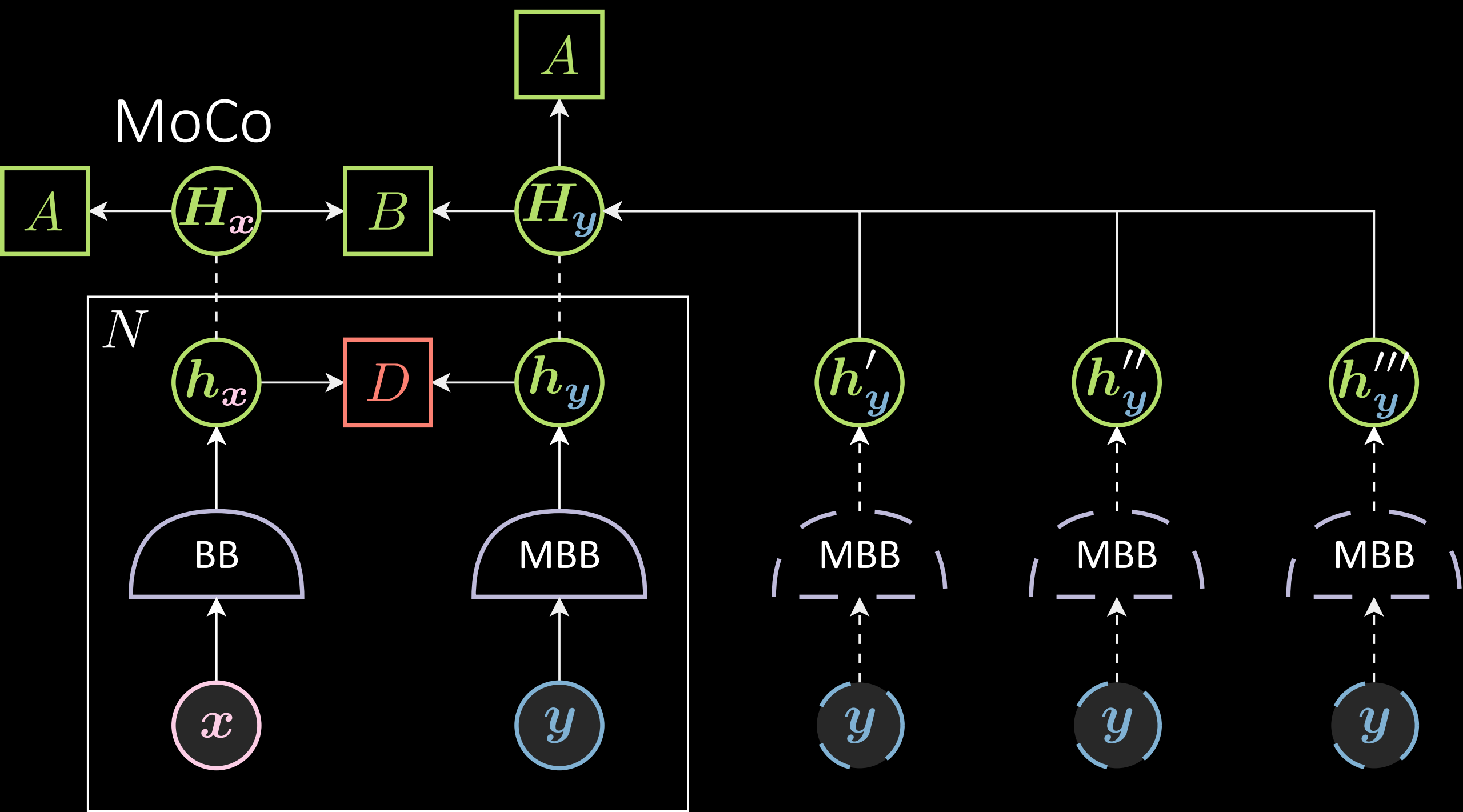
Batch size

Very large!

