

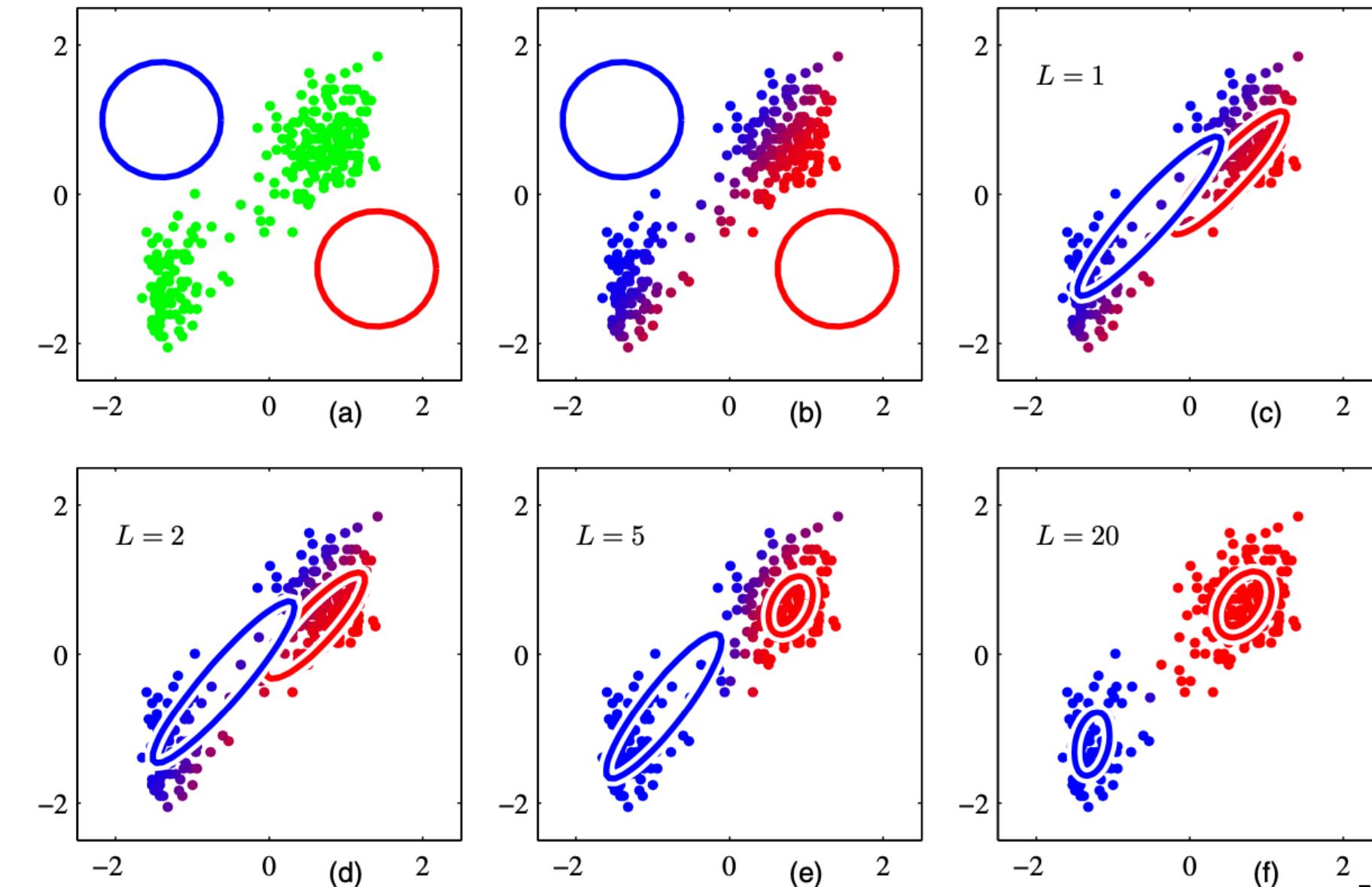
SELF-SUPERVISED REPRESENTATION LEARNING

Jonathan Gordon and Javier Antorán

Self Supervised Learning (SSL)

Unsupervised Learning:

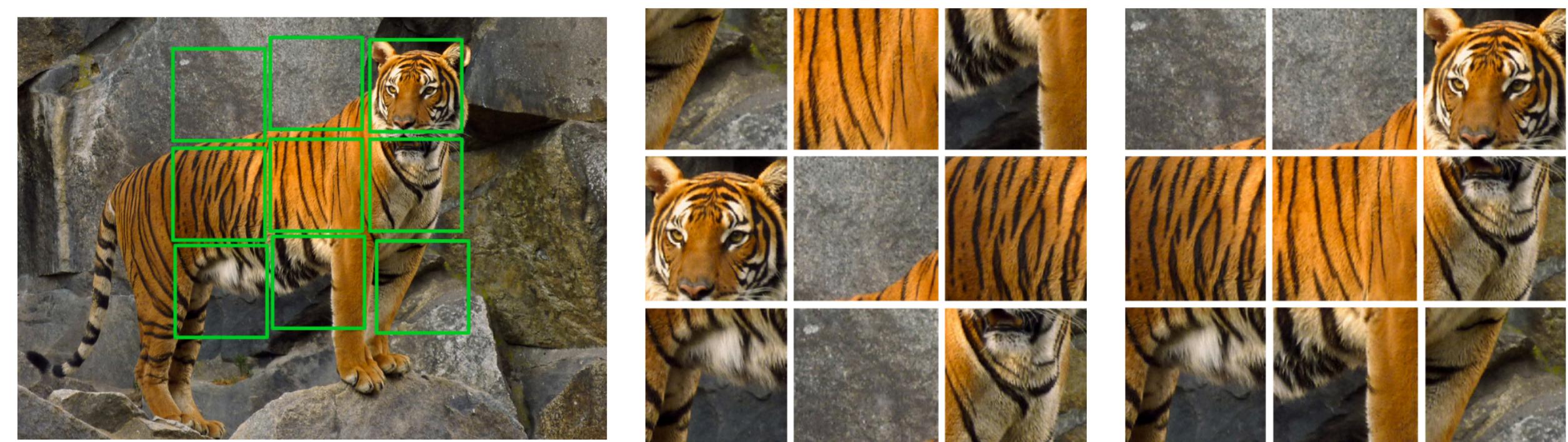
Learning to Model the Observed Data



[C. Bishop]

Self Supervised Learning:

Learning by Predicting Proxy Targets
Extracted From Unlabelled Data



[Larsson et. al.]

But Why Self Supervised Learning?

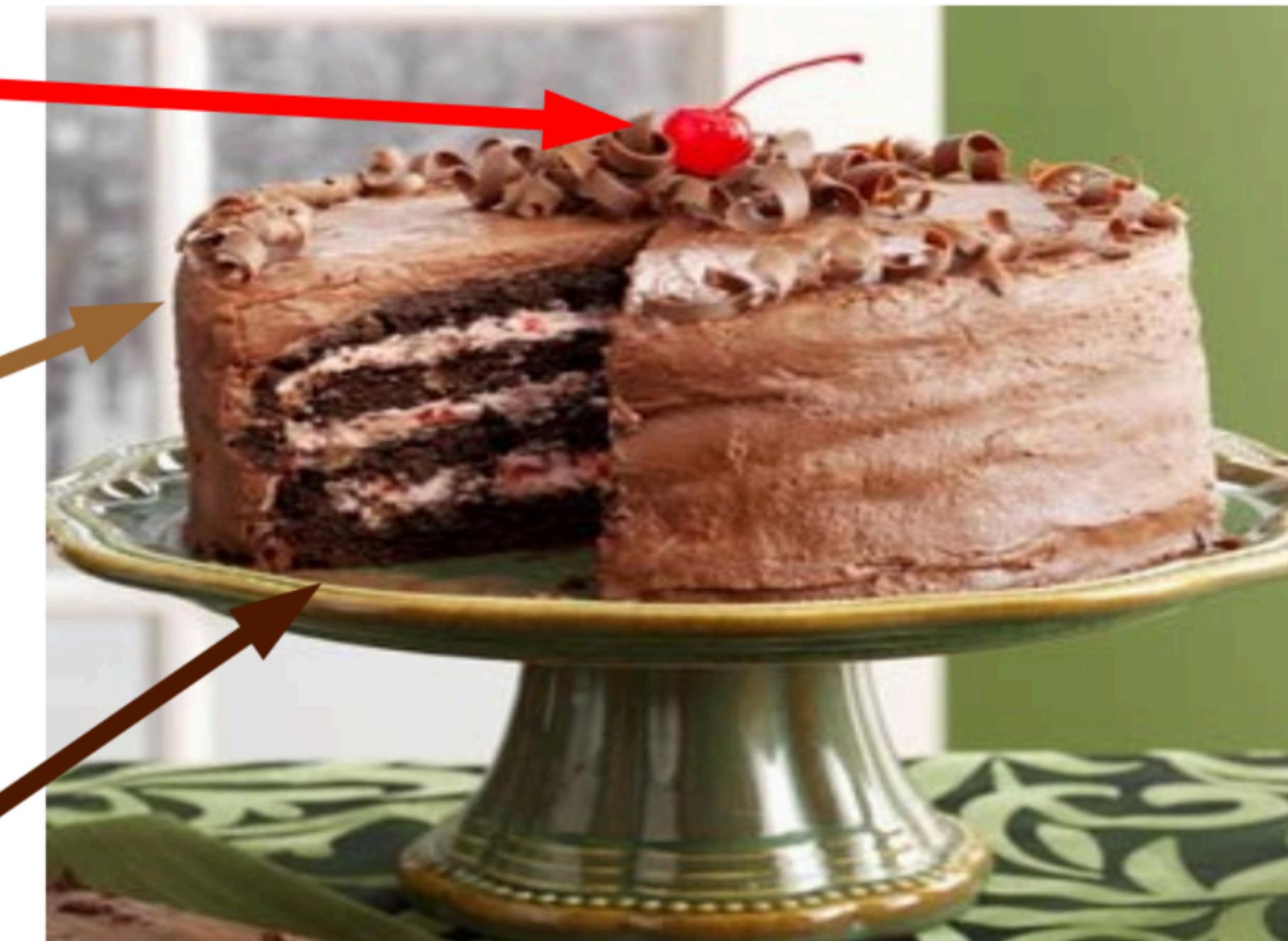
Y. LeCun

How Much Information is the Machine Given during Learning?

- ▶ “Pure” Reinforcement Learning (**cherry**)
 - ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

- ▶ Supervised Learning (**icing**)
 - ▶ The machine predicts a category or a few numbers for each input
 - ▶ Predicting human-supplied data
 - ▶ **10 → 10,000 bits per sample**

- ▶ Self-Supervised Learning (**cake génoise**)
 - ▶ The machine predicts any part of its input for any observed part.
 - ▶ Predicts future frames in videos
 - ▶ **Millions of bits per sample**

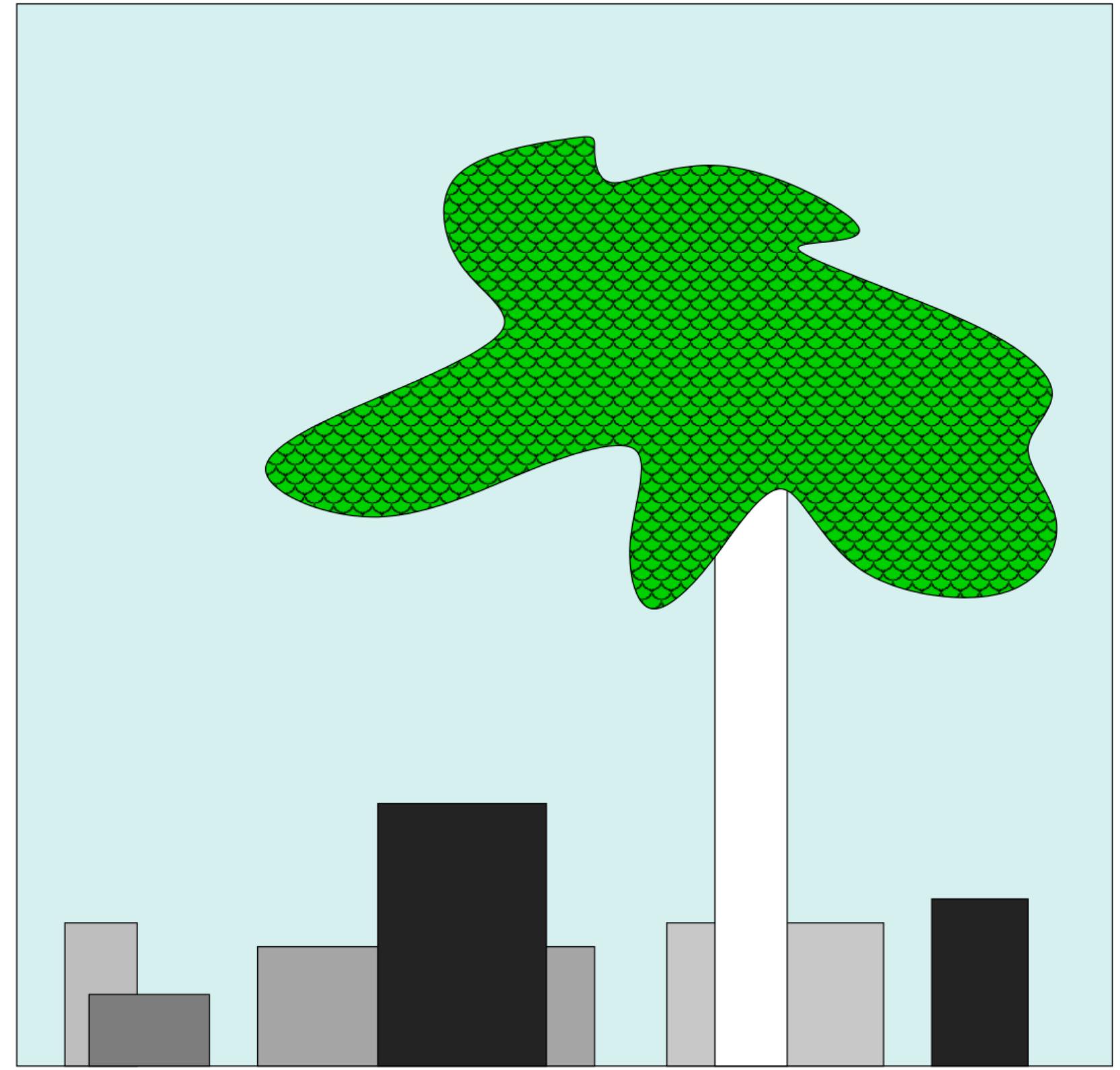


[Y. LeCun]

But Why Self Supervised Learning?



[Y. LeCun]



[D. MacKay]

In Other Words:



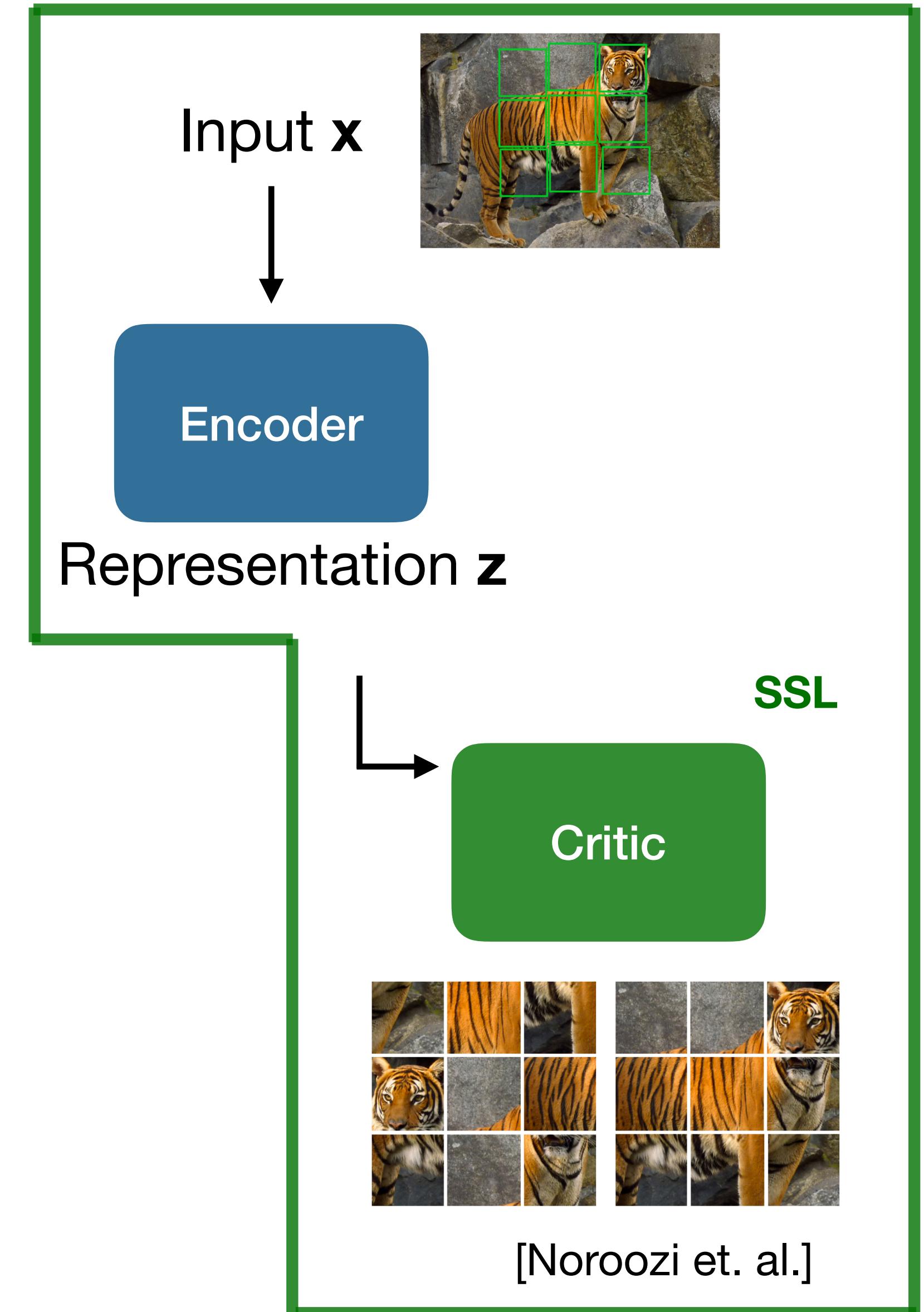
[Y. LeCun]

Talk Layout

1. Examples of Self Supervised Learning (SSL) Approaches
 - Information Theoretic Interpretation
2. Mutual Information (MI) Motivated SSL approaches
 - Deep InfoMax
 - Contrastive Predictive Coding
3. Is MI the real reason behind the success of SSL?
4. Generative Models for Representation Learning
5. Self Supervised Learning for Identifiability in Non-Linear ICA

Typical SSL Approaches:

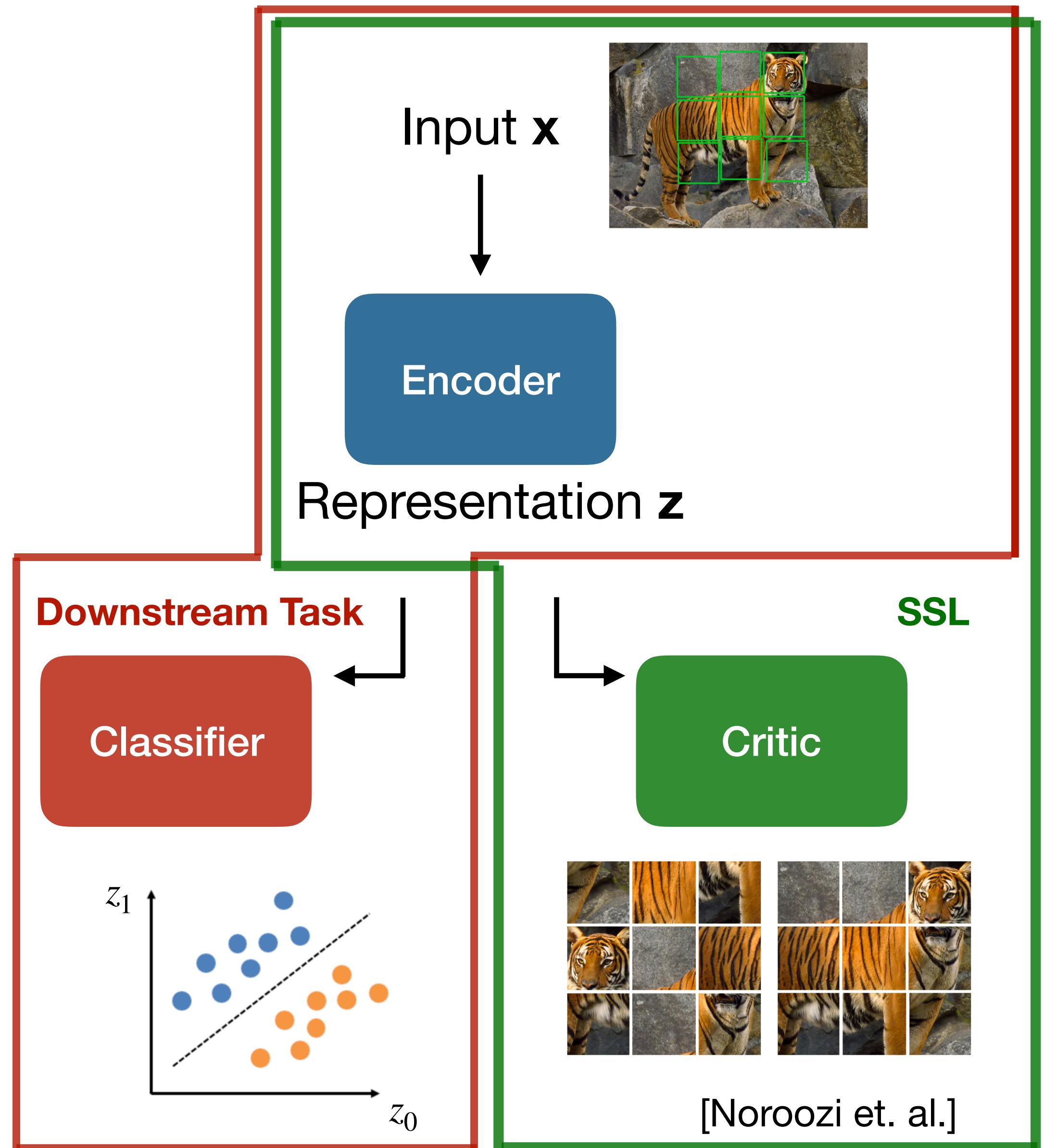
- Predict Future From Past
- Predict Adjacent Sections in Structured Data
- Predict Occluded Area from Non-Occluded One
- Undo Data-Augmentation Transformations



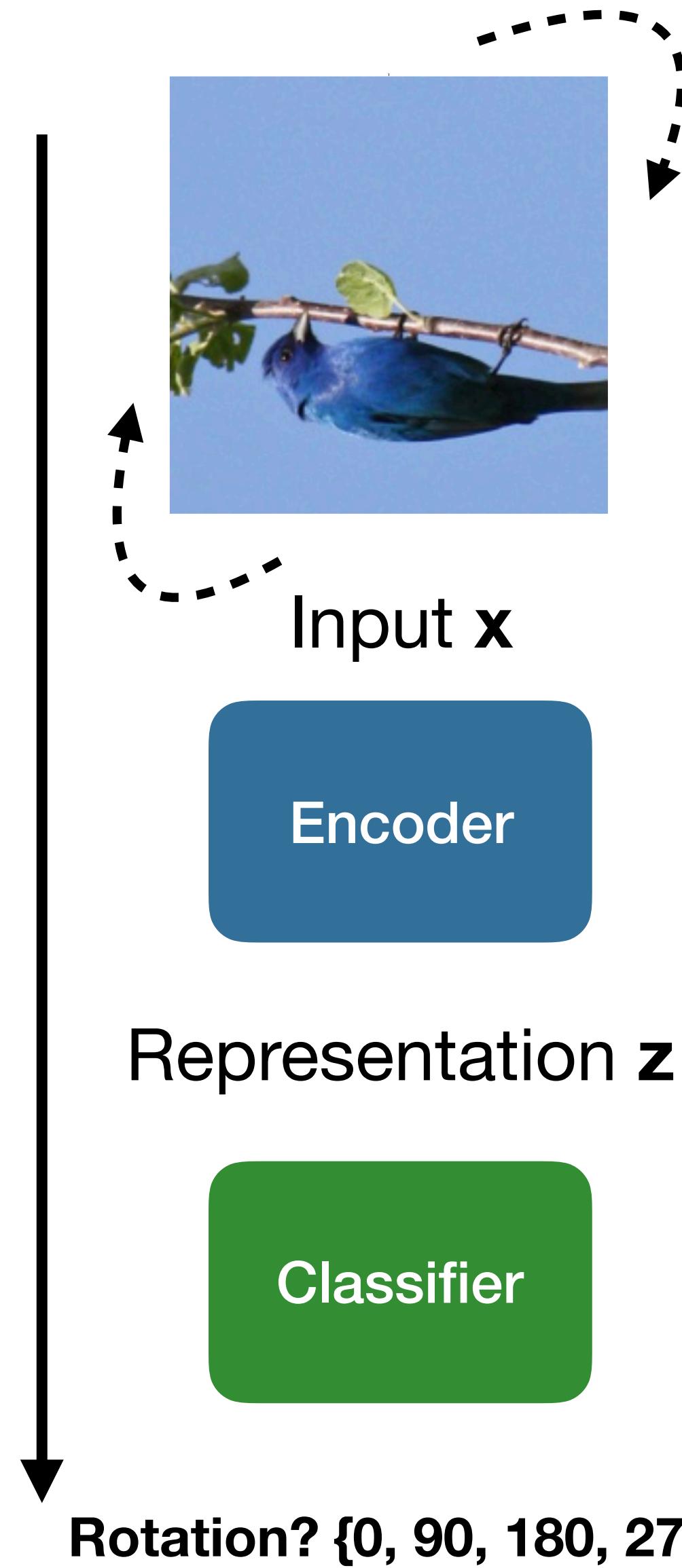
Typical SSL Approaches:

- Predict Future From Past
- Predict Adjacent Sections in Structured Data
- Predict Occluded Area from Non-Occluded One
- Undo Data-Augmentation Transformations

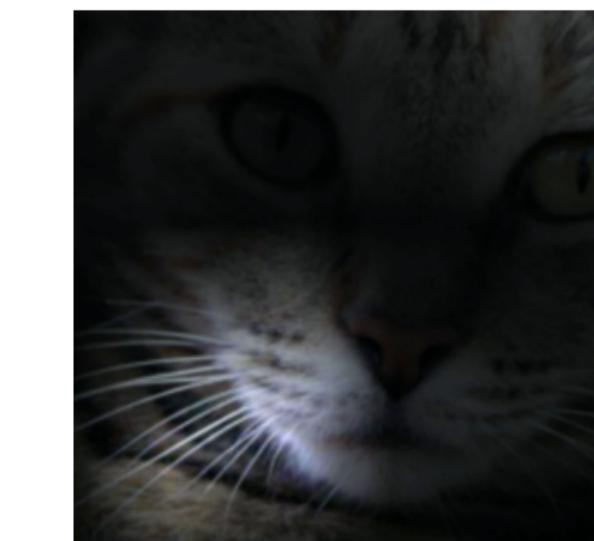
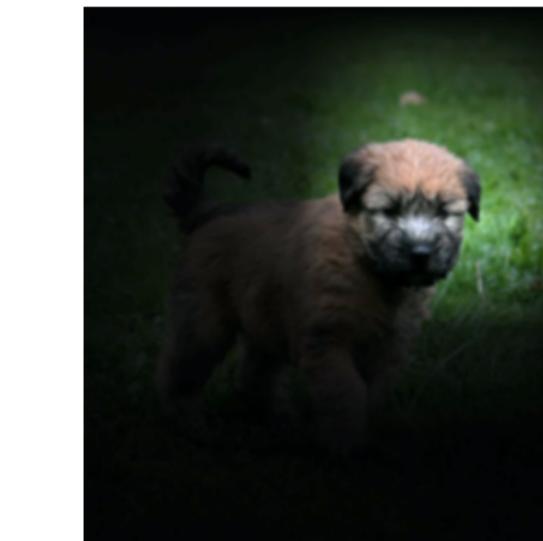
Obtain gains in downstream tasks (usually classification)



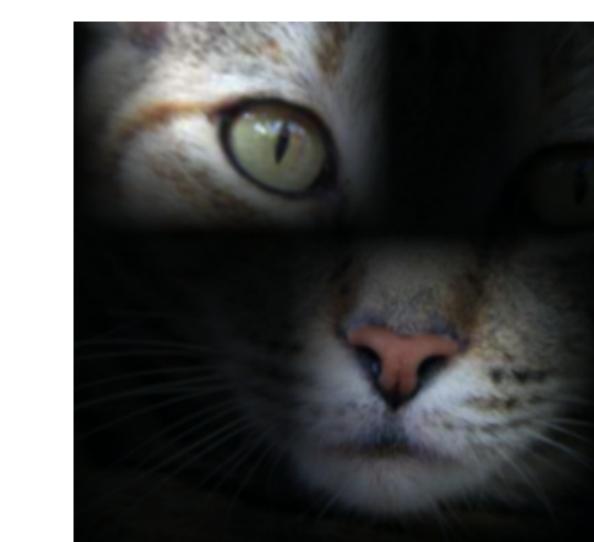
Predicting Rotations



CNNs are not rotation invariant.
Will need to learn part-whole relationships.

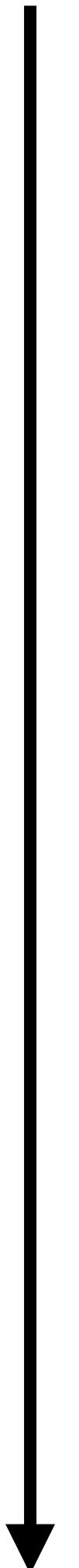


Supervised Attention



Self-Supervised Attention

Multimodal SSL: Audio + Video



Video x_v

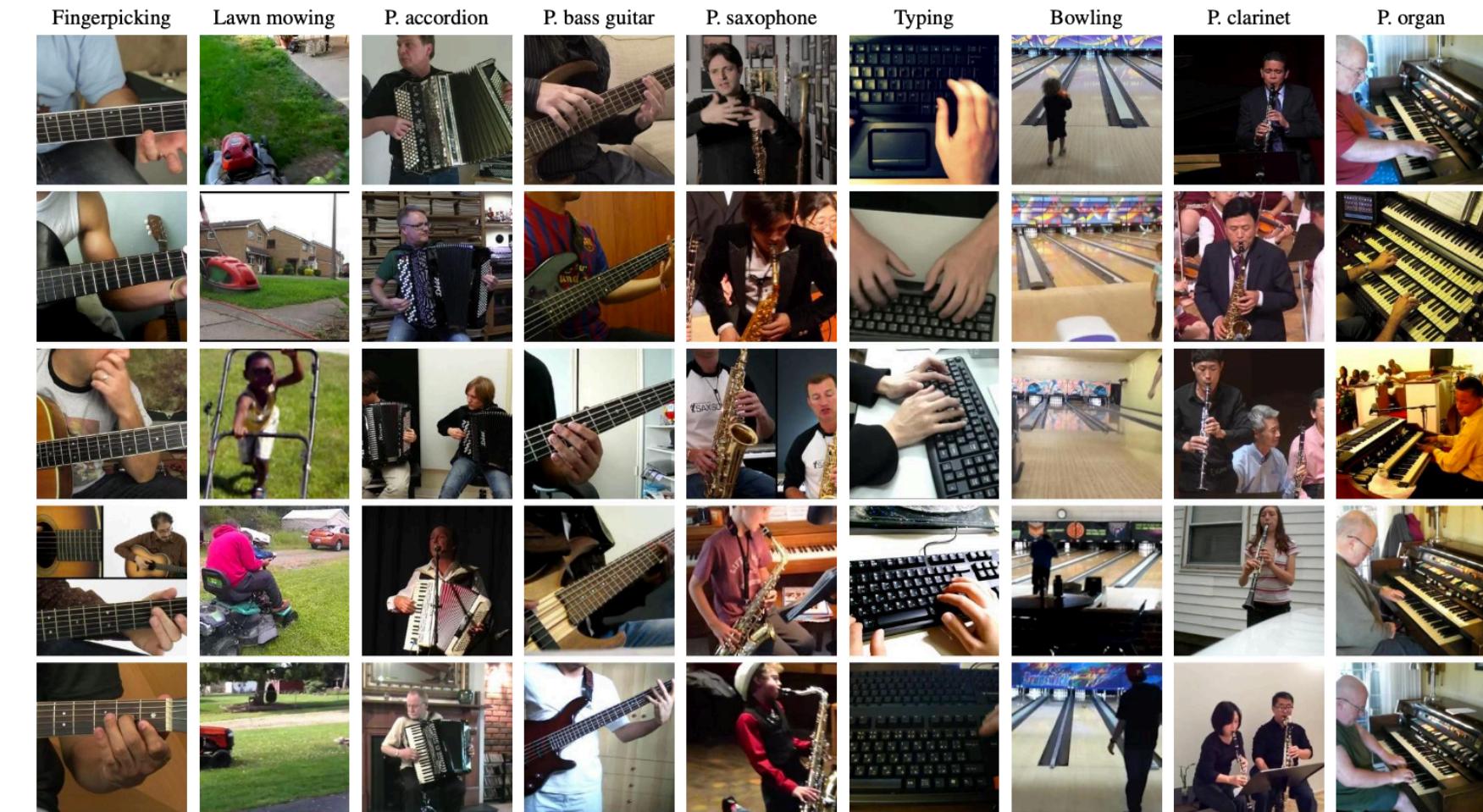


Encoder v

Audio x_a



Encoder v



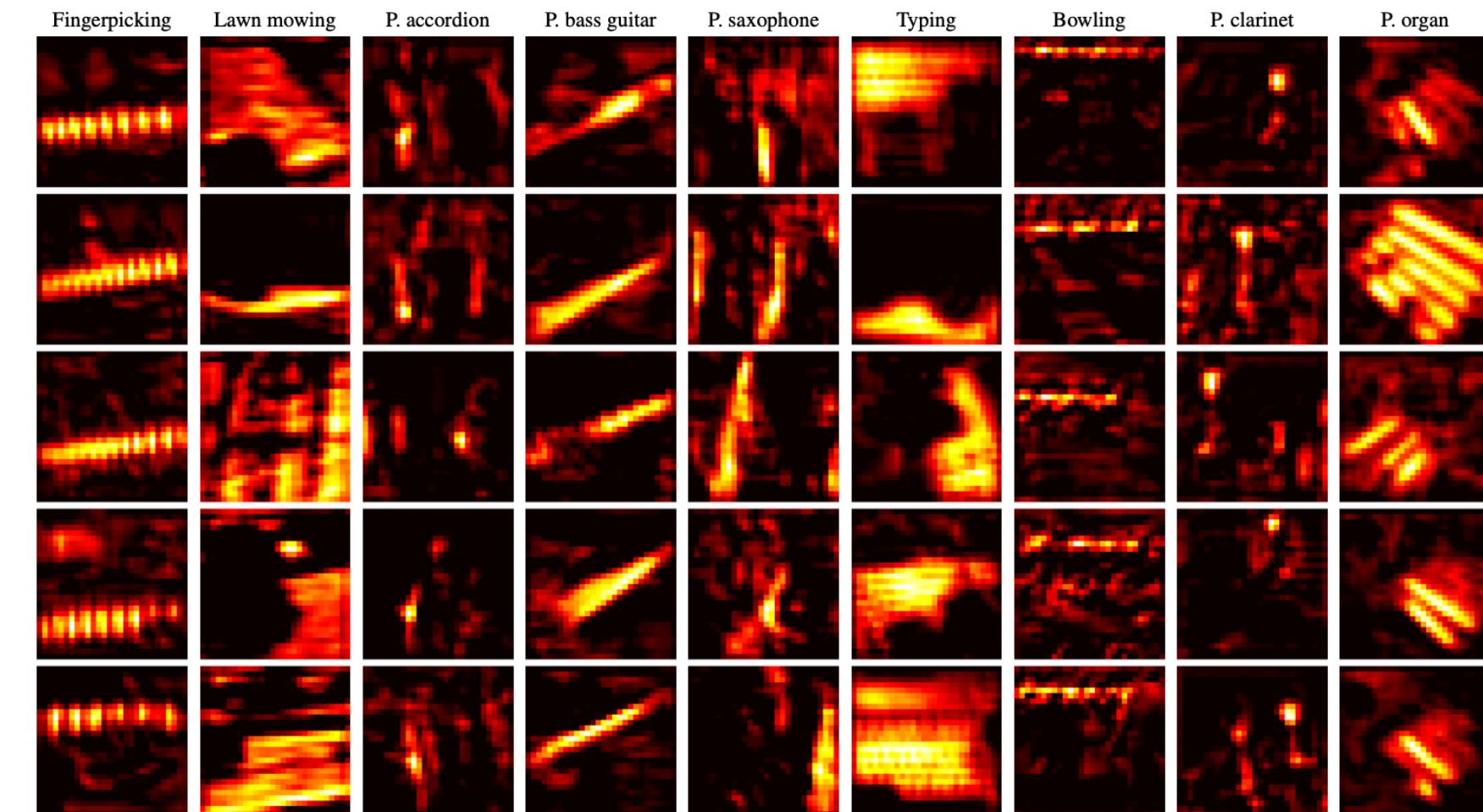
z_v

z_a

Classifier

a-v Correspondence?: {Y, N}

Separate $p(a, v)$ from $p(a)p(v)$



[Arandjelovic et. al.]

Word Embeddings

Input $\mathbf{x}^{(t)}$

Encoder

Embedding \mathbf{z}

NCE Critic

Unsupervised Analogical Reasoning

Newspapers			
New York	New York Times	Baltimore	Baltimore Sun
San Jose	San Jose Mercury News ■	Cincinnati	Cincinnati Enquirer ■
NHL Teams			
Boston	Boston Bruins	Montreal	Montreal Canadiens
Phoenix	Phoenix Coyotes ■	Nashville	Nashville Predators ■
NBA Teams			
Detroit	Detroit Pistons	Toronto	Toronto Raptors
Oakland	Golden State Warriors ■	Memphis	Memphis Grizzlies ■
Airlines			
Austria	Austrian Airlines	Spain	Spainair
Belgium	Brussels Airlines ■	Greece	Aegean Airlines ■
Company executives			
Steve Ballmer	Microsoft	Larry Page	Google
Samuel J. Palmisano	IBM ■	Werner Vogels	Amazon ■

Predict among negative samples:

$$\{\mathbf{x}^{(t-2)}, \mathbf{x}^{(t-1)}, \mathbf{x}^{(t+1)}, \mathbf{x}^{(t+2)}\}$$

Skip Gram Model: Separating $p(\mathbf{x}^{(t)}, \mathbf{x}^{(t+k)})$ from $p(\mathbf{x}^{(t)})p(\mathbf{x}^{(t+k)})$

[Mikolov et. al.]

Typical SSL Approaches:

- Predict Future From Past
- Predict Adjacent Sections in Structured Data
- Predict Occluded Area from Non-Occluded One
- Undo Data-Augmentation Transformations

Maximise Mutual Information Between Inputs and Representations?

$$I(\mathbf{x}, \mathbf{z}) = E_{p(\mathbf{x}, \mathbf{z})}[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})p(\mathbf{z})}]$$

But What is Mutual Information?

But What is Mutual Information?

$$I(\mathbf{x}, \mathbf{z}) = KL(p(\mathbf{x}, \mathbf{z}) || p(\mathbf{x})p(\mathbf{z})) = E_{p(\mathbf{x}, \mathbf{z})}[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})p(\mathbf{z})}]$$

But What is Mutual Information?

$$I(\mathbf{x}, \mathbf{z}) = KL(p(\mathbf{x}, \mathbf{z}) \parallel p(\mathbf{x})p(\mathbf{z})) = E_{p(\mathbf{x}, \mathbf{z})}[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})p(\mathbf{z})}]$$

$$I(\mathbf{x}, \mathbf{z}) = H(\mathbf{x}) - H(\mathbf{x} \mid \mathbf{z})$$

But What is Mutual Information?

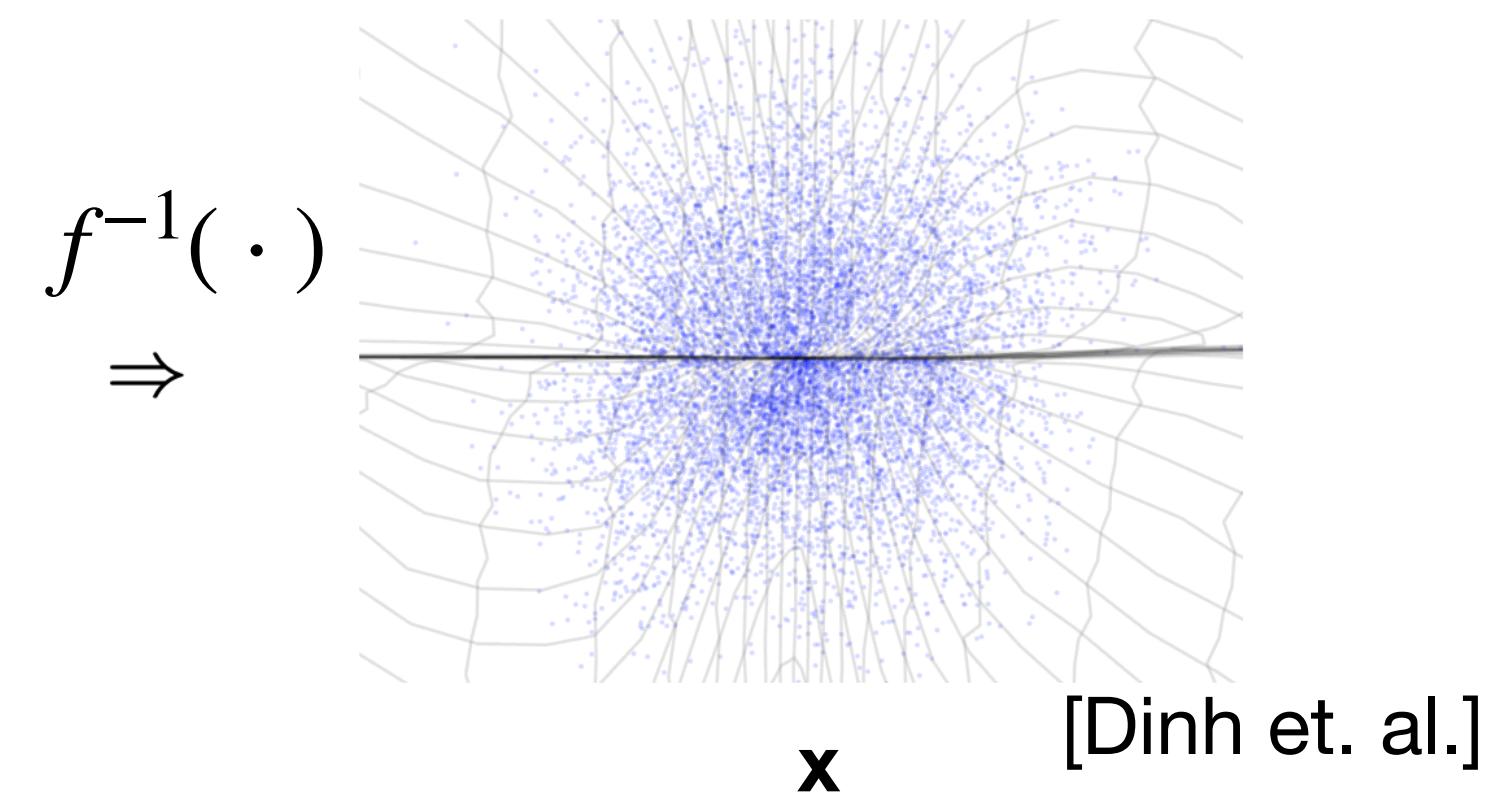
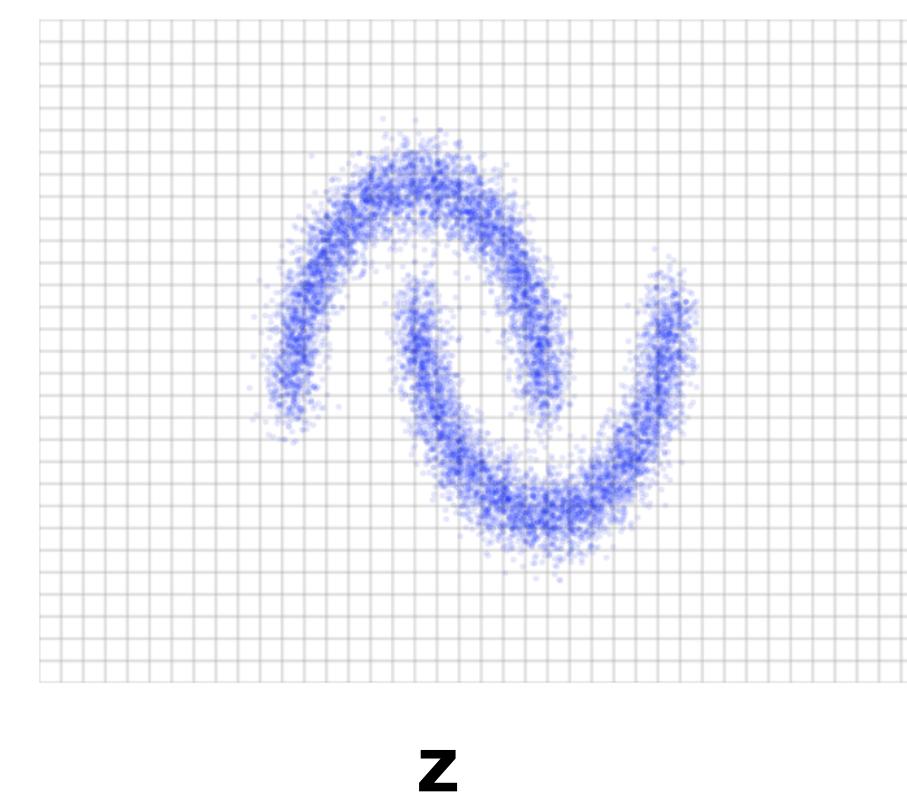
$$I(\mathbf{x}, \mathbf{z}) = KL(p(\mathbf{x}, \mathbf{z}) \parallel p(\mathbf{x})p(\mathbf{z})) = E_{p(\mathbf{x}, \mathbf{z})}[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})p(\mathbf{z})}]$$

$$I(\mathbf{x}, \mathbf{z}) = H(\mathbf{x}) - H(\mathbf{x} \mid \mathbf{z})$$

- Symmetric, Bounded: $0 \leq I(\mathbf{x}, \mathbf{y}) = I(\mathbf{y}, \mathbf{x}) \leq \min(H(\mathbf{x}), H(\mathbf{y}))$

- Invariant to Reparametrisations!

Here: $I(\mathbf{x}, \mathbf{z}) = H(\mathbf{x}) = H(\mathbf{z}) \longrightarrow$



But What is Mutual Information?

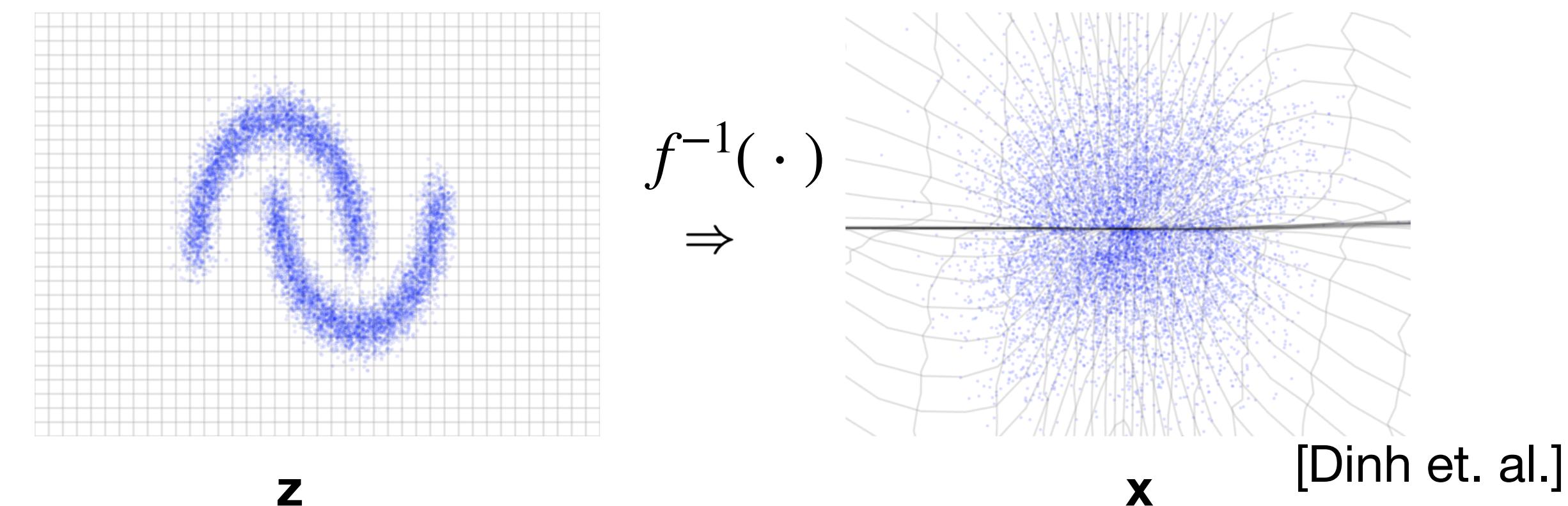
$$I(\mathbf{x}, \mathbf{z}) = KL(p(\mathbf{x}, \mathbf{z}) \parallel p(\mathbf{x})p(\mathbf{z})) = E_{p(\mathbf{x}, \mathbf{z})}[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})p(\mathbf{z})}]$$

$$I(\mathbf{x}, \mathbf{z}) = H(\mathbf{x}) - H(\mathbf{x} \mid \mathbf{z})$$

- Symmetric, Bounded: $0 \leq I(\mathbf{x}, \mathbf{y}) = I(\mathbf{y}, \mathbf{x}) \leq \min(H(\mathbf{x}), H(\mathbf{y}))$

- Invariant to Reparametrisations!