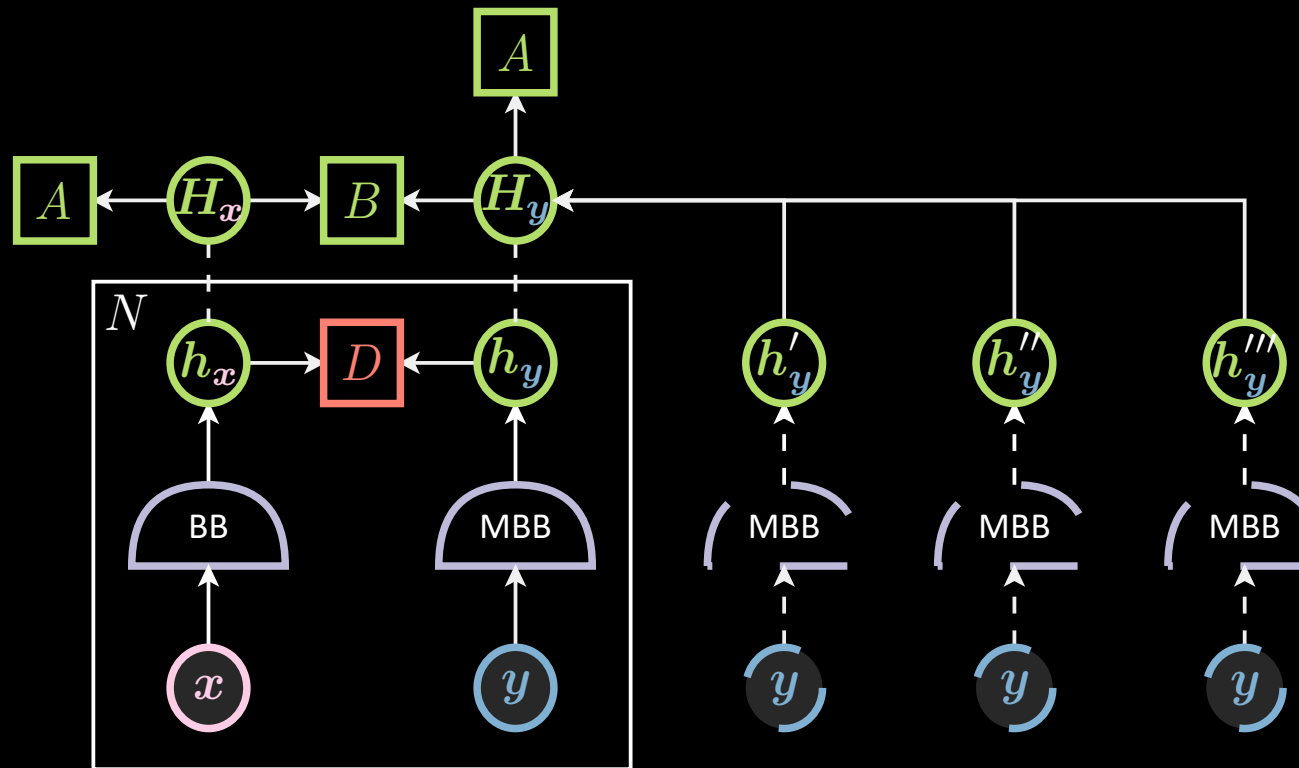
SimCLR  $N = 8192$







# Contrastive Methods – MoCo



MoCo n = 256

$$\theta_{t+1} = \theta_t - \eta \nabla \theta_t$$

$$\vartheta_{t+1} = m\vartheta_t + (1 - m)\theta_{t+1}$$

$\theta$  : backbone parameters

$\vartheta$  : momentum backbone parameters

# Recap

What we've learnt so far

# Quick Recap

1. Visual representation learning: pretraining + evaluation
2. Generative vs. pretext task vs. joint-embedding methods
3. Joint-embedding methods:
  1. Invariant to data augmentation
  2. Prevent trivial solution
4. Contrastive methods
  1. Hard negative mining
  2. Large negative sample pool (SimCLR vs. MoCo)

# Non-contrastive methods

Prevent trivial solution without negative samples

# The Disadvantages of Contrastive Methods

Contrastive methods require techniques such as:

1. weight sharing between the branches
2. batch normalization
3. feature-wise normalization
4. output quantization
5. stop gradient
6. memory banks
7. ...