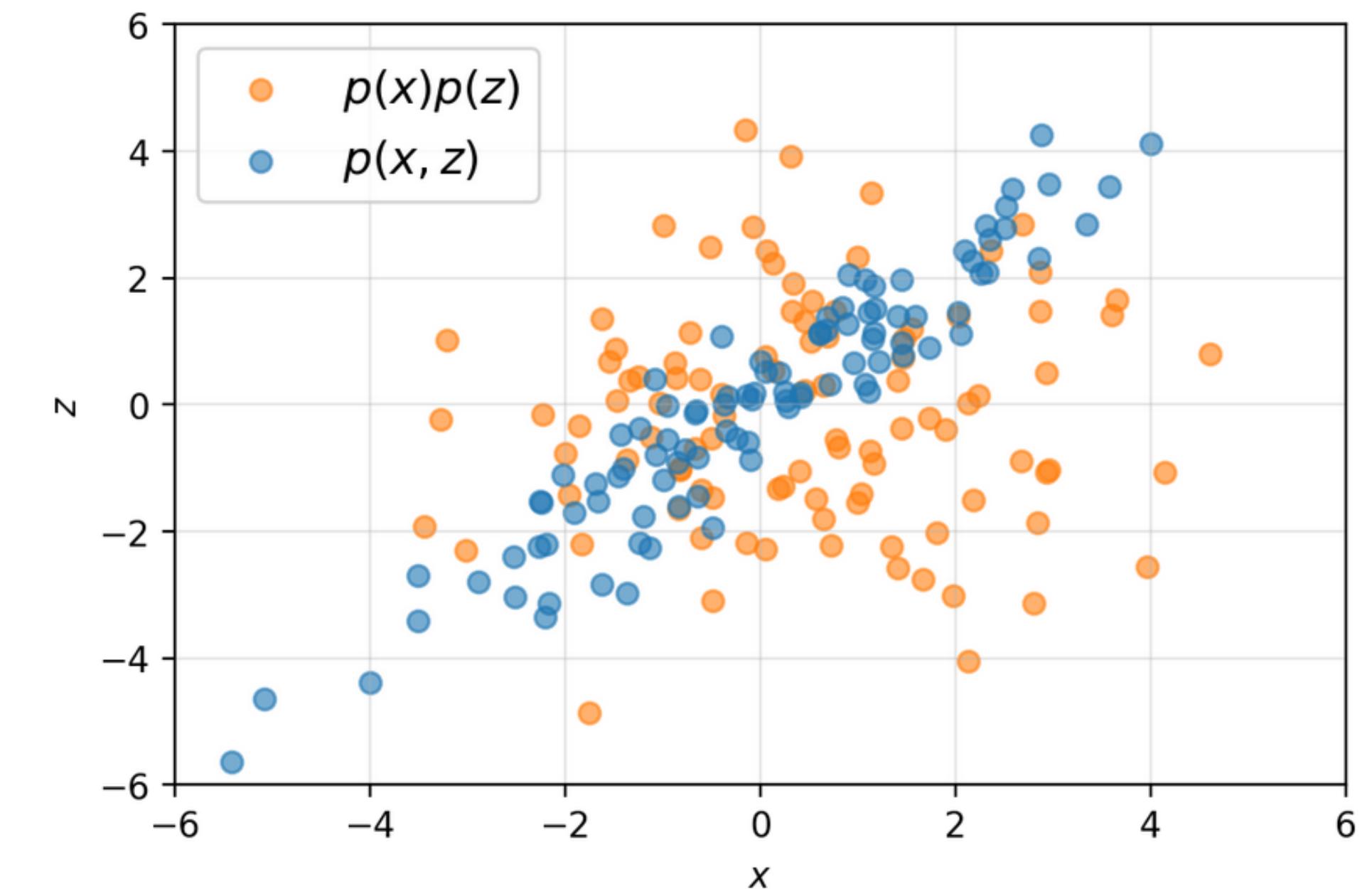
Here: $I(\mathbf{x}, \mathbf{z}) = H(\mathbf{x}) = H(\mathbf{z}) \longrightarrow$



ISSUE: No closed form for distributions. Only have samples!

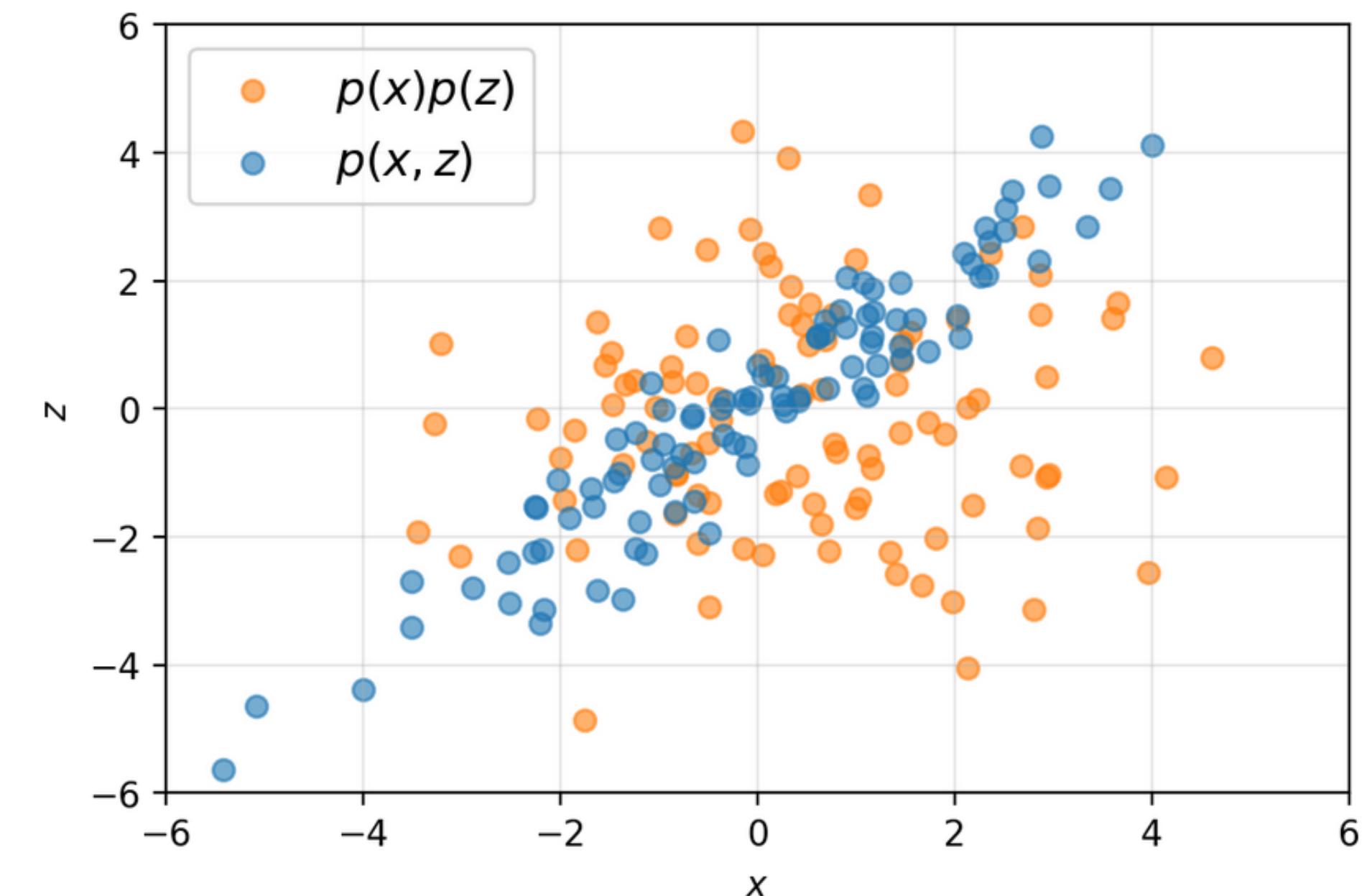
Density Ratio Estimation

$$I(\mathbf{x}, \mathbf{z}) = KL(p(\mathbf{x}, \mathbf{z}) \parallel p(\mathbf{x})p(\mathbf{z})) = E_{p(\mathbf{x}, \mathbf{z})}[\log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})p(\mathbf{z})}]$$



Density Ratio Estimation

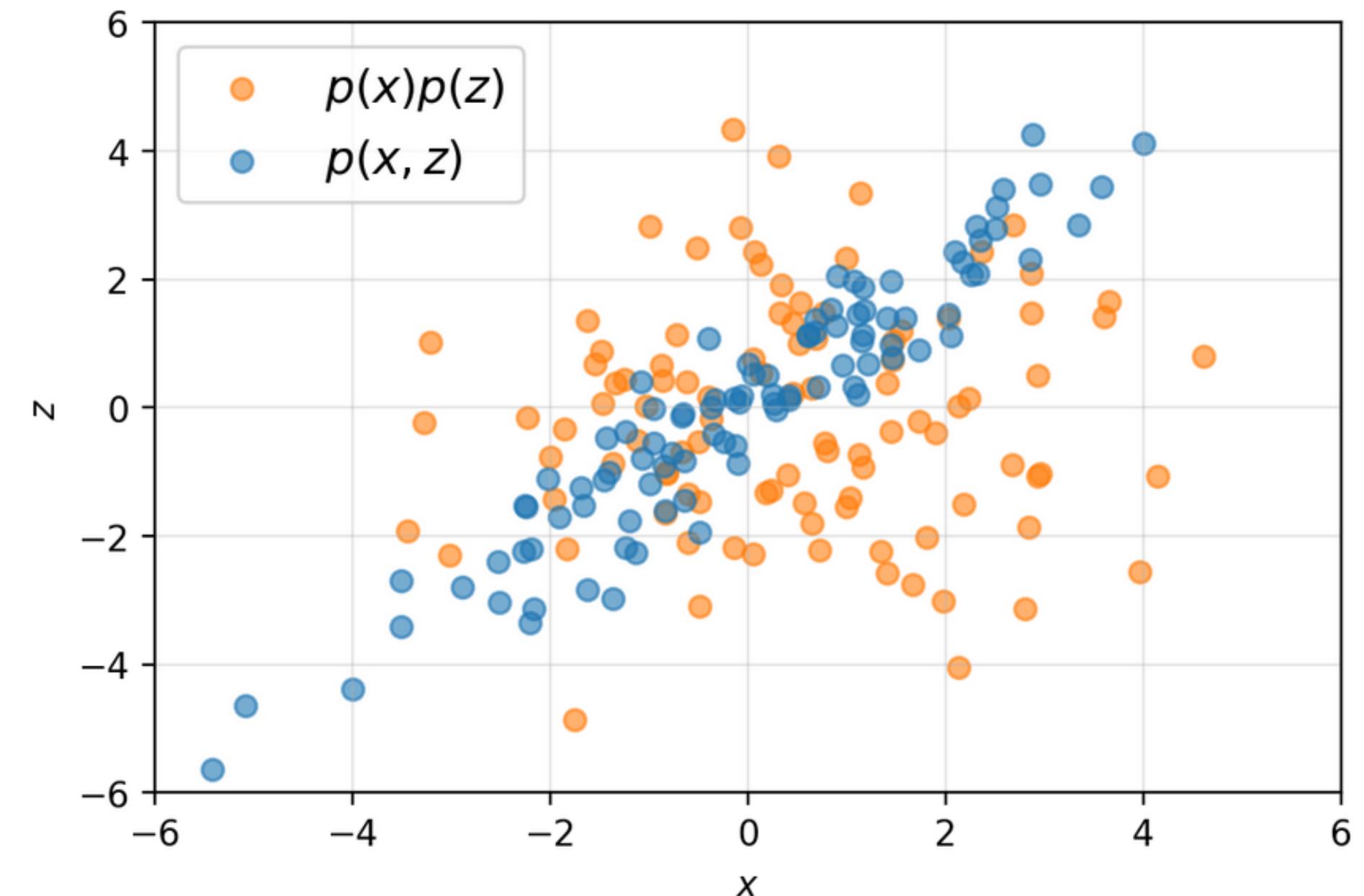
- Lets write $\mathbf{a} = (\mathbf{x}, \mathbf{z})$; $p(\mathbf{a} | b = 1) = p(\mathbf{x})p(\mathbf{z})$; $p(\mathbf{a} | b = 0) = p(\mathbf{x}, \mathbf{z})$



Density Ratio Estimation

- Lets write $\mathbf{a} = (\mathbf{x}, \mathbf{z})$; $p(\mathbf{a} | b = 1) = p(\mathbf{x})p(\mathbf{z})$; $p(\mathbf{a} | b = 0) = p(\mathbf{x}, \mathbf{z})$

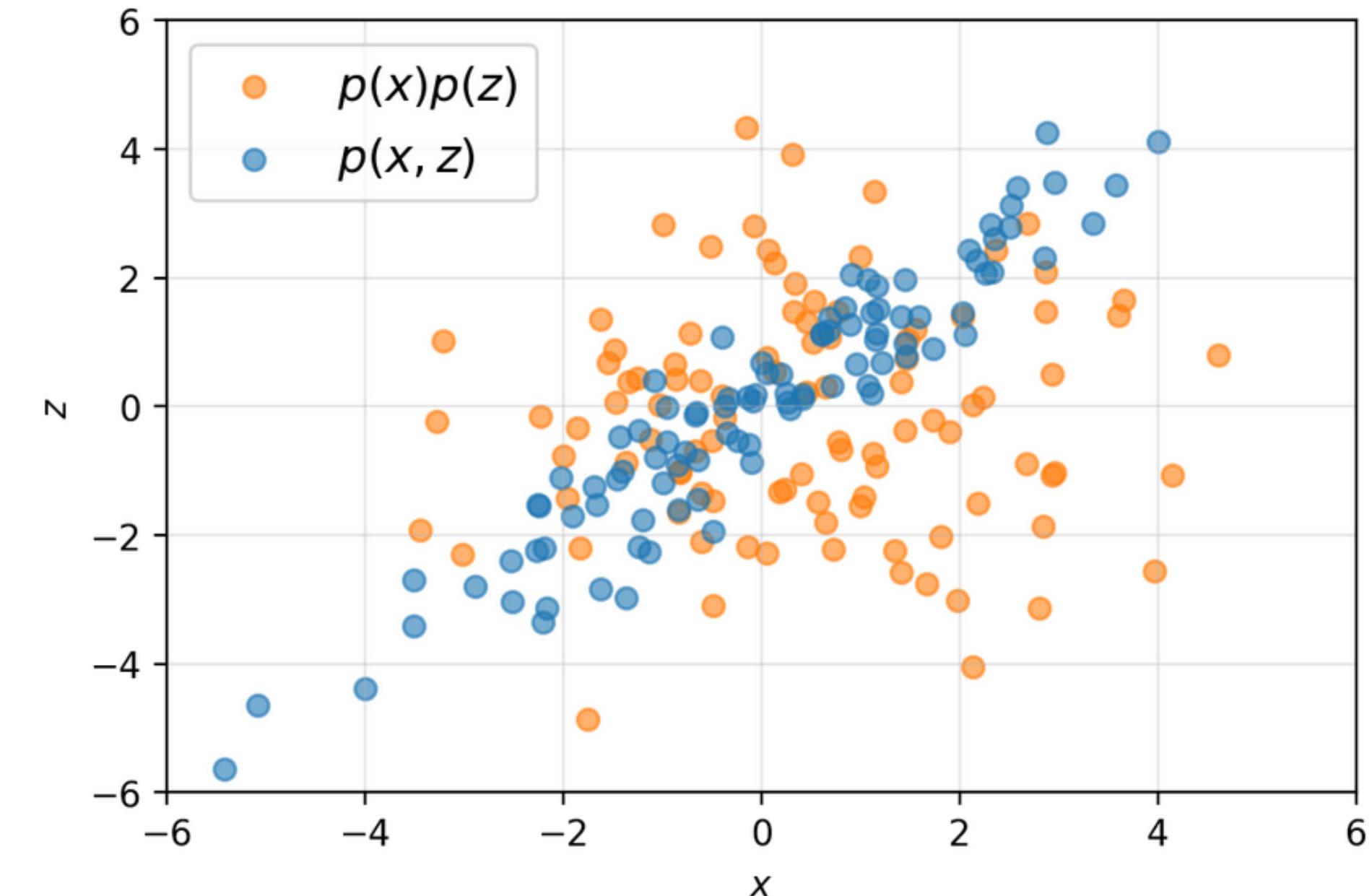
$$r(\mathbf{a}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})p(\mathbf{z})} = \frac{p(\mathbf{a} | b = 0)}{p(\mathbf{a} | b = 1)}$$



Density Ratio Estimation

- Lets write $\mathbf{a} = (\mathbf{x}, \mathbf{z})$; $p(\mathbf{a} | b = 1) = p(\mathbf{x})p(\mathbf{z})$; $p(\mathbf{a} | b = 0) = p(\mathbf{x}, \mathbf{z})$

$$r(\mathbf{a}) = \frac{p(\mathbf{a} | b = 0)}{p(\mathbf{a} | b = 1)} = \frac{p(b = 0 | \mathbf{a})p(\mathbf{a})p(b = 1)}{p(b = 1 | \mathbf{a})p(\mathbf{a})p(b = 0)} = \frac{p(b = 0 | \mathbf{a})}{p(b = 1 | \mathbf{a})} = \frac{p(b = 0 | \mathbf{a})}{1 - p(b = 0 | \mathbf{a})}$$



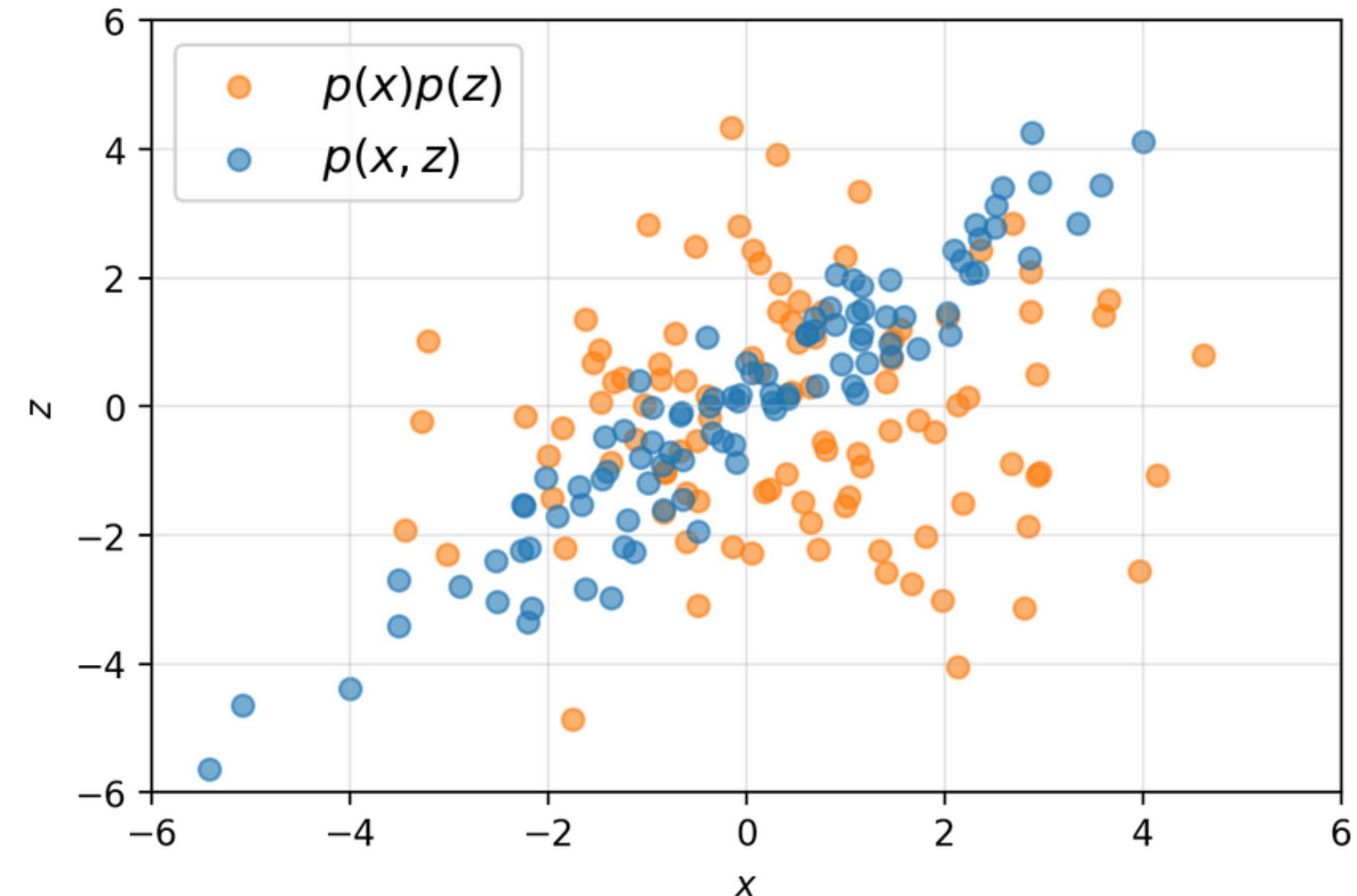
Density Ratio Estimation

- Lets write $\mathbf{a} = (\mathbf{x}, \mathbf{z})$; $p(\mathbf{a} | b = 1) = p(\mathbf{x})p(\mathbf{z})$; $p(\mathbf{a} | b = 0) = p(\mathbf{x}, \mathbf{z})$

$$r(\mathbf{a}) = \frac{p(\mathbf{a} | b = 0)}{p(\mathbf{a} | b = 1)} = \frac{p(b = 0 | \mathbf{a})p(\mathbf{a})p(b = 1)}{p(b = 1 | \mathbf{a})p(\mathbf{a})p(b = 0)} = \frac{p(b = 0 | \mathbf{a})}{p(b = 1 | \mathbf{a})} = \frac{p(b = 0 | \mathbf{a})}{1 - p(b = 0 | \mathbf{a})}$$

We often model conditional probabilities with parametric functions:

$$f_{NN}(\mathbf{a}) = p(b = 0 | \mathbf{a})$$



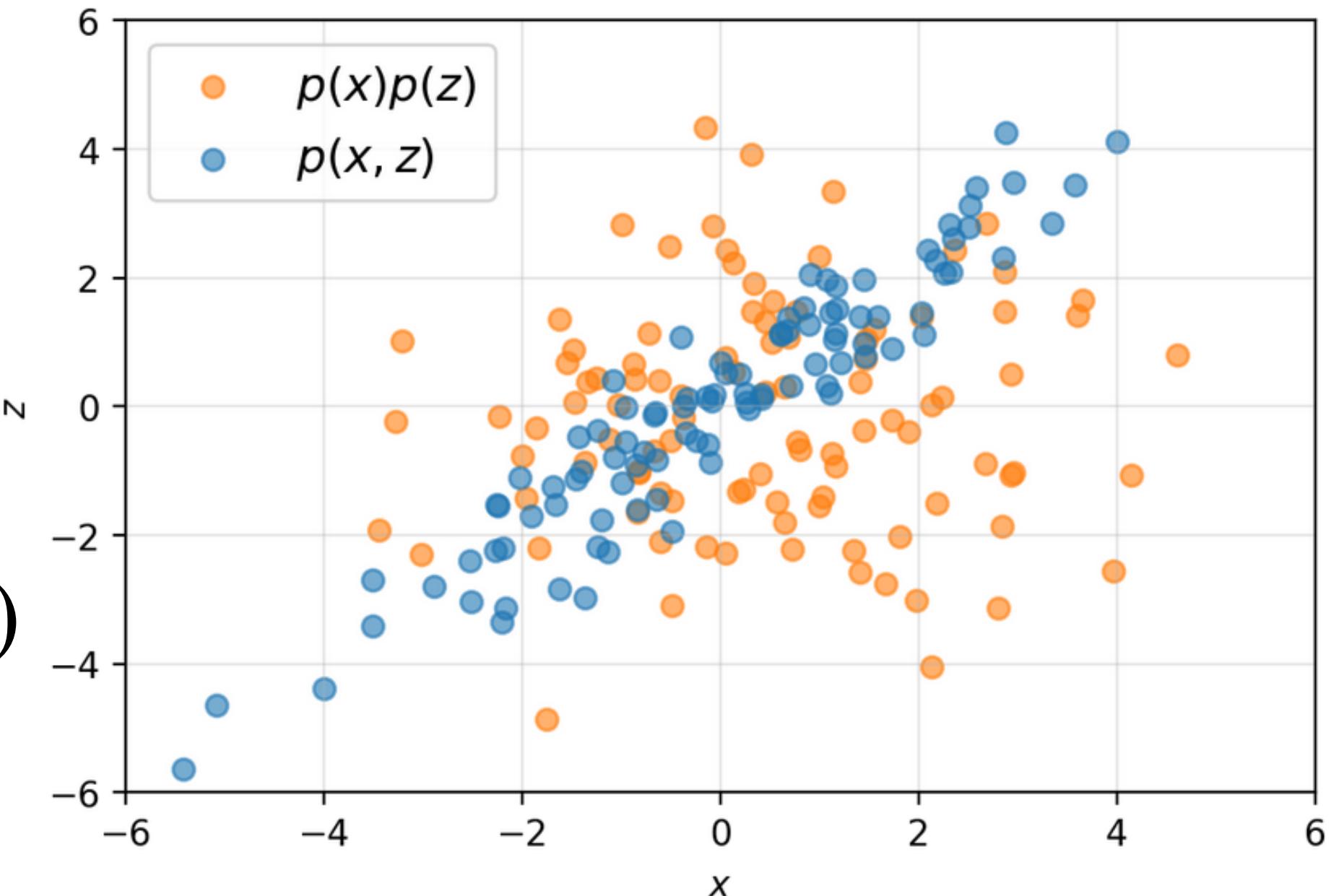
Density Ratio Estimation

- Lets write $\mathbf{a} = (\mathbf{x}, \mathbf{z})$; $p(\mathbf{a} | b = 1) = p(\mathbf{x})p(\mathbf{z})$; $p(\mathbf{a} | b = 0) = p(\mathbf{x}, \mathbf{z})$

$$r(\mathbf{a}) = \frac{p(\mathbf{a} | b = 0)}{p(\mathbf{a} | b = 1)} = \frac{p(b = 0 | \mathbf{a})p(\mathbf{a})p(b = 1)}{p(b = 1 | \mathbf{a})p(\mathbf{a})p(b = 0)} = \frac{p(b = 0 | \mathbf{a})}{p(b = 1 | \mathbf{a})} = \frac{p(b = 0 | \mathbf{a})}{1 - p(b = 0 | \mathbf{a})}$$