# Non-contrastive methods and information theory

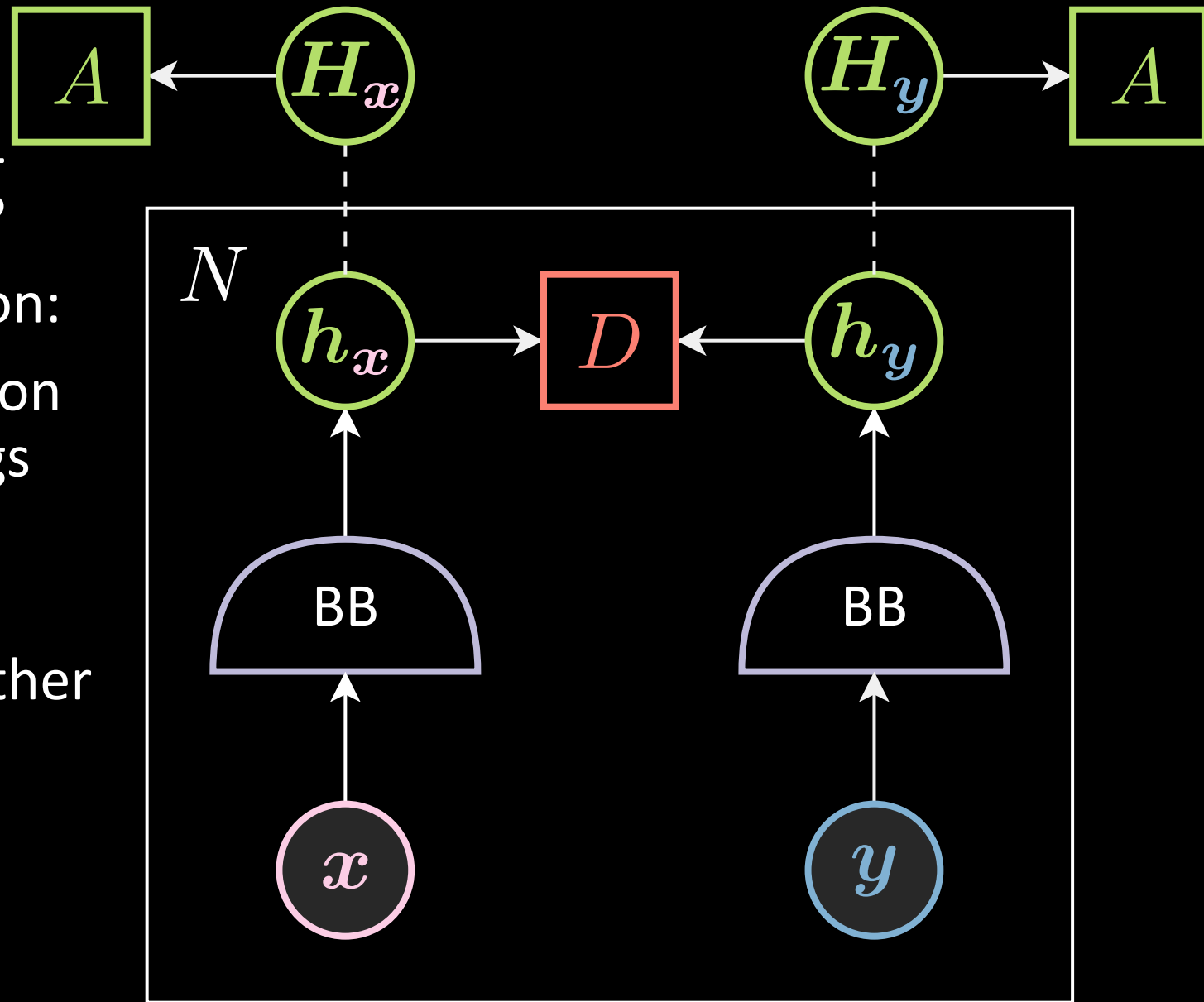
## Most of non-contrastive methods

- based on information theory
  - Redundancy reduction (Barlow Twins)
  - Information maximization (VICReg)
- Don't require special architectures