We can learn a parametric function
that estimates log density ratios from samples!

$$f_{NN}(\mathbf{a}) = \log r(\mathbf{a}) = \log \frac{p(b = 0 | \mathbf{a})}{1 - p(b = 0 | \mathbf{a})} = \log \left(\frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})p(\mathbf{z})} \right)$$



Representation Learning through Maximising $I(x, z)$

1. Sample from joint distribution:

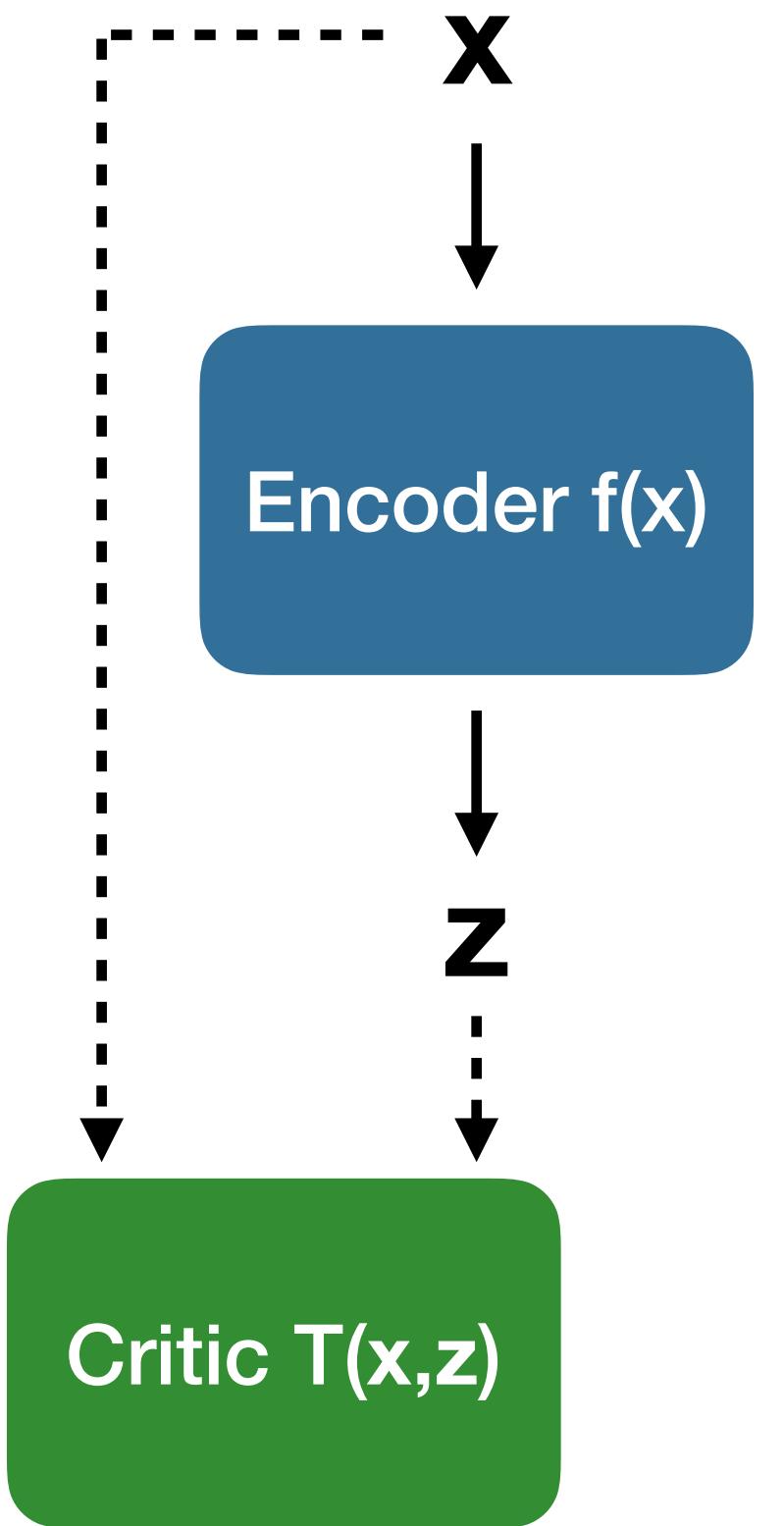
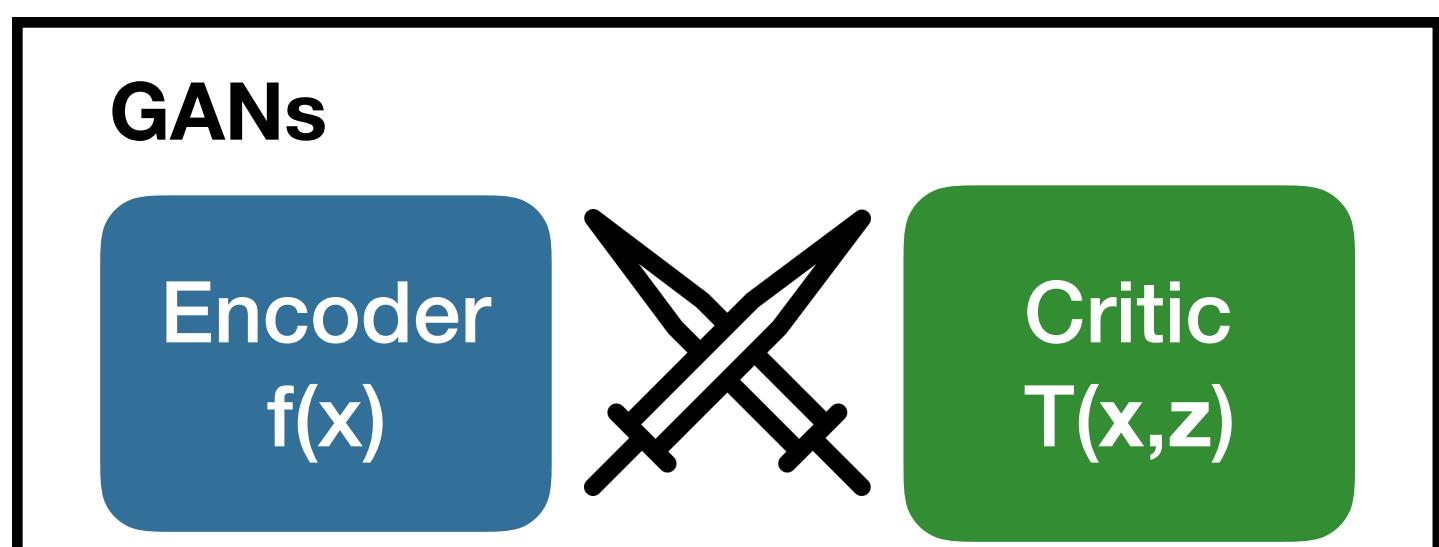
$$(x, z) \sim p(x, z) \quad \text{as} \quad x \sim p_{\mathcal{D}}(x); \quad z = f_{\phi}(x)$$

2. Negative sample from factorised distribution:

$$(x, z) \sim p(x)p(z) \quad \text{as} \quad x, x' \sim p_{\mathcal{D}}(x); \quad z = f_{\phi}(x')$$

3. Estimate $I(x, z)$ with $T_{\theta}(x, z)$

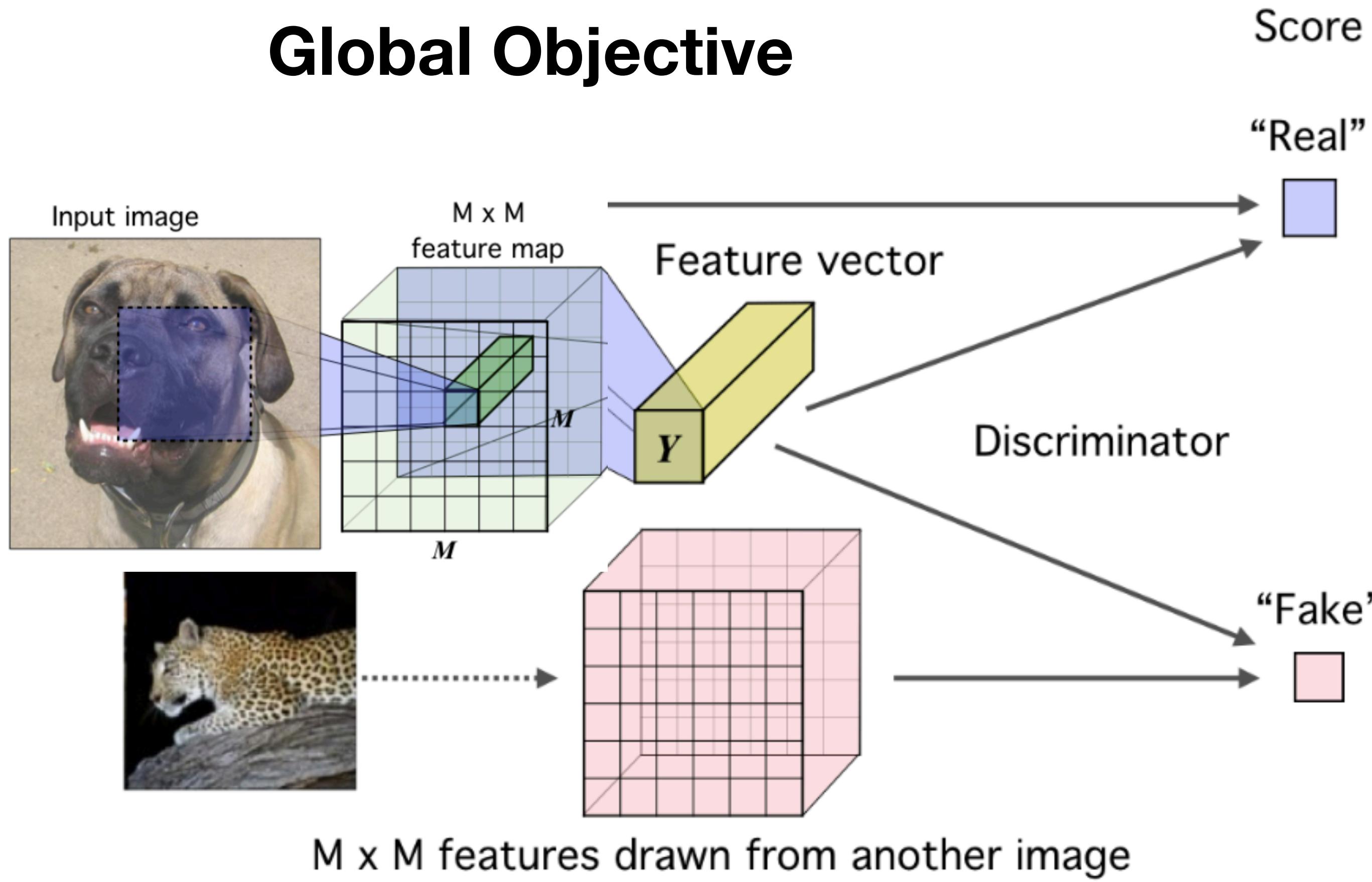
4. Optimise $\arg \max_{\theta, \phi} I(x, z)$



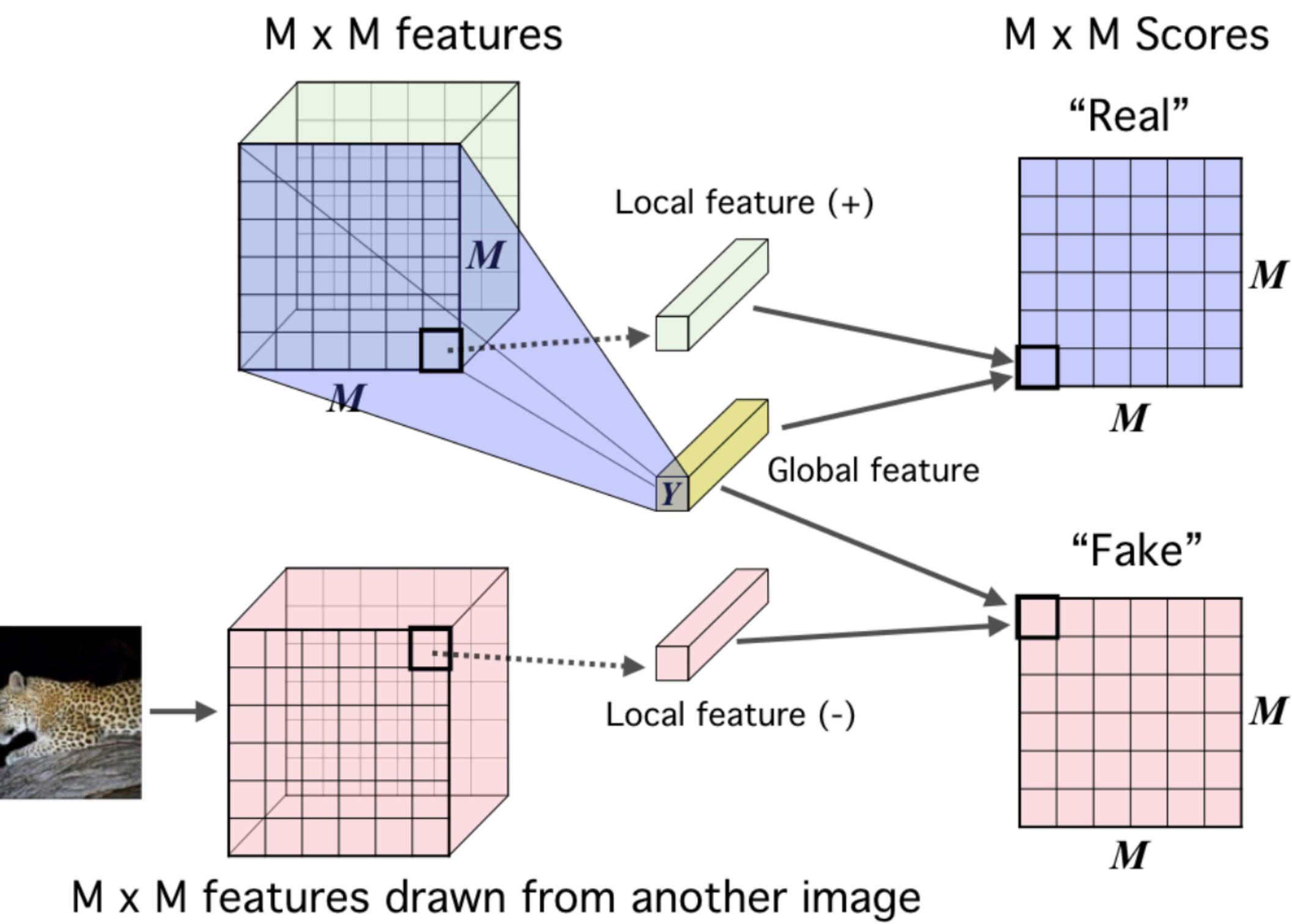
$$T_{\theta}(x, z) \approx \log \frac{p(x, z)}{p(x)p(z)}$$

Deep InfoMax

Global Objective



Local Objective



[Hjelm et. al.]

Deep InfoMax (DIM)

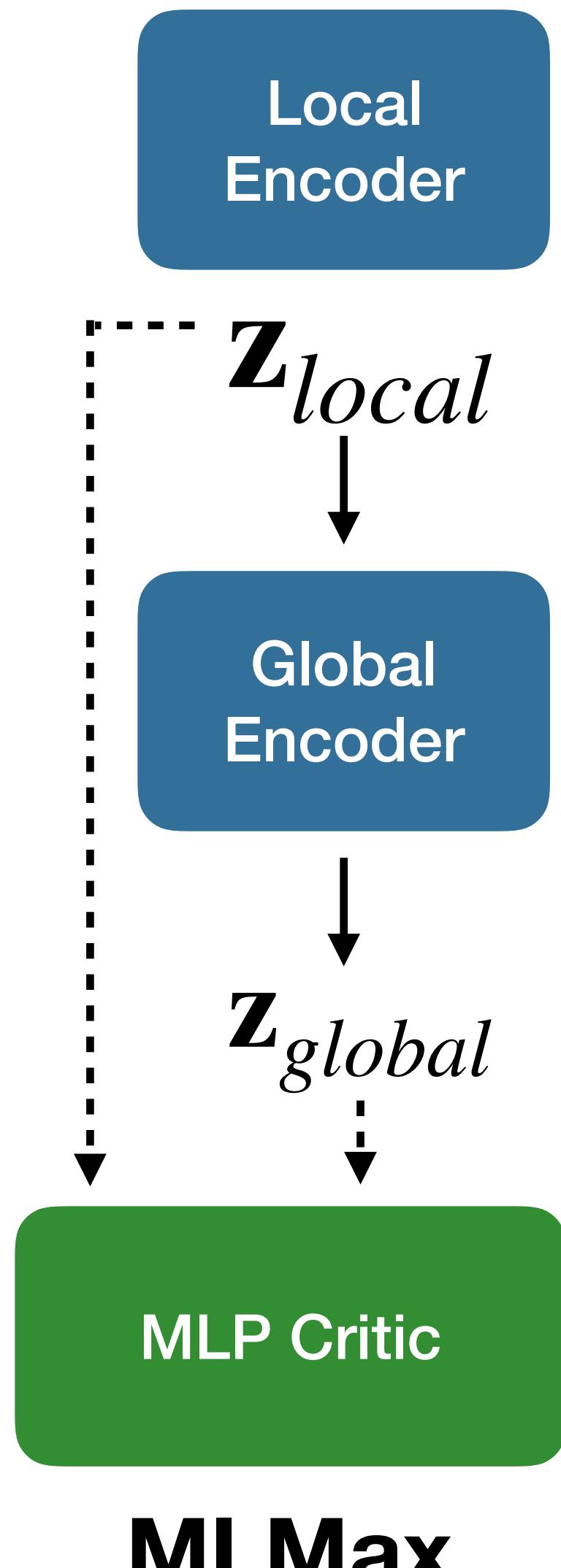
- Max MI between Local and Global Representations
- Critic is an MLP:

Operation	Size	Activation
Input → Linear layer	512	ReLU
Linear layer	512	ReLU
Linear layer	1	

*

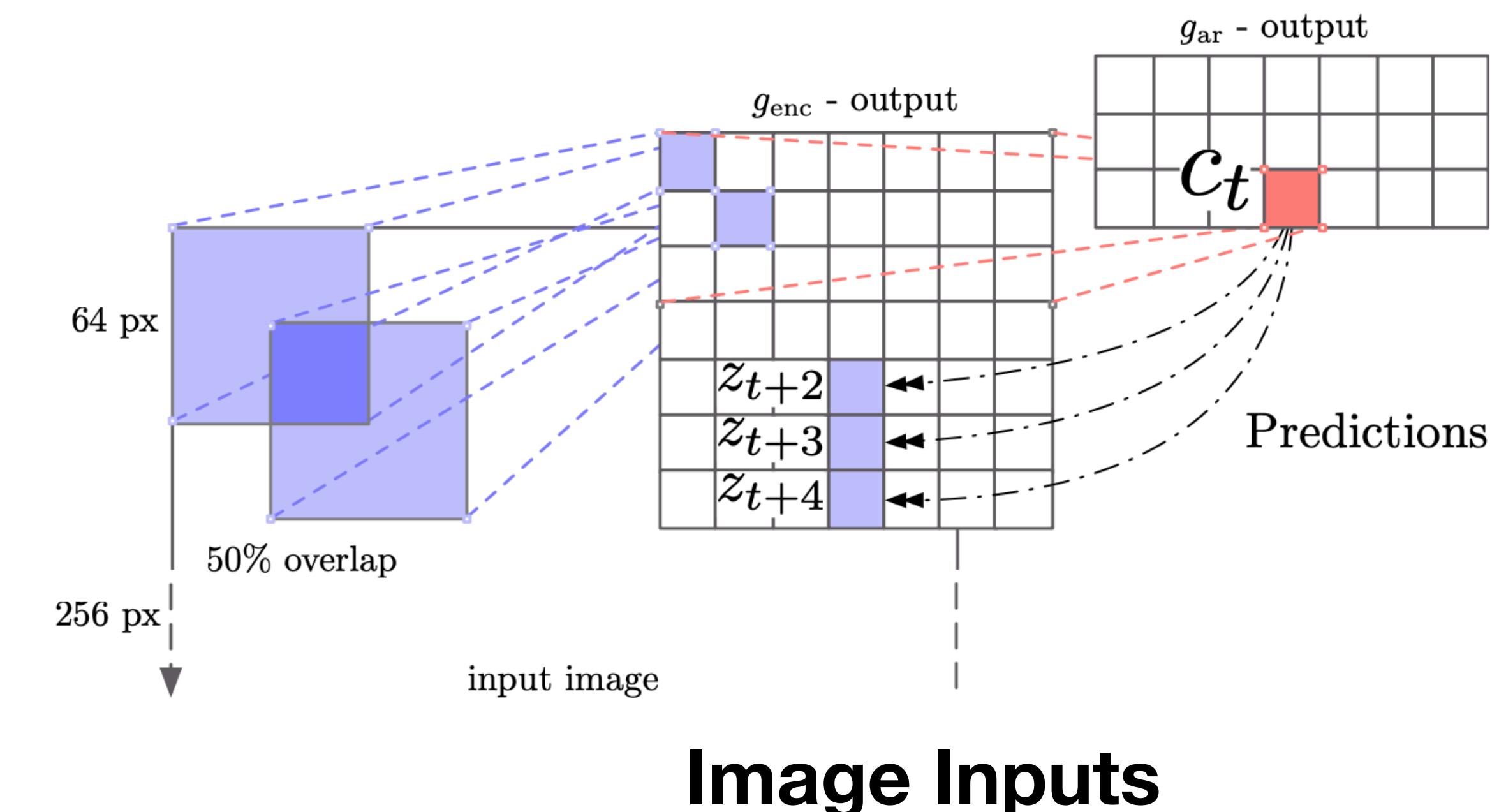
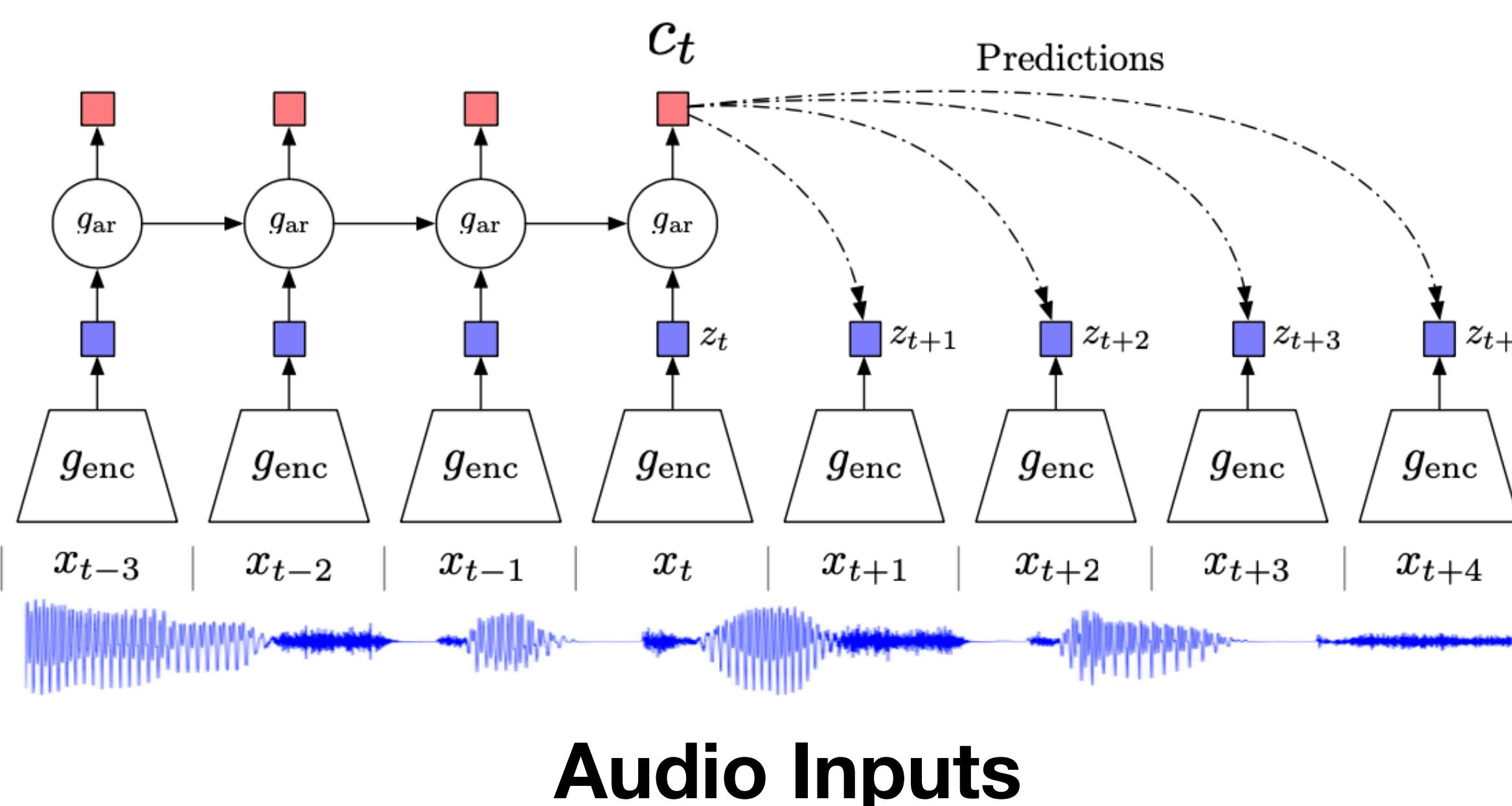
Linear SVM
Classification Results:

Model	CIFAR10			CIFAR100		
	conv	fc (1024)	Y(64)	conv	fc (1024)	Y(64)
Fully supervised		75.39			42.27	
VAE	60.71	60.54	54.61	37.21	34.05	24.22
AE	62.19	55.78	54.47	31.50	23.89	27.44
β -VAE	62.4	57.89	55.43	32.28	26.89	28.96
AAE	59.44	57.19	52.81	36.22	33.38	23.25
BiGAN	62.57	62.74	52.54	37.59	33.34	21.49
NAT	56.19	51.29	31.16	29.18	24.57	9.72
DIM(G)	52.2	52.84	43.17	27.68	24.35	19.98
DIM(L) (DV)	72.66	70.60	64.71	48.52	44.44	39.27
DIM(L) (JSD)	73.25	73.62	66.96	48.13	45.92	39.60
DIM(L) (infoNCE)	75.21	75.57	69.13	49.74	47.72	41.61



Contrastive Predictive Coding

- Max MI between representations of spatially/temporally similar inputs
- $\max I(\mathbf{z}^{(t)}, \mathbf{z}^{(t+k)})$ with **(AR NN + Bilinear)** Critic $T(\mathbf{z}^{(1\dots k)}, \mathbf{z}^{(k+j)}) = g_{AR}(\mathbf{z}^{(1\dots k)})^\top W^{(j)} \mathbf{z}^{(k+j)}$



Audio-Video matching and Word2Vec fit within this framework!

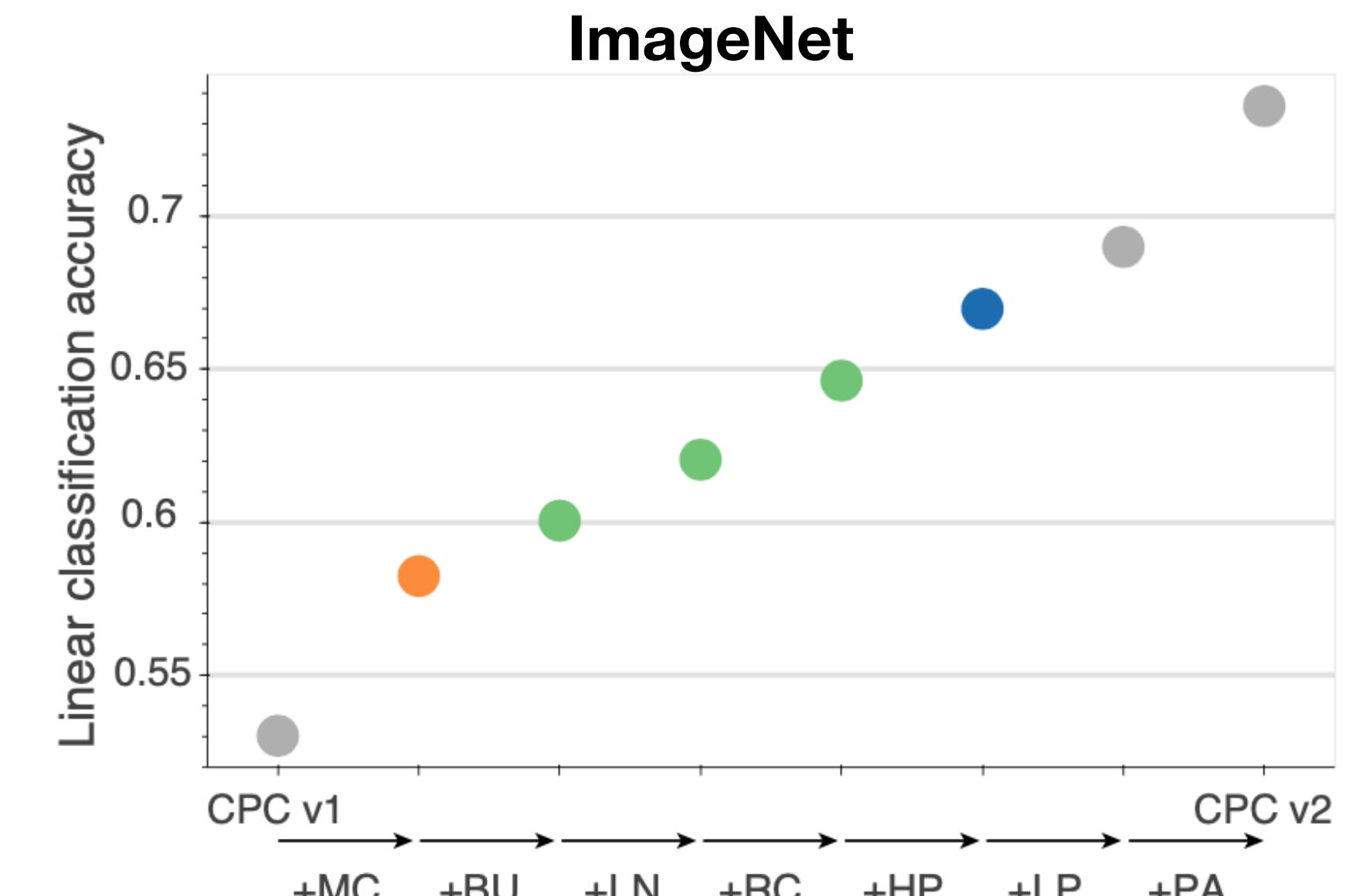
[van den Oord et. al.]

Better Contrastive Predictive Coding

- Large Batch Sizes
- Heavy Patch Augmentation
- Predict in every direction (not just forward)
- Larger Capacity NNs than when using labels
- Layer Norm (**Not Batch Norm**)

Text Document Classification

Method	MR	CR	Subj	MPQA	TREC
Paragraph-vector [40]	74.8	78.1	90.5	74.2	91.8
Skip-thought vector [26]	75.5	79.3	92.1	86.9	91.4
Skip-thought + LN [41]	79.5	82.6	93.4	89.0	-
CPC	76.9	80.1	91.2	87.7	96.8



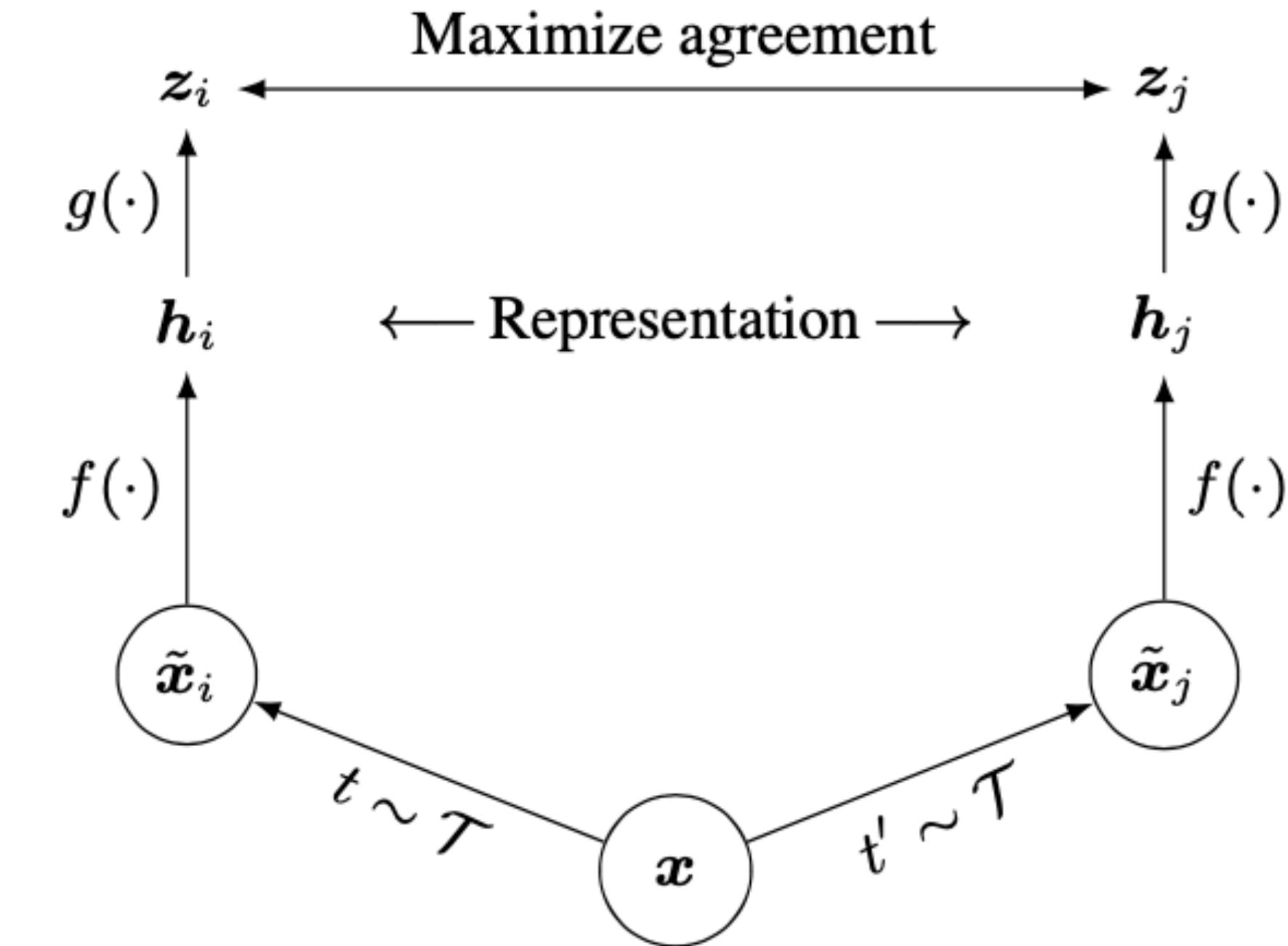
A Unifying Framework for Contrastive Learning

[Chen et. al.]

A Unifying Framework for Contrastive Learning

- Use ~InfoNCE to **max MI between augmented inputs**

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}$$

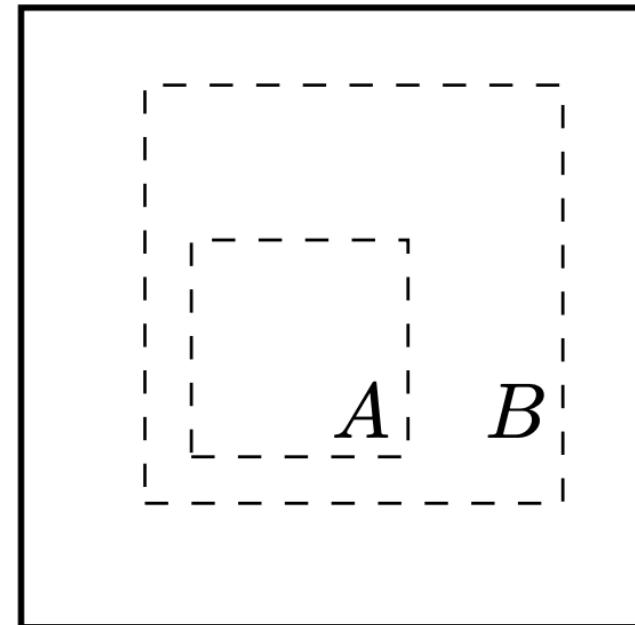


A Unifying Framework for Contrastive Learning

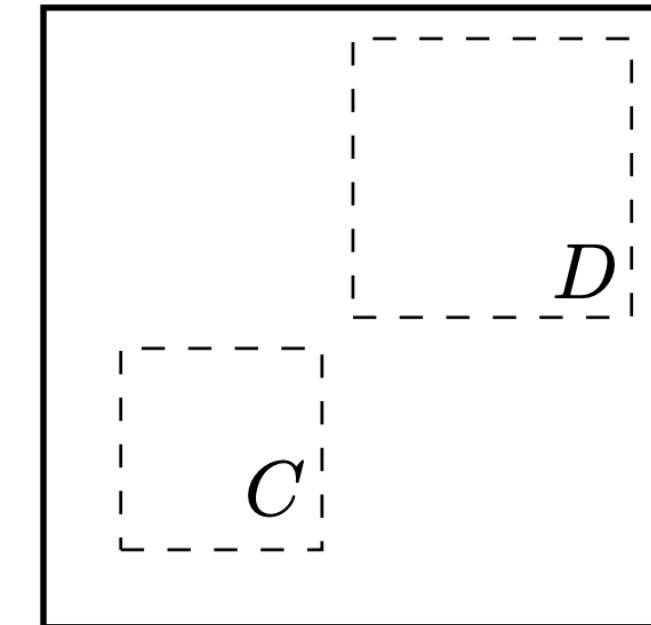
- Use \sim InfoNCE to **max MI between augmented inputs**

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}$$

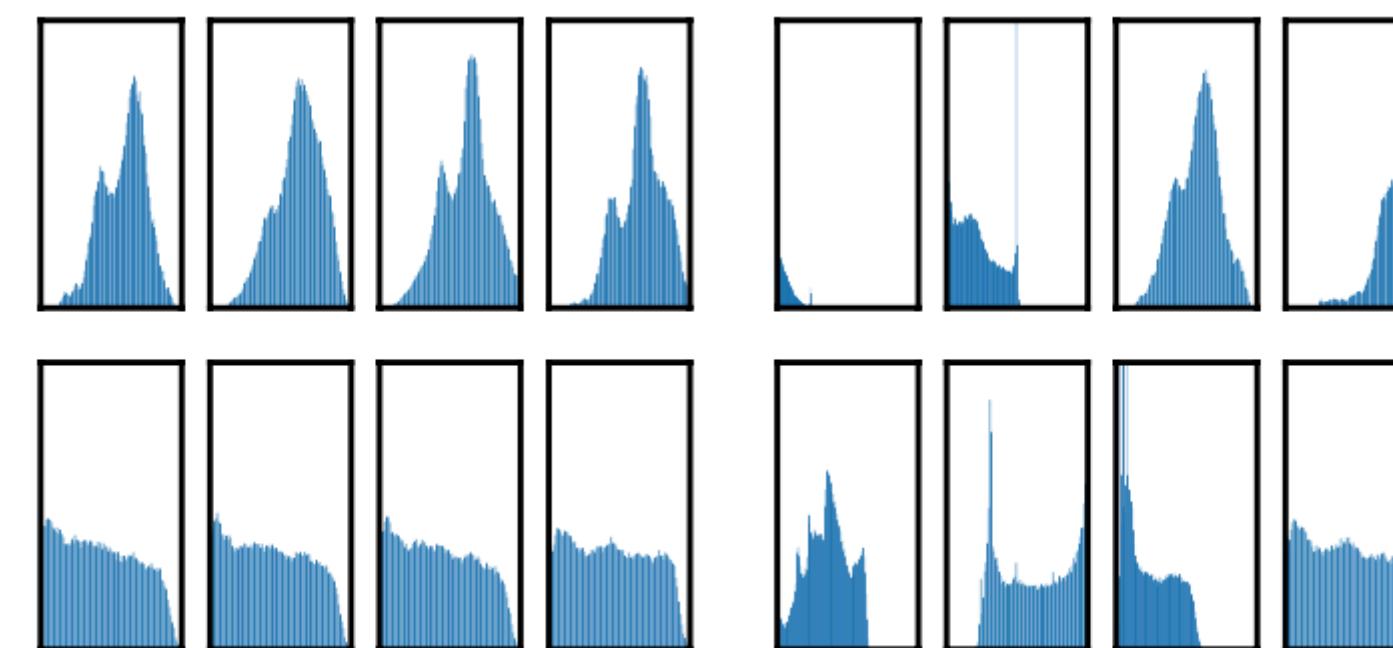
- Random Crop + Color Distortion Augmentation



(a) Global and local views.

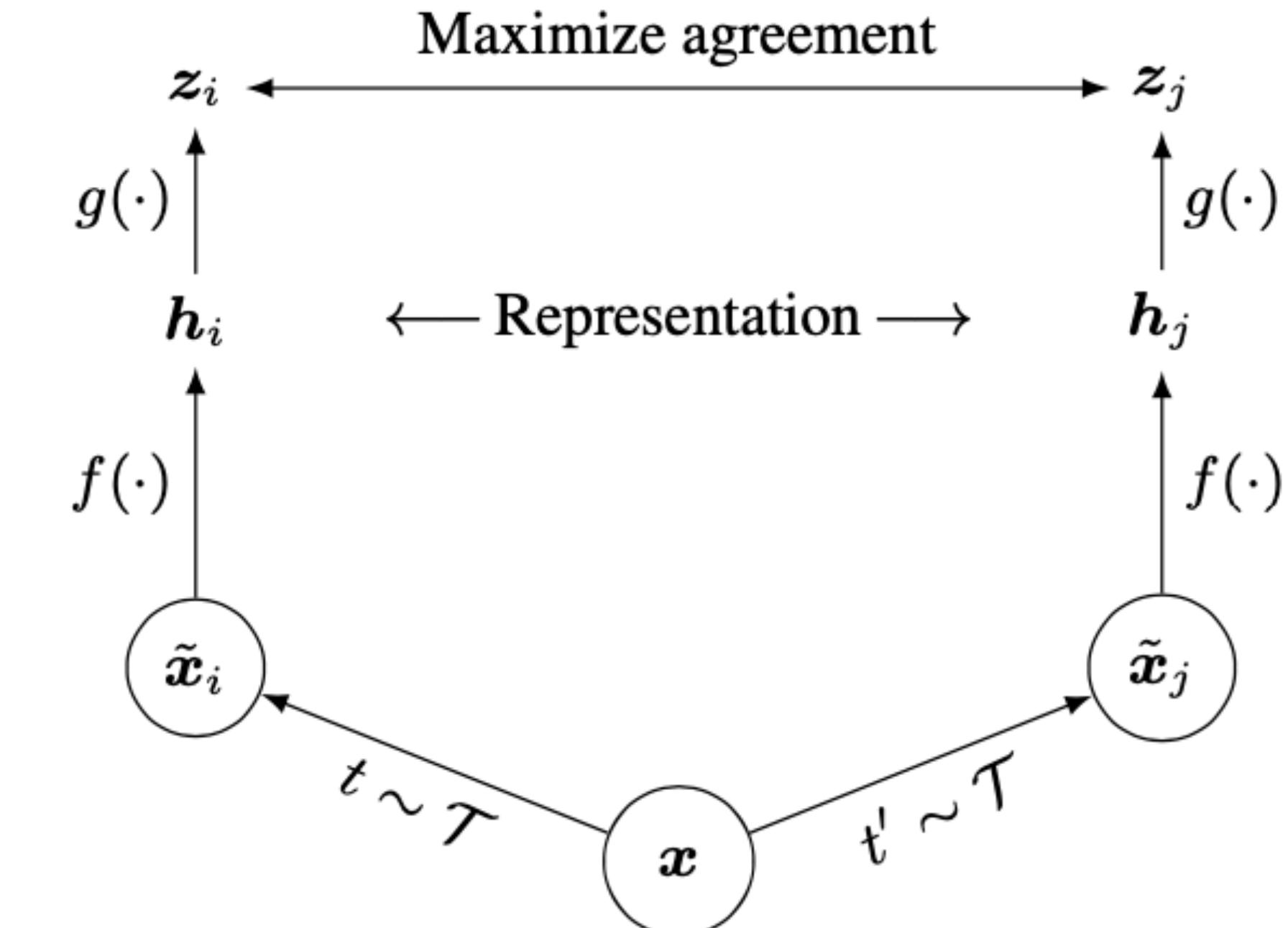


(b) Adjacent views.



(a) Without color distortion.

(b) With color distortion.

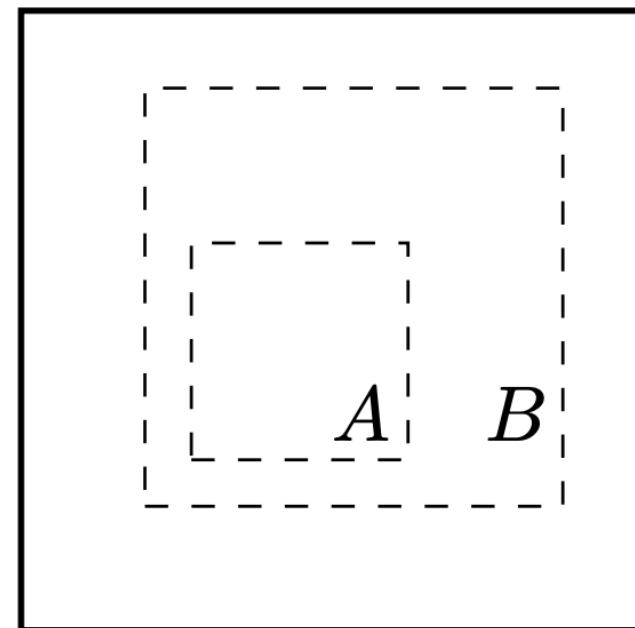


A Unifying Framework for Contrastive Learning

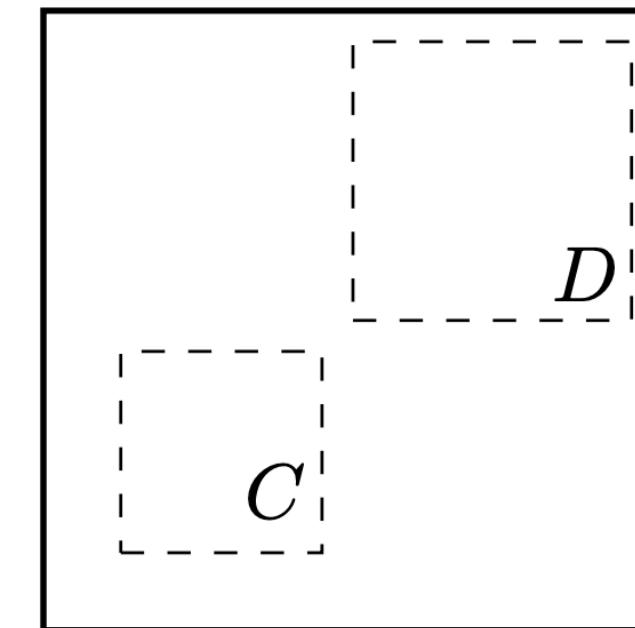
- Use \sim InfoNCE to **max MI between augmented inputs**

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k)/\tau)}$$

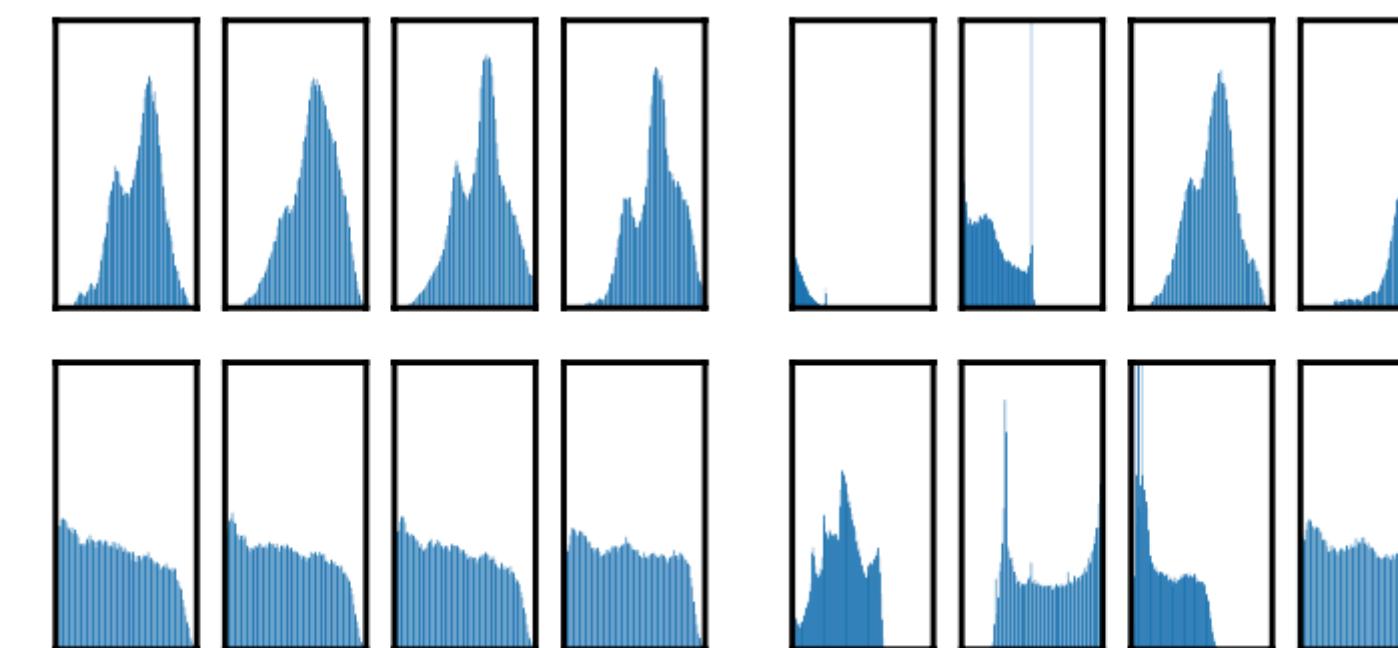
- Random Crop + Color Distortion Augmentation



(a) Global and local views.