# Non-Contrastive Methods – VICReg

The Information maximization:

1. to maximize the information content of the embeddings
2. produce embedding variables that are decorrelated from each other
3. prevent an informational collapse in which the variables carry redundant information



# Non-Contrastive Methods – VICReg

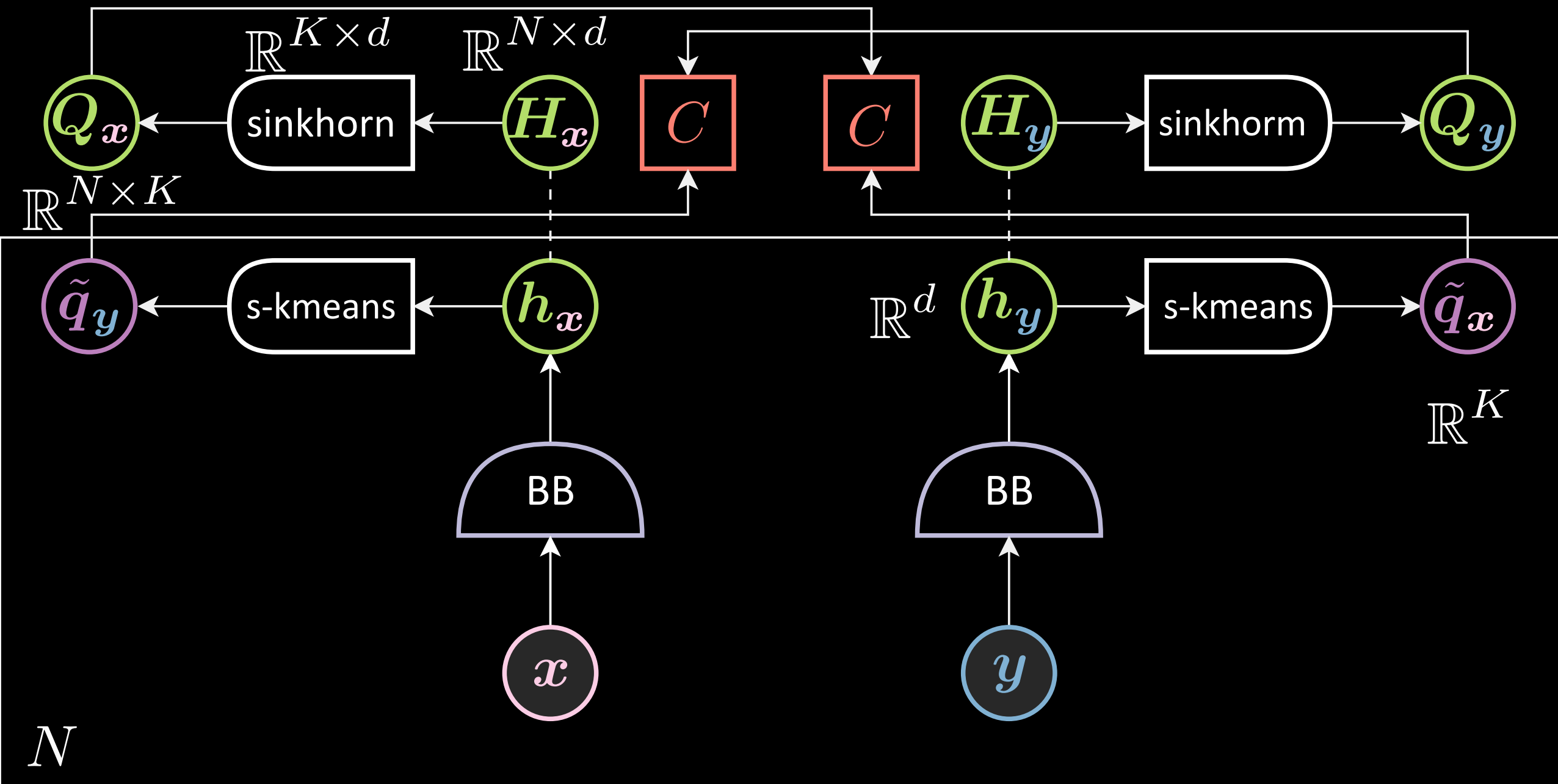
The loss function is pushing

1. the positive pairs closer ( $\mathbf{h}_x^i, \mathbf{h}_y^i$ )
2. the variance of the embeddings large
3. The covariance of the embeddings small

$$\begin{aligned} L(\mathbf{w}, \mathbf{x}, \mathbf{y}) &= \|\mathbf{h}_x - \mathbf{h}_y\|^2 \\ &+ \frac{1}{d} \left[ \sum_i^d (\gamma - \mathbf{x} \mathbf{C}_{ii})^+ + (\gamma - \mathbf{y} \mathbf{C}_{ii})^+ \right] \\ &+ \frac{1}{d} \left[ \sum_i^d \sum_{j \neq i}^d \mathbf{x} \mathbf{C}_{ij}^2 + \mathbf{y} \mathbf{C}_{ij}^2 \right] \end{aligned} \quad \mathbf{C} = \frac{1}{N} \mathbf{H}^\top \mathbf{H}.$$

# Clustering Methods- SwAV

Prevent trivial solution by quantizing the embedding space



## SwAV (II)

$$Q_x = \text{sinkhorn}_W(\mathbf{H}_x) \in \mathbb{R}^{N \times K}$$

$$Q_x = [\mathbf{q}_x^1, \dots, \mathbf{q}_x^N]^\top$$

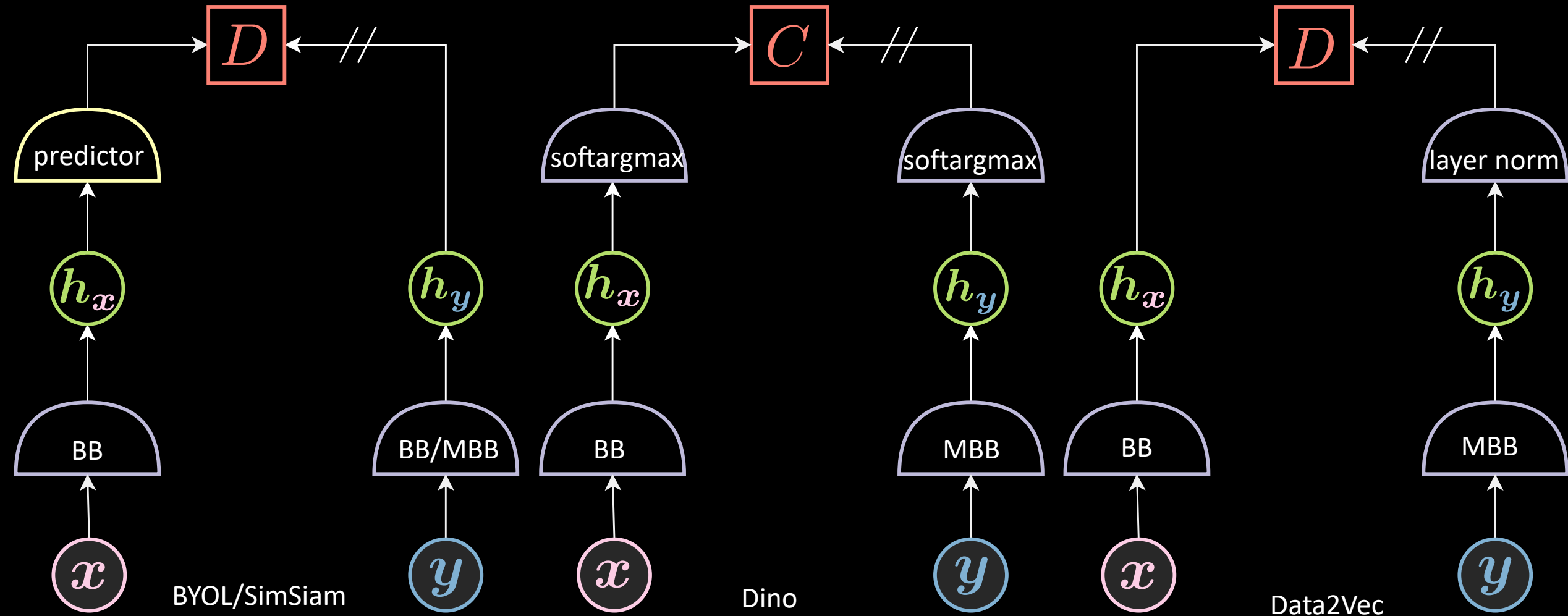
$$\mathbf{W} \in \mathbb{R}^{K \times d} : \text{dictionary}$$