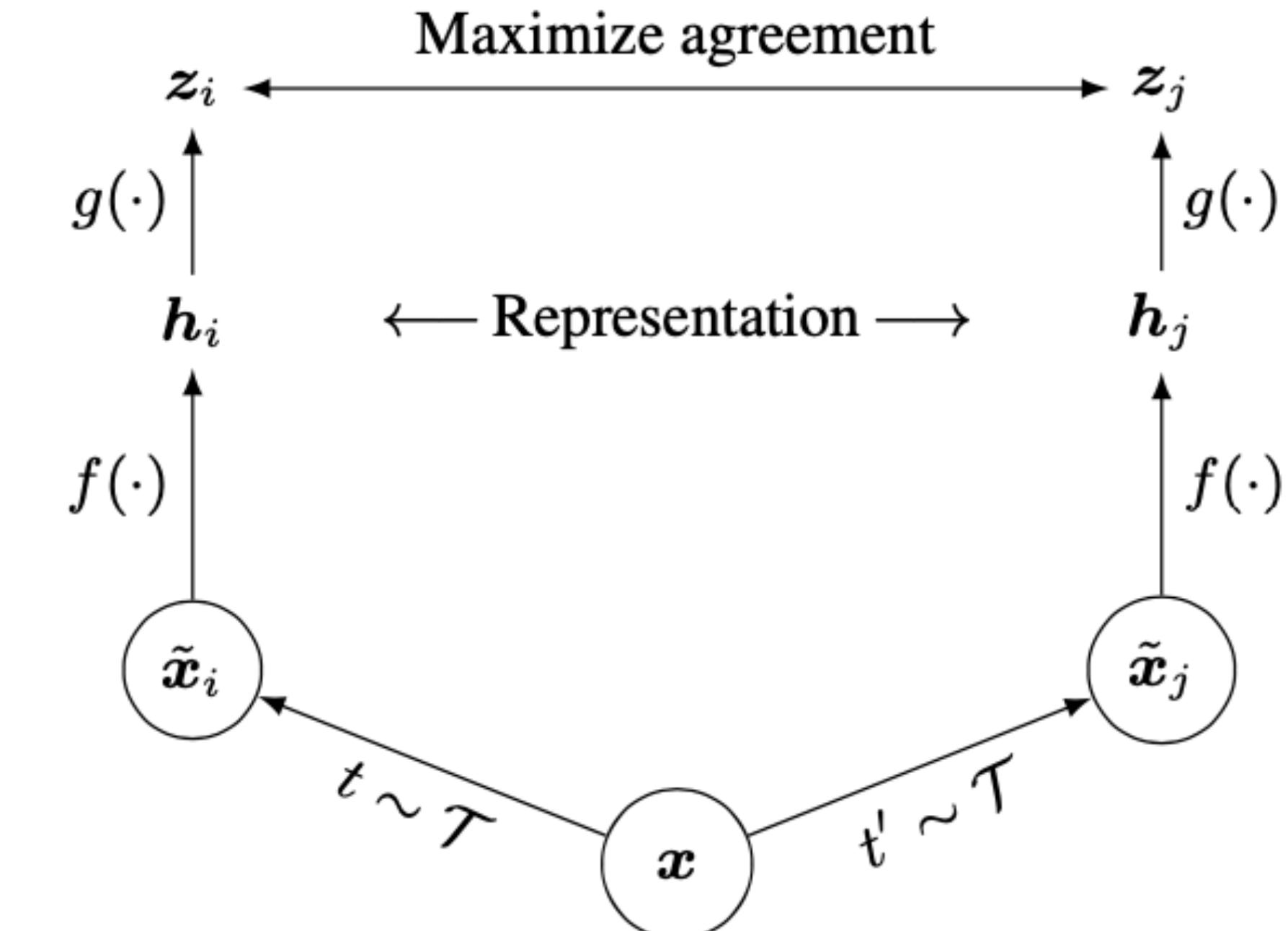


(b) Adjacent views.



(a) Without color distortion.

(b) With color distortion.



- Big Batches, Big Computers

¹With 128 TPU v3 cores, it takes ~ 1.5 hours to train our ResNet-50 with a batch size of 4096 for 100 epochs.

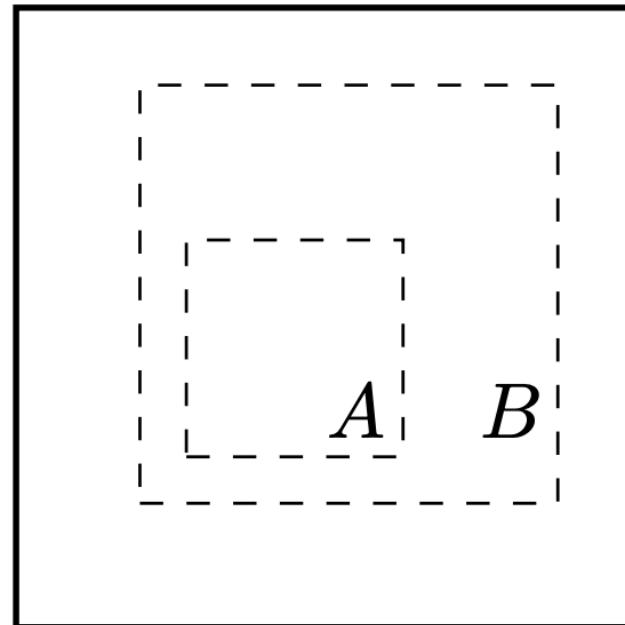
[Chen et. al.]

A Unifying Framework for Contrastive Learning

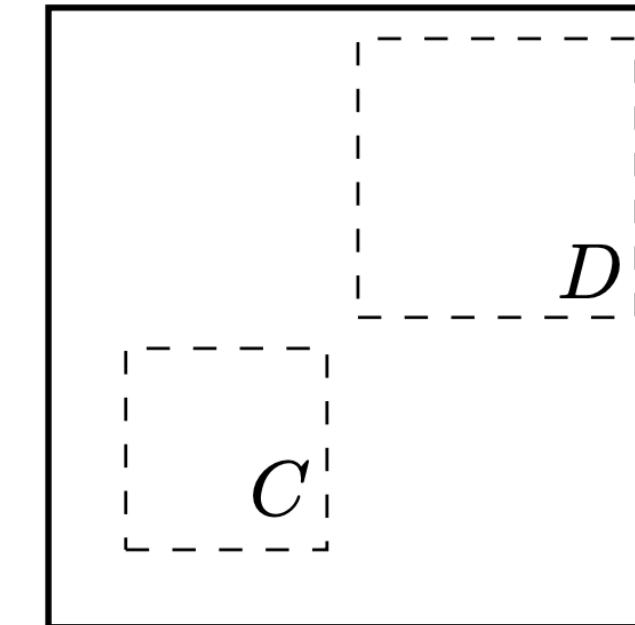
- Use \sim InfoNCE to **max MI between augmented inputs**

$$\ell_{i,j} = -\log \frac{\exp(\text{sim}(z_i, z_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(z_i, z_k)/\tau)}$$

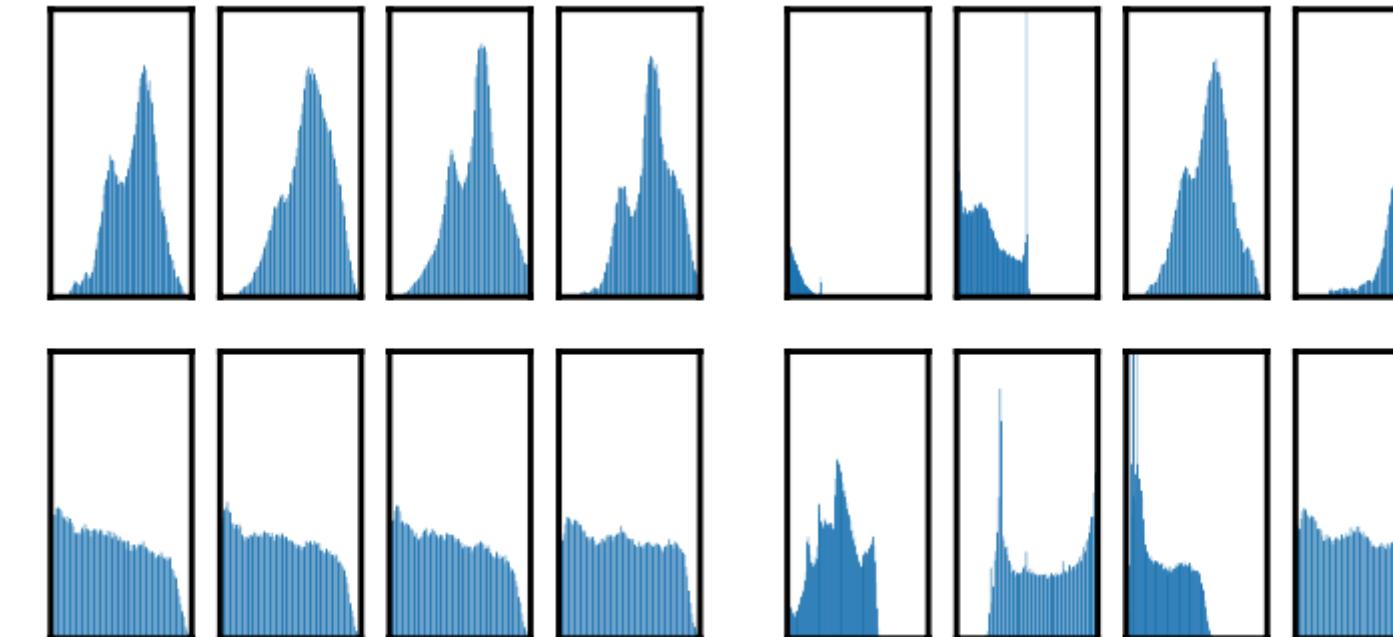
- Random Crop + Color Distortion Augmentation



(a) Global and local views.

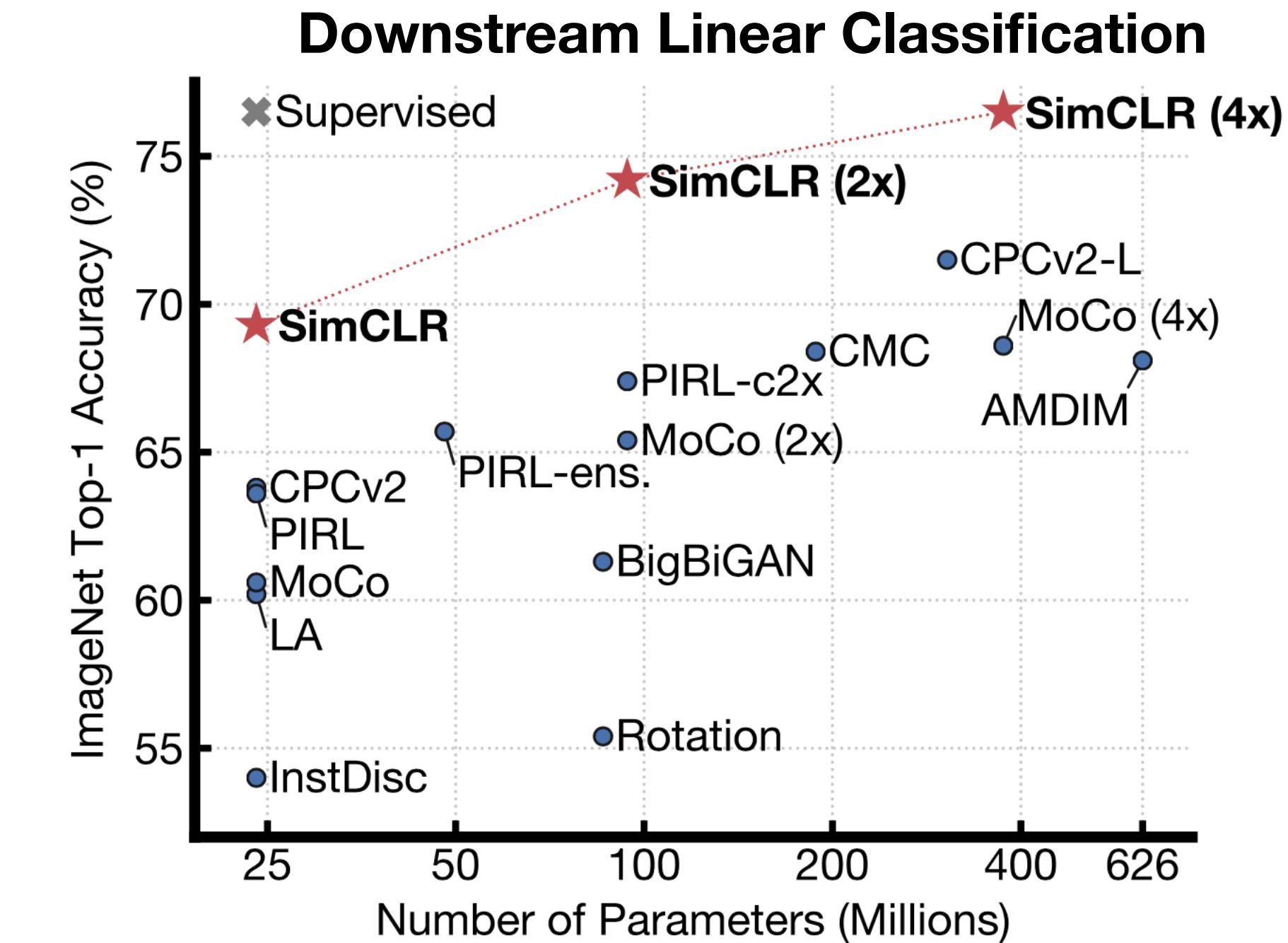


(b) Adjacent views.



(a) Without color distortion.

(b) With color distortion.



- Big Batches, Big Computers

¹With 128 TPU v3 cores, it takes \sim 1.5 hours to train our ResNet-50 with a batch size of 4096 for 100 epochs.

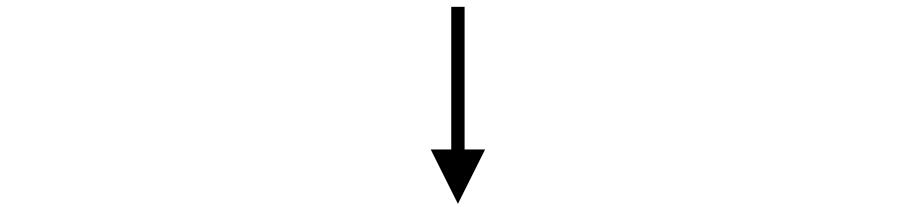
Mutual Information Neural Estimation (Deep InfoMax)

Mutual Information Neural Estimation (Deep InfoMax)

$$I(\mathbf{x}, \mathbf{z}) = E_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{p(\mathbf{x} | \mathbf{z})q(\mathbf{x} | \mathbf{z})}{p(\mathbf{x})q(\mathbf{x} | \mathbf{z})} \right] = \underbrace{E_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{q(\mathbf{x} | \mathbf{z})}{p(\mathbf{x})} \right]}_{+ E_{p(\mathbf{z})}[KL(p(\mathbf{x} | \mathbf{z}) || q(\mathbf{x} | \mathbf{z})]$$

Mutual Information Neural Estimation (Deep InfoMax)

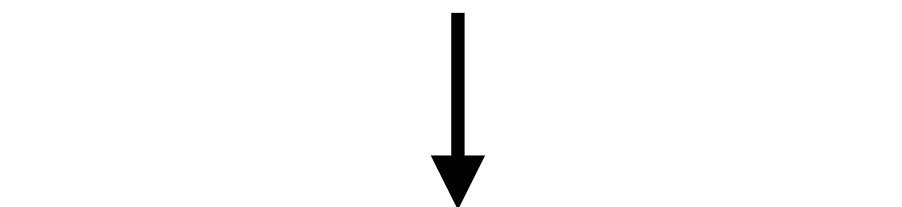
$$I(\mathbf{x}, \mathbf{z}) = E_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{p(\mathbf{x} | \mathbf{z})q(\mathbf{x} | \mathbf{z})}{p(\mathbf{x})q(\mathbf{x} | \mathbf{z})} \right] = \underbrace{E_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{q(\mathbf{x} | \mathbf{z})}{p(\mathbf{x})} \right]} + E_{p(\mathbf{z})}[KL(p(\mathbf{x} | \mathbf{z}) || q(\mathbf{x} | \mathbf{z})]$$



$$I(\mathbf{x}, \mathbf{z}) \geq H(x) + E_{p(\mathbf{x}, \mathbf{z})}[\log q(\mathbf{x} | \mathbf{z})]$$

Mutual Information Neural Estimation (Deep InfoMax)

$$I(\mathbf{x}, \mathbf{z}) = E_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{p(\mathbf{x} | \mathbf{z})q(\mathbf{x} | \mathbf{z})}{p(\mathbf{x})q(\mathbf{x} | \mathbf{z})} \right] = \underbrace{E_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{q(\mathbf{x} | \mathbf{z})}{p(\mathbf{x})} \right]} + E_{p(\mathbf{z})}[KL(p(\mathbf{x} | \mathbf{z}) || q(\mathbf{x} | \mathbf{z})]$$

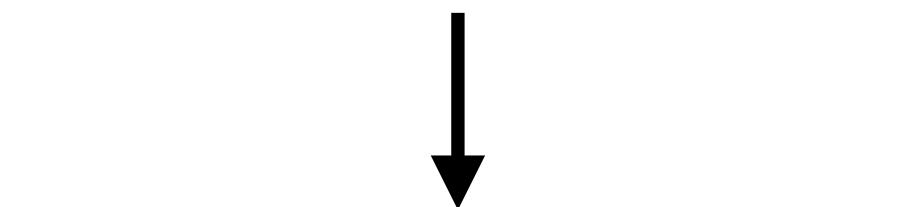


$$I(\mathbf{x}, \mathbf{z}) \geq H(x) + E_{p(\mathbf{x}, \mathbf{z})}[\log q(\mathbf{x} | \mathbf{z})]$$

We choose $q(\mathbf{x} | \mathbf{z}) = \frac{p(\mathbf{x})}{E_{p(\mathbf{x})}[e^{T(\mathbf{x}, \mathbf{z})}]} e^{T(\mathbf{x}, \mathbf{z})}$

Mutual Information Neural Estimation (Deep InfoMax)

$$I(\mathbf{x}, \mathbf{z}) = E_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{p(\mathbf{x} | \mathbf{z})q(\mathbf{x} | \mathbf{z})}{p(\mathbf{x})q(\mathbf{x} | \mathbf{z})} \right] = \underbrace{E_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{q(\mathbf{x} | \mathbf{z})}{p(\mathbf{x})} \right]} + E_{p(\mathbf{z})}[KL(p(\mathbf{x} | \mathbf{z}) || q(\mathbf{x} | \mathbf{z})]$$



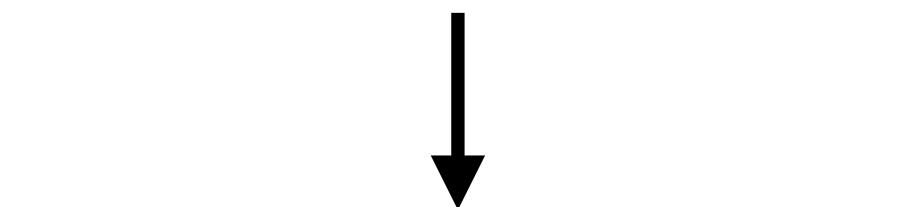
$$I(\mathbf{x}, \mathbf{z}) \geq H(x) + E_{p(\mathbf{x}, \mathbf{z})}[\log q(\mathbf{x} | \mathbf{z})]$$

We choose $q(\mathbf{x} | \mathbf{z}) = \frac{p(\mathbf{x})}{E_{p(\mathbf{x})}[e^{T(\mathbf{x}, \mathbf{z})}]} e^{T(\mathbf{x}, \mathbf{z})}$

$$I(\mathbf{x}, \mathbf{z}) \geq \max_{\theta} E_{p(\mathbf{x}, \mathbf{z})}[T_{\theta}] - \log(E_{p(\mathbf{x})p(\mathbf{z})}[e^{T_{\theta}}])$$

Mutual Information Neural Estimation (Deep InfoMax)

$$I(\mathbf{x}, \mathbf{z}) = E_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{p(\mathbf{x} | \mathbf{z})q(\mathbf{x} | \mathbf{z})}{p(\mathbf{x})q(\mathbf{x} | \mathbf{z})} \right] = \underbrace{E_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{q(\mathbf{x} | \mathbf{z})}{p(\mathbf{x})} \right]} + E_{p(\mathbf{z})}[KL(p(\mathbf{x} | \mathbf{z}) || q(\mathbf{x} | \mathbf{z})]$$



$$I(\mathbf{x}, \mathbf{z}) \geq H(x) + E_{p(\mathbf{x}, \mathbf{z})}[\log q(\mathbf{x} | \mathbf{z})]$$

We choose $q(\mathbf{x} | \mathbf{z}) = \frac{p(\mathbf{x})}{E_{p(\mathbf{x})}[e^{T(\mathbf{x}, \mathbf{z})}]} e^{T(\mathbf{x}, \mathbf{z})}$

$$I(\mathbf{x}, \mathbf{z}) \geq \max_{\theta} E_{p(\mathbf{x}, \mathbf{z})}[T_{\theta}] - \log(E_{p(\mathbf{x})p(\mathbf{z})}[e^{T_{\theta}}])$$

Bound tight if $T = \log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})p(\mathbf{z})}$

InfoNCE (Contrastive Predictive Coding)

[van den Oord et. al.]

InfoNCE (Contrastive Predictive Coding)

Similar to MINE, but:

We draw an additional K-1 samples from x to a normalisation constant:

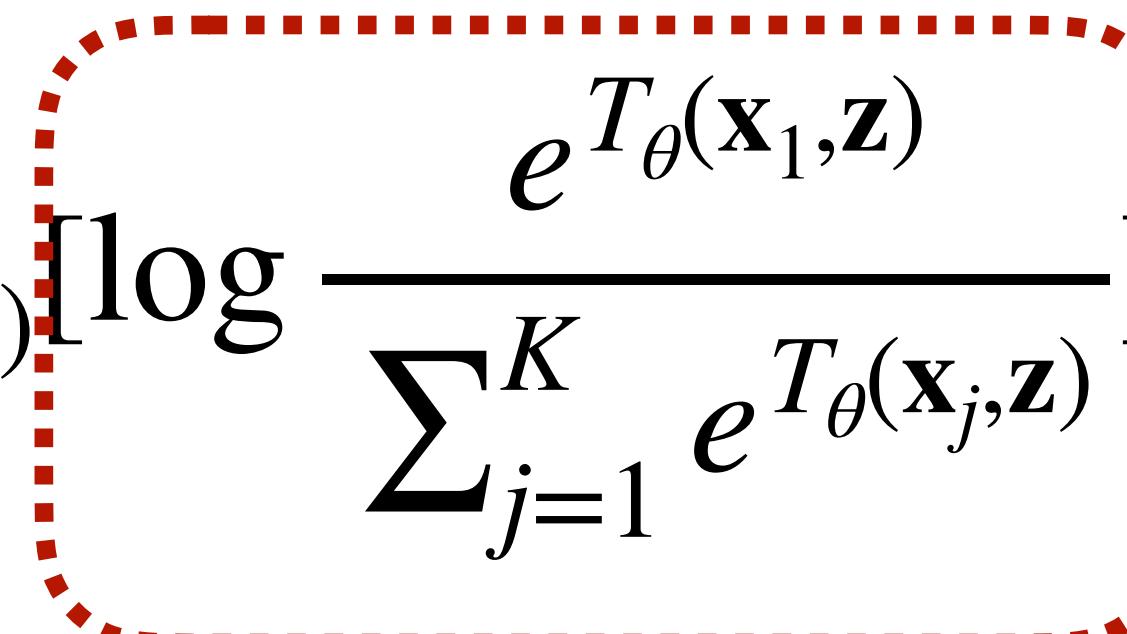
$$E_{p(\mathbf{x})}[e^T] \approx m(\mathbf{z} \mid \mathbf{x}_1 \dots \mathbf{x}_k) = \frac{1}{K} \sum_{i=1}^K e^{T(\mathbf{x}_i, \mathbf{z})}$$

InfoNCE (Contrastive Predictive Coding)

Similar to MINE, but:

We draw an additional K-1 samples from x to a normalisation constant:

$$E_{p(\mathbf{x})}[e^T] \approx m(\mathbf{z} | \mathbf{x}_1 \dots \mathbf{x}_k) = \frac{1}{K} \sum_{i=1}^K e^{T(\mathbf{x}_i, \mathbf{z})}$$

$$I(\mathbf{x}, \mathbf{z}) \geq \max_{\theta} E_{p(\mathbf{z}|\mathbf{x}_1)p(\mathbf{x}_1 \dots \mathbf{x}_K)} \left[\log \frac{e^{T_{\theta}(\mathbf{x}_1, \mathbf{z})}}{\sum_{j=1}^K e^{T_{\theta}(\mathbf{x}_j, \mathbf{z})}} \right] := I^{NCE}$$


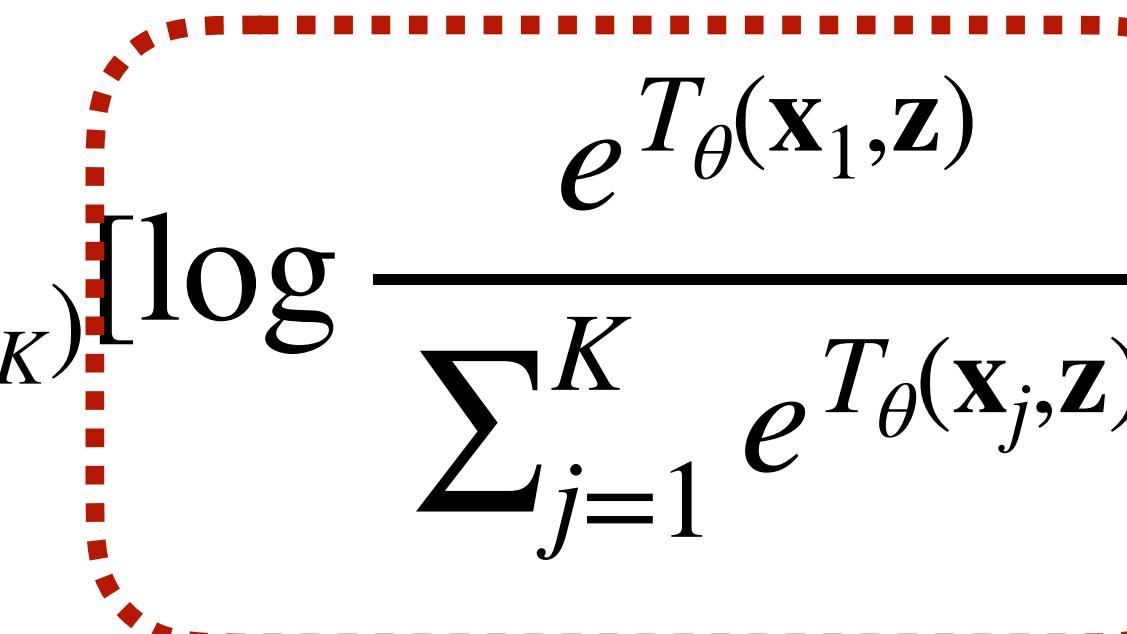
A red arrow points from the text "softmax cross-entropy" to the term $\frac{e^{T_{\theta}(\mathbf{x}_1, \mathbf{z})}}{\sum_{j=1}^K e^{T_{\theta}(\mathbf{x}_j, \mathbf{z})}}$ in the equation.

InfoNCE (Contrastive Predictive Coding)

Similar to MINE, but:

We draw an additional K-1 samples from x to a normalisation constant:

$$E_{p(\mathbf{x})}[e^T] \approx m(\mathbf{z} | \mathbf{x}_1 \dots \mathbf{x}_k) = \frac{1}{K} \sum_{i=1}^K e^{T(\mathbf{x}_i, \mathbf{z})}$$

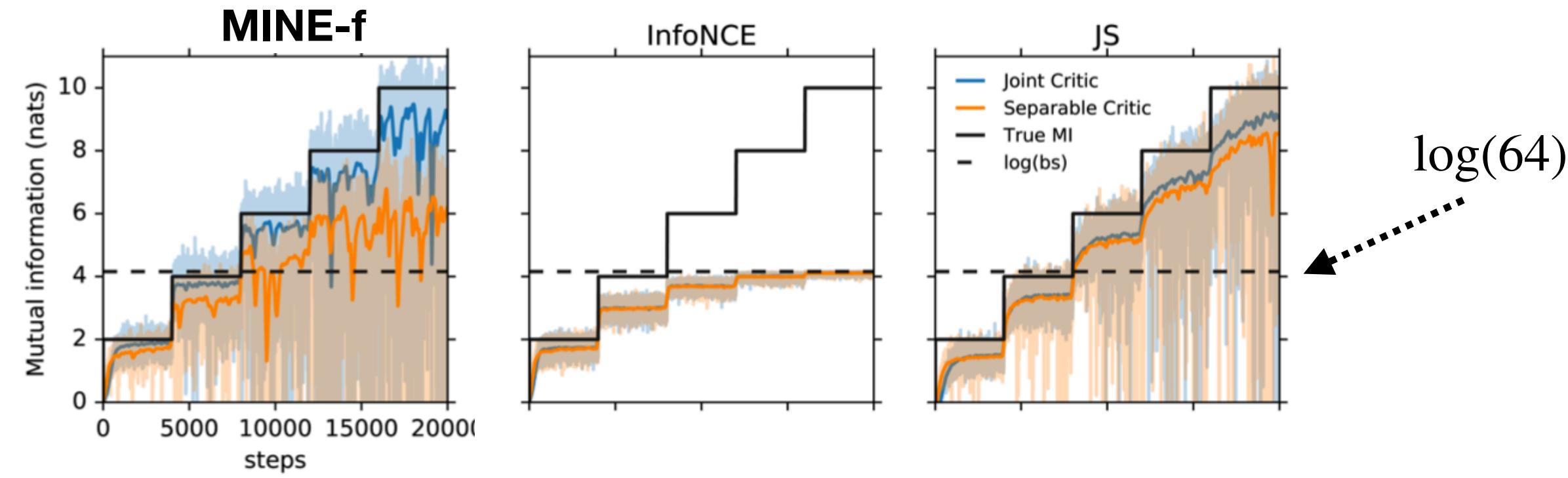
$$I(\mathbf{x}, \mathbf{z}) \geq \max_{\theta} E_{p(\mathbf{z}|\mathbf{x}_1)p(\mathbf{x}_1 \dots \mathbf{x}_K)} \left[\log \frac{e^{T_{\theta}(\mathbf{x}_1, \mathbf{z})}}{\sum_{j=1}^K e^{T_{\theta}(\mathbf{x}_j, \mathbf{z})}} \right] := I^{NCE}$$


Note that: $I(\mathbf{x}, \mathbf{z}) \geq I^{NCE} + \log K$

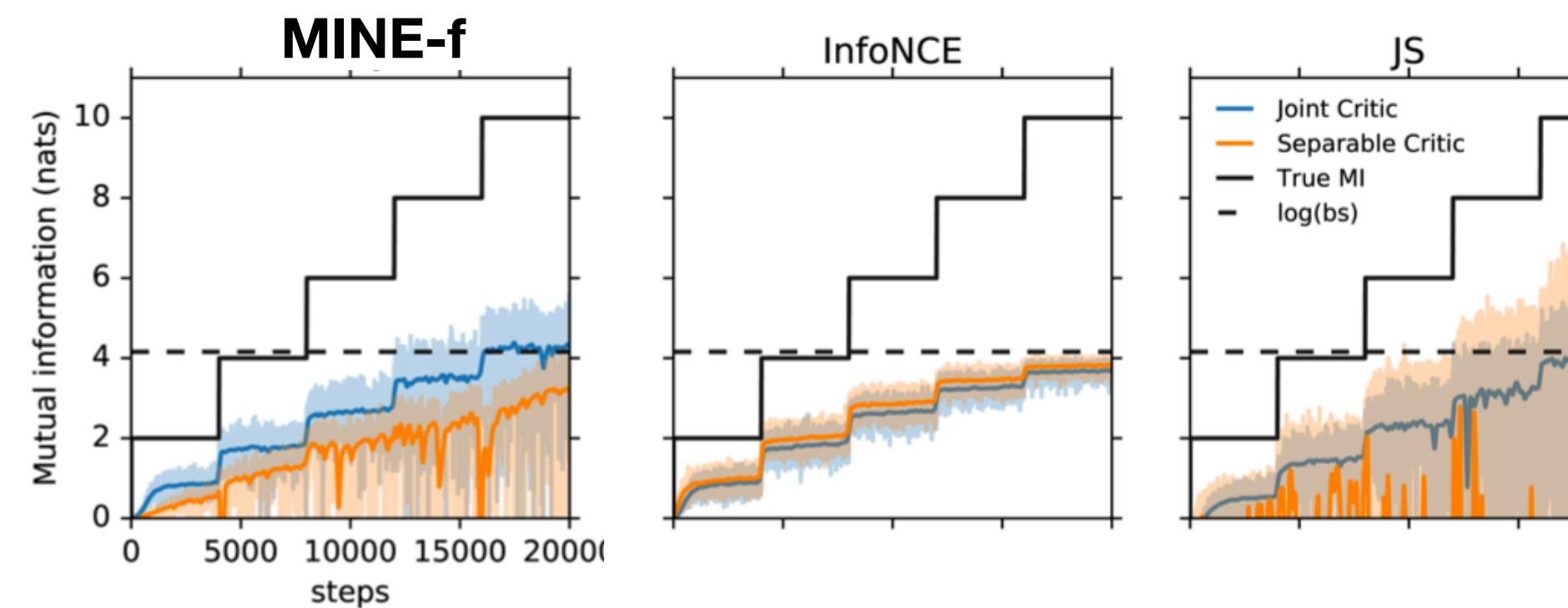
Larger K (batch sizes) will yield less biased estimates

Comparing MI estimators

Correlated Gaussian



$\mathbf{z} = (W\mathbf{x})^3; \quad \mathbf{x} \sim \mathcal{N}(0,1)$



Which One Should I Use?

Batch-Size: N

2N Critic Evaluations

- JSD
- MINE
- MINE-f

N^2 Critic Evaluations

- InfoNCE

Not an estimator

- JSD

Bias

-InfoNCE

-MINE

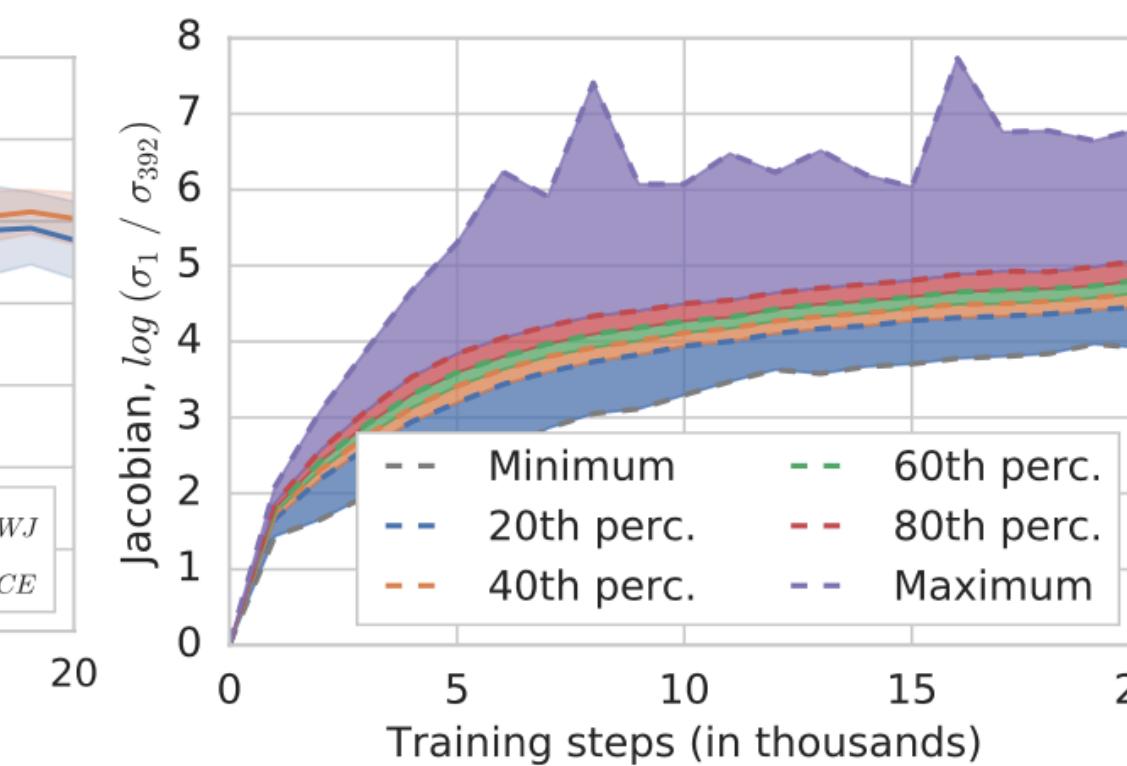
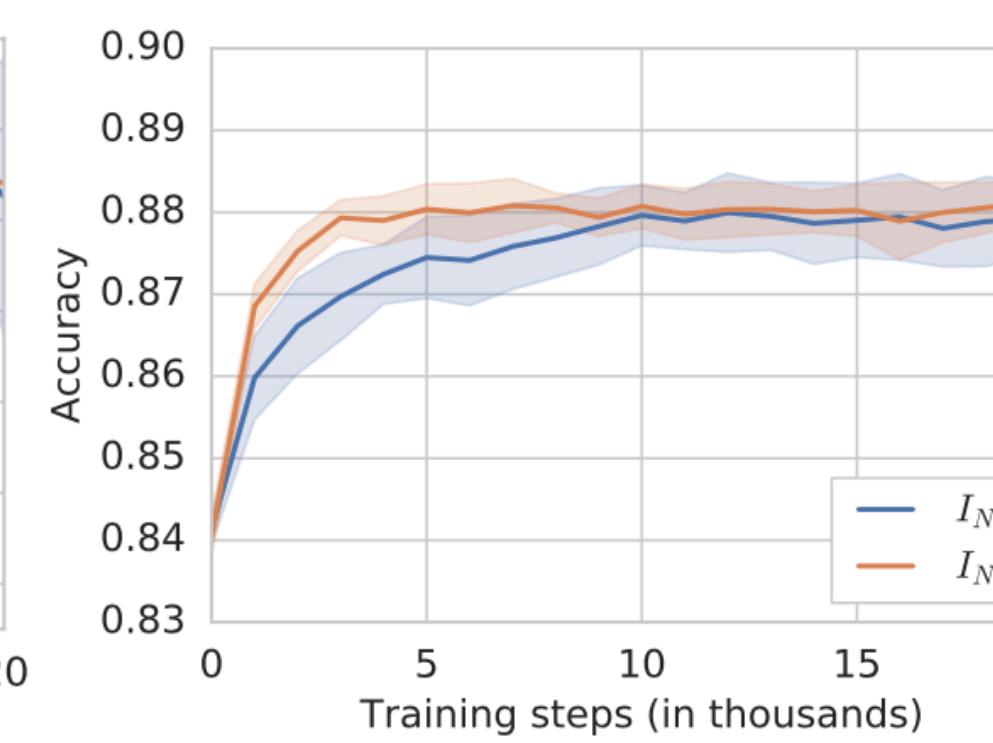
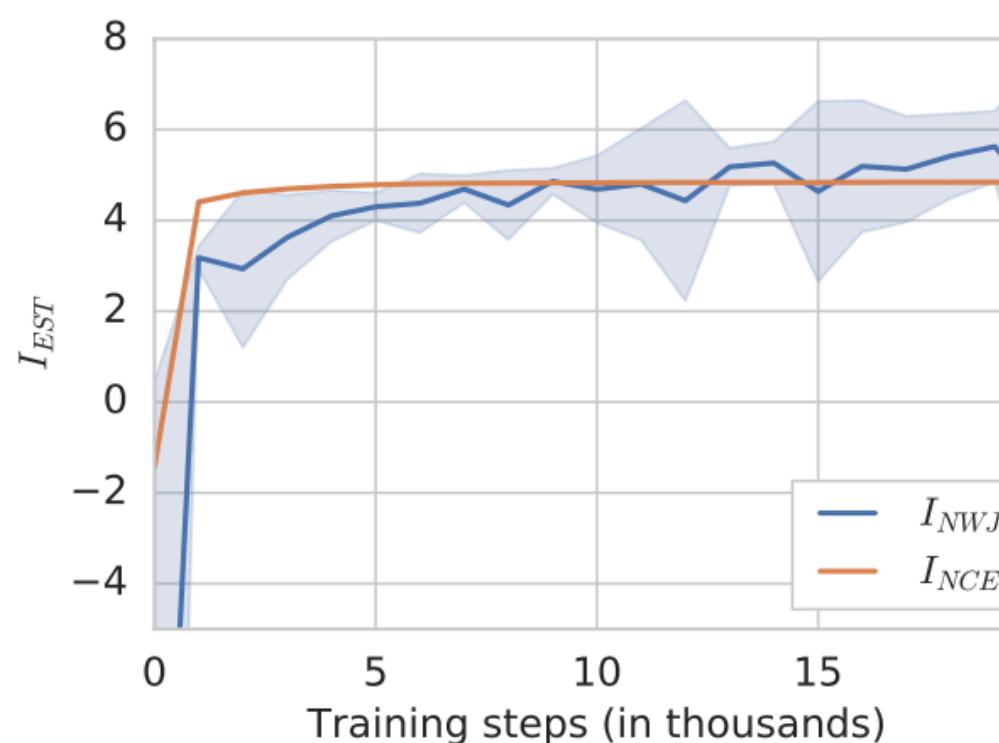
-MINE-f

Variance

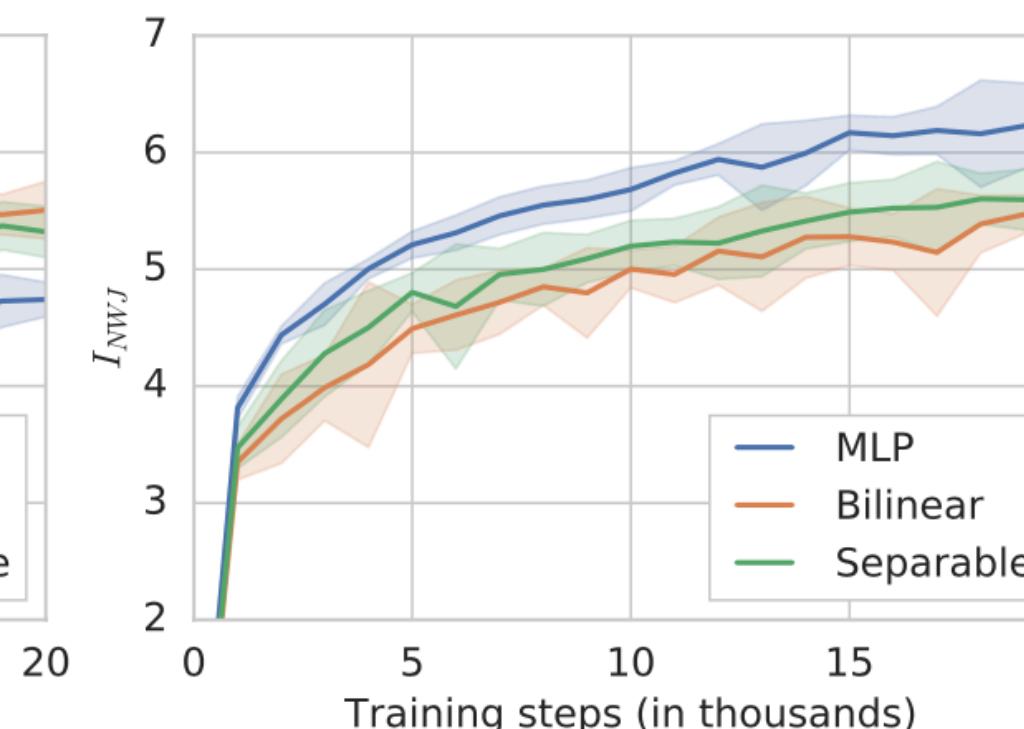
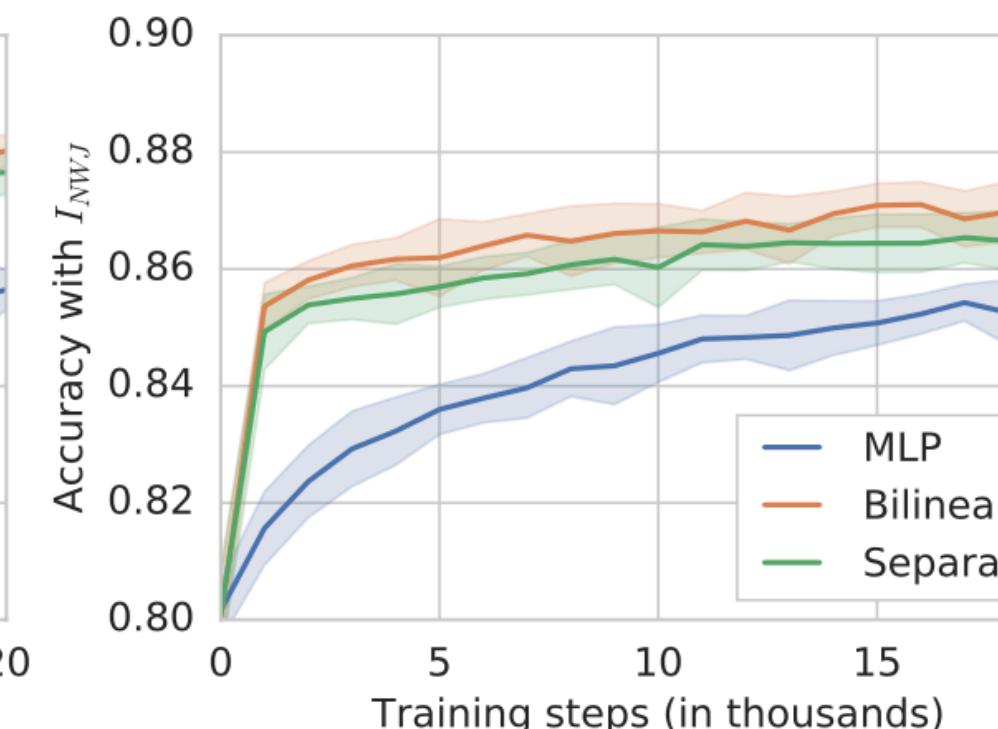
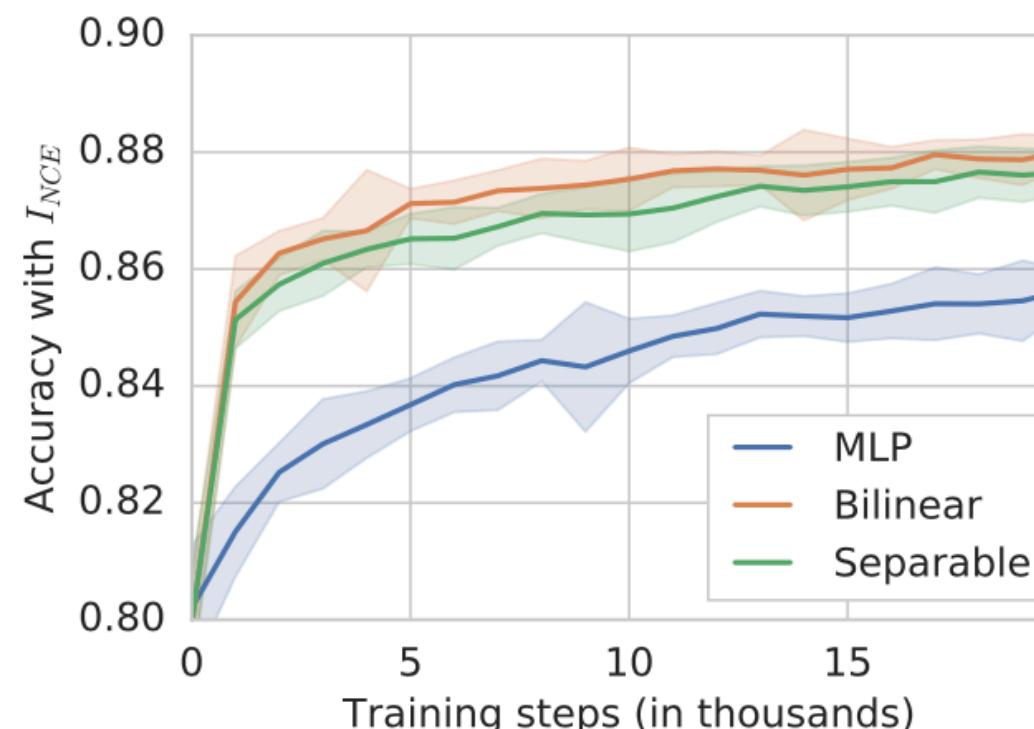
* If possible use InfoNCE with large batch size

Is MI Maximisation Really the Solution?