$$\tilde{\mathbf{q}}_{\textcircled{x}} = \text{softargmax}_\beta(\mathbf{W} \mathbf{h}_{\textcircled{y}}) \in \mathbb{R}^K$$

$$F(x, y) = C(\mathbf{q}_x, \tilde{\mathbf{q}}_x) + C(\mathbf{q}_y, \tilde{\mathbf{q}}_y)$$

Other methods

# Other Methods – BYOL, SimSiam, Dino, Data2Vec



Grill, Strub, Altché, Tallec, Richemond, Buchatskaya, ... & Valko (2020). Bootstrap your own latent-a new approach to self-supervised learning.

Chen, X., & He, K. (2021). Exploring simple siamese representation learning.

Caron, Touvron, Misra, Jégou, Mairal, Bojanowski & Joulin (2021). Emerging properties in self-supervised vision transformers.

Baevski, Hsu, Xu, Babu, Gu & Auli (2022). Data2vec: A general framework for self-supervised learning in speech, vision and language.

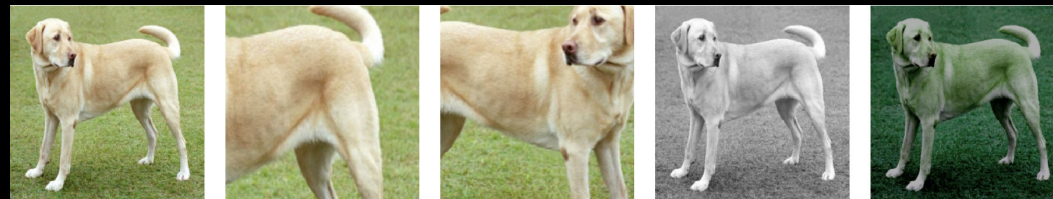


# Data augmentation and network architecture

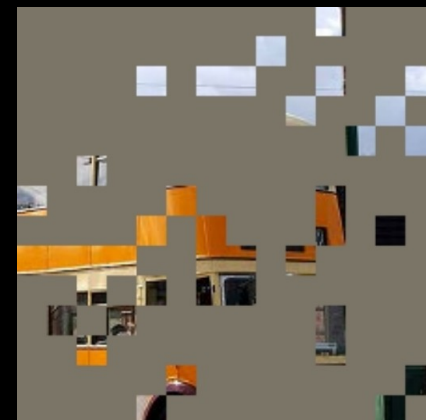
# Data augmentation

The SimCLR/BYOL data augmentation:

1. Random Crop (the most critical one)
2. Flip
3. Color Jitter
4. Gaussian Blur



For traditional augmentation to  
masking augmentation



# Network architecture

It is always better to add a two/three-layer projector/expander

Even without memory bank, momentum encoder usually helps the performance of the downstream tasks, especially with weak data augmentation