- During training encoders become less invertible



- More flexible critics can make for worse representations



$$\text{Bilinear: } T(\mathbf{z}, \mathbf{x}) = \mathbf{z}^\top \mathbf{W} \mathbf{x}$$

$$\text{Separable: } T(\mathbf{z}, \mathbf{x}) = \phi(\mathbf{z}^\top) \phi(\mathbf{x})$$

[Tschannen et. al.]

Is MI Maximisation Really the Solution?

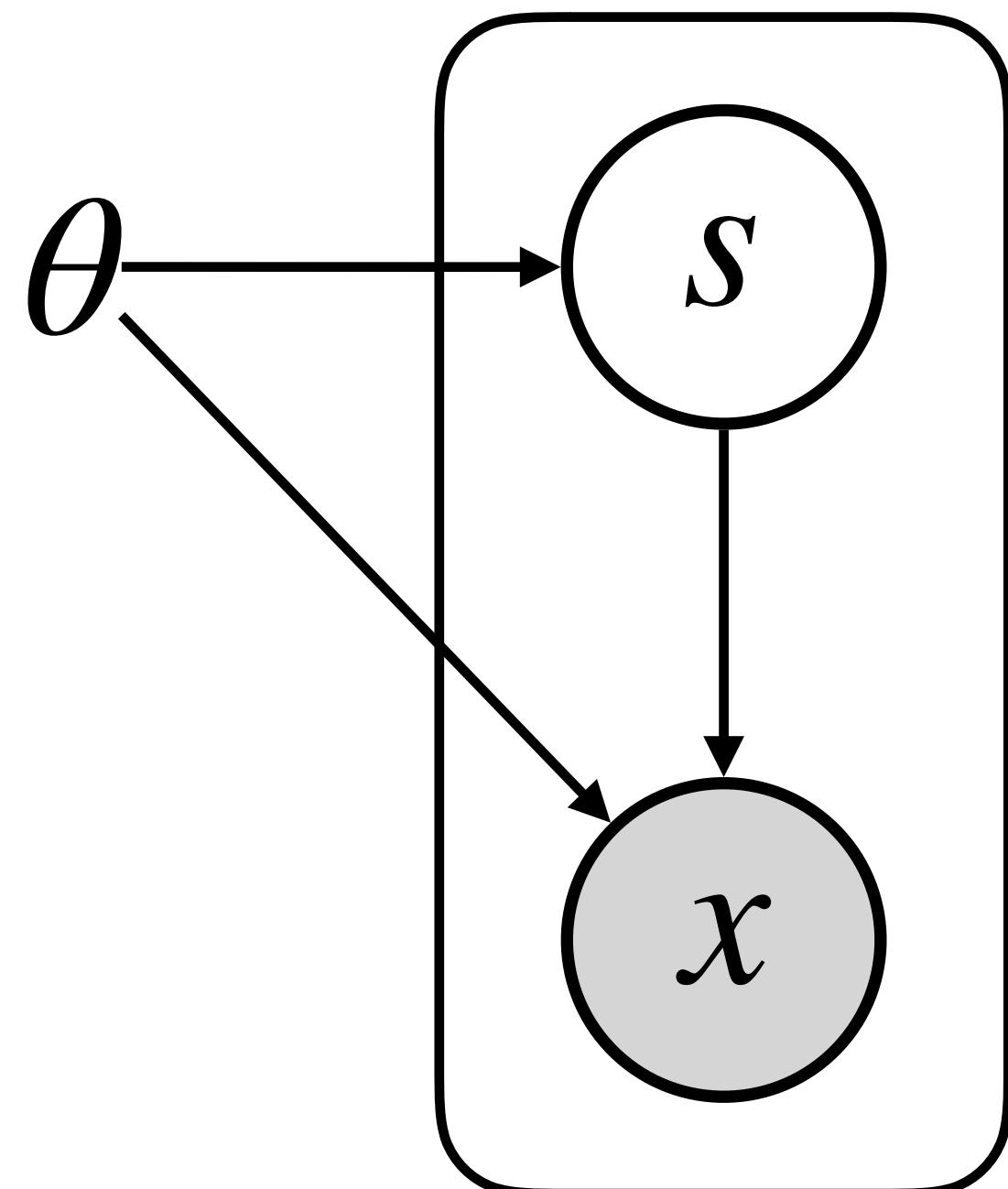
- A bijection can be applied to a useful representation making it not useful, maintaining $I(\mathbf{x}, \mathbf{z})$
- Utility of representation is a function of "**decodable information**", not MI
- This depends on **inductive biases** from:
 - MI estimator
 - Critic Function + Encoder Function
 - Objects between which MI is maximised

References

- [Larsson et. al.] Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles <https://arxiv.org/pdf/1603.09246.pdf>
- [Gidaris et. al.] Unsupervised Representation Learning by Predicting Image Rotations <https://openreview.net/forum?id=S1v4N2I0->
- [Arandjelovic et. al.] Look, Listen and Learn <https://arxiv.org/abs/1705.08168>
- [Mikolov et. al.] Distributed Representations of Words and Phrases and their Compositionality <https://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases-and-their-compositionality.pdf>
- [Hjelm et. al.] Learning deep representations by mutual information estimation and maximization <https://arxiv.org/abs/1808.06670>
- [Belghazi et. al.] MINE: Mutual Information Neural Estimation <https://arxiv.org/abs/1801.04062>
- [van den Oord et. al.] Representation Learning with Contrastive Predictive Coding <https://arxiv.org/abs/1807.03748>
- [Henáff et. al.] Data-Efficient Image Recognition with Contrastive Predictive Coding <https://arxiv.org/abs/1905.09272>
- [Chen et. al.] A Simple Framework for Contrastive Learning of Visual Representations <https://arxiv.org/abs/2002.05709>
- [Poole et. al.] On Variational Bounds of Mutual Information <https://arxiv.org/abs/1905.06922>
- [Tschannen et. al.] On Mutual Information Maximization for Representation Learning <https://arxiv.org/abs/1907.13625>

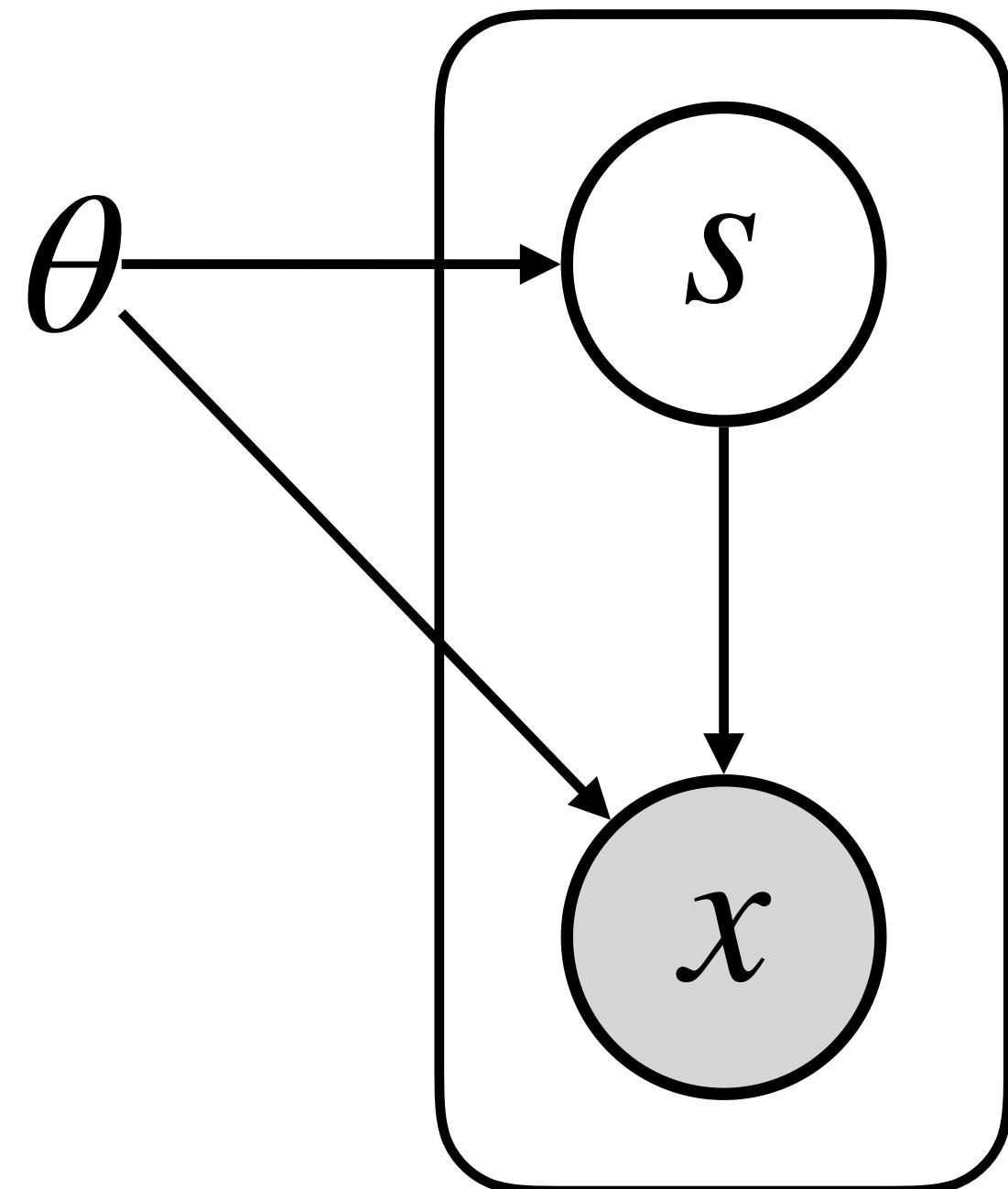
A Generative Modelling Perspective

Latent Variable Generative Models

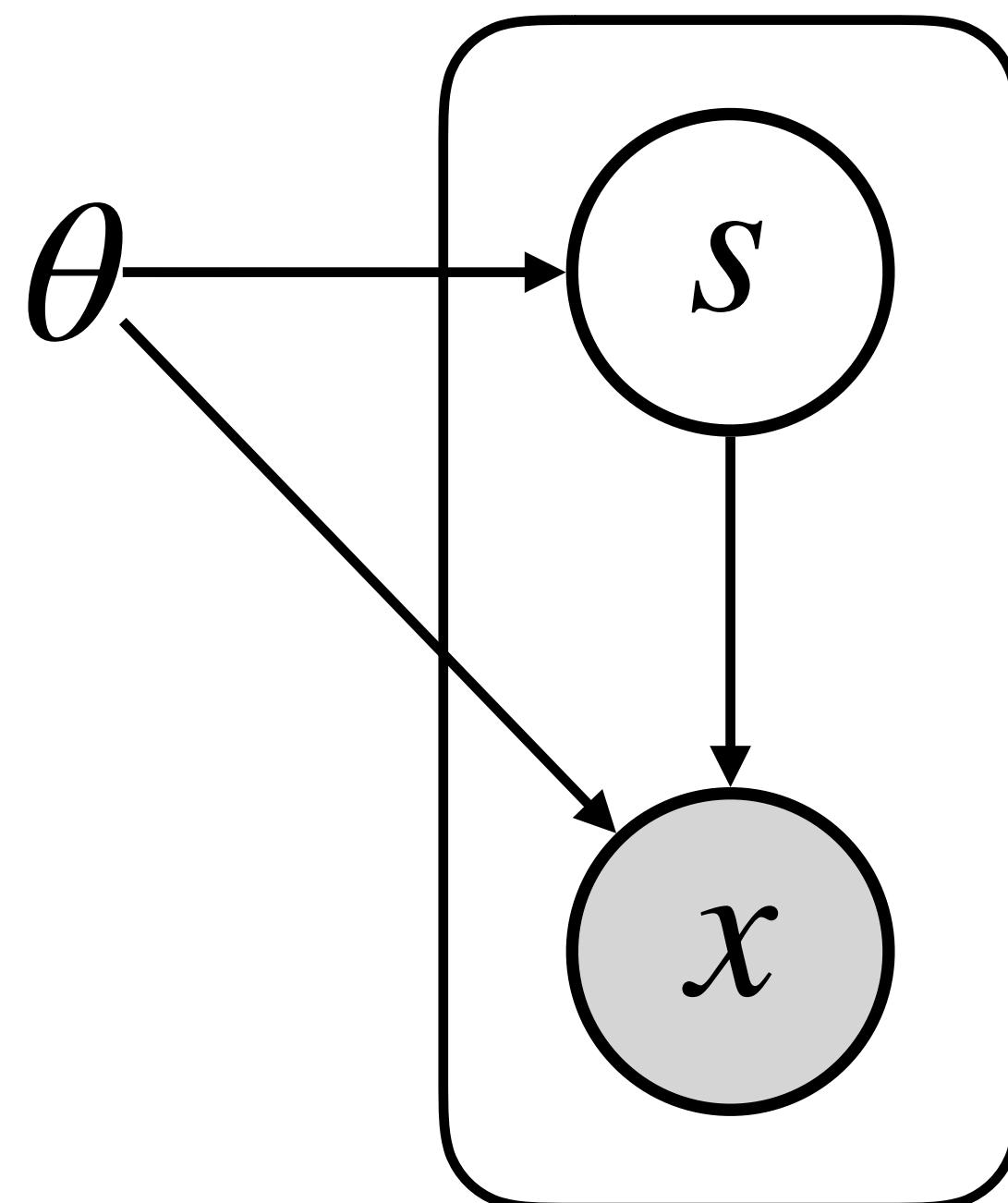


Latent Variable Generative Models

Goal: $\theta^* = \arg \max_{\theta \in \Theta} \log p_\theta(D)$



Latent Variable Generative Models



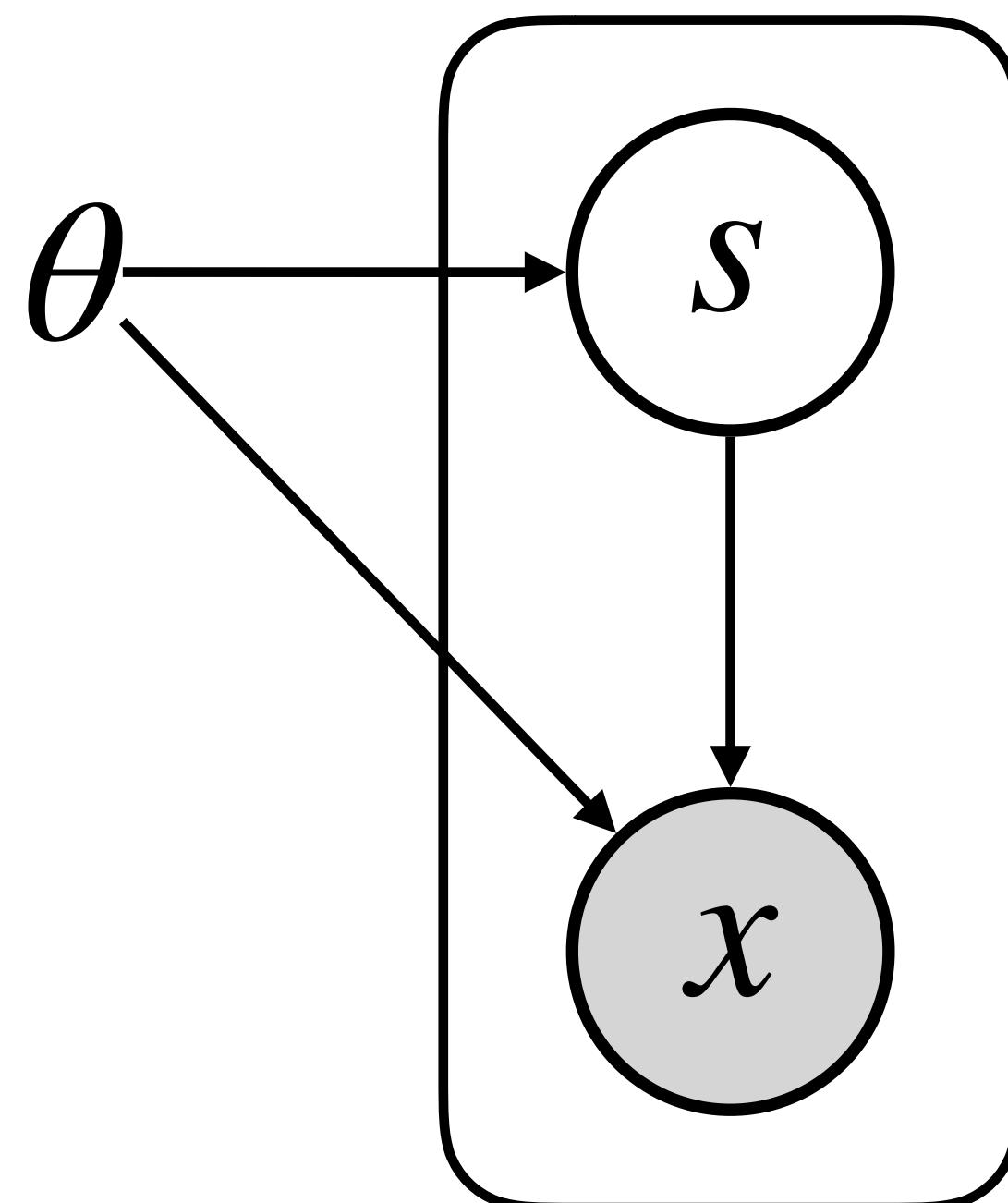
Goal: $\theta^* = \arg \max_{\theta \in \Theta} \log p_\theta(D)$

Probabilistic PCA

$$p_\theta(s) = \mathcal{N}(s; 0, I)$$

$$p_\theta(x | s) = \mathcal{N}(x; Ws + \mu, \sigma^2 I)$$

Latent Variable Generative Models



Goal: $\theta^* = \arg \max_{\theta \in \Theta} \log p_\theta(D)$

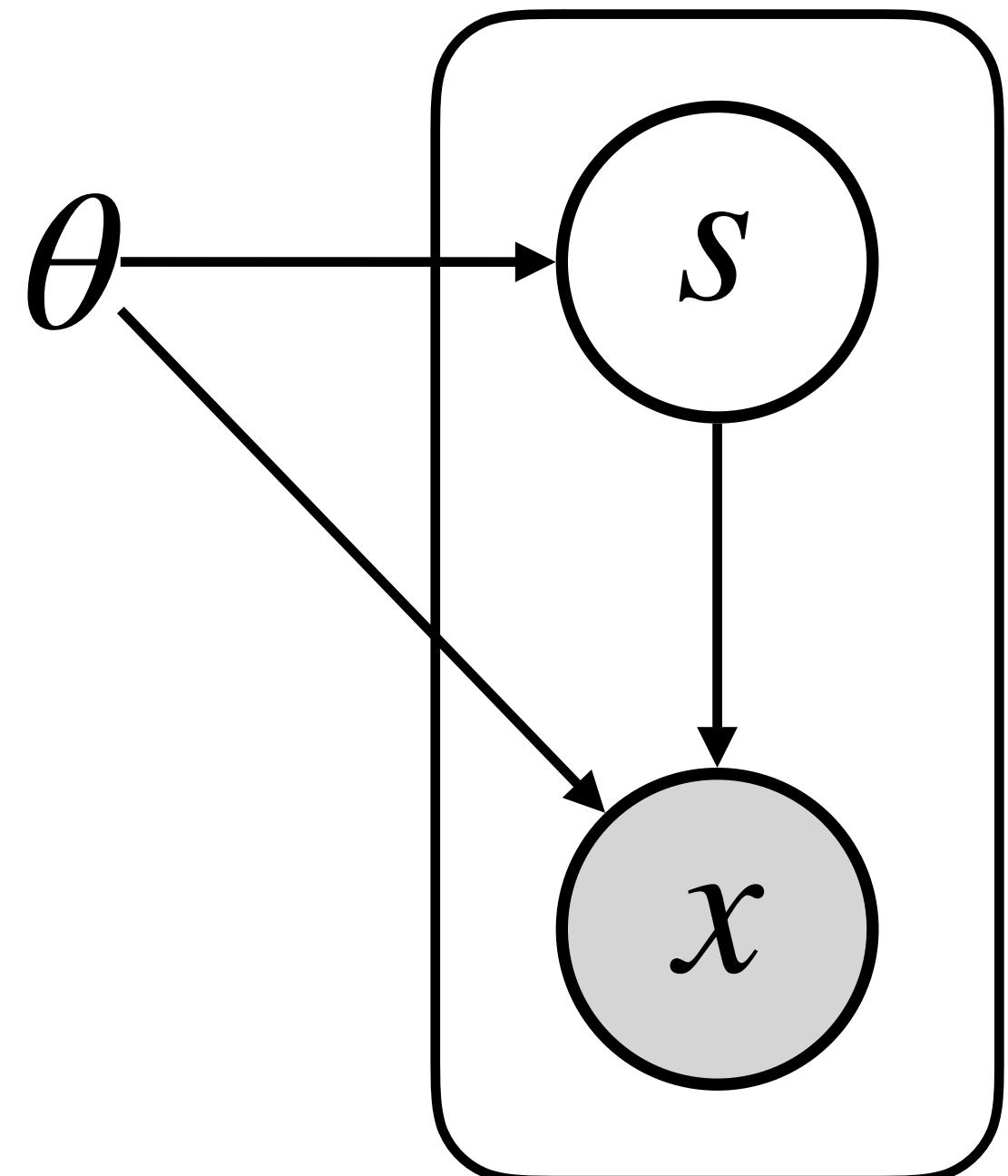
Factor Analysis

$$p_\theta(s) = \mathcal{N}(s; 0, I)$$

$$p_\theta(x | s) = \mathcal{N}(x; Ws + \mu, \Sigma)$$

Latent Variable Generative Models

Goal: $\theta^* = \arg \max_{\theta \in \Theta} \log p_\theta(D)$



VAE

$$p_\theta(s) = \mathcal{N}(s; 0, I)$$

$$p_\theta(x | s) = \mathcal{N}\left(x; f_{nn}^\mu(s), f_{nn}^{\sigma^2}(s)\right)$$

Identifiability

Identifiability

Identifiability: the “true” parameters will be recovered given infinite observations

Identifiability

Identifiability: the “true” parameters will be recovered given infinite observations

For generative models: $p_{\theta}(x) = p_{\theta'}(x) \implies \theta = \theta'$

Identifiability

Identifiability: the “true” parameters will be recovered given infinite observations

For generative models: $p_\theta(x) = p_{\theta'}(x) \implies \theta = \theta'$

Identifiable models provide:

- A principled approach to representation learning (as opposed to “just” generative modelling).
- Links to other unsupervised desiderata (e.g., “disentanglement”)

PCA Non-Identifiability

PCA Non-Identifiability

PCA / Factor Analysis:

PCA Non-Identifiability

PCA / Factor Analysis:

Scaling: $WS = \frac{\alpha}{\alpha} WS = (\alpha^{-1} W)(\alpha S)$

PCA Non-Identifiability

PCA / Factor Analysis:

Scaling: $WS = \frac{\alpha}{\alpha} WS = (\alpha^{-1} W)(\alpha S)$

Rotation: $WS = WRR^{-1}S = W'S'$

VAE Non-identifiability

VAE Non-identifiability

Theorem. Let $p(s) = \prod p_d(s_d)$, then there exists an infinite family of bijections of the form $f: \mathcal{S} \rightarrow \mathcal{S}$ such that

$$\int p_\theta(x | s)p(s)ds = \int p_{\theta'}(x | f(s))p(f(s))ds$$

for some alternative generative model with parameters θ' .

VAE Non-identifiability

Theorem. Let $p(s) = \prod p_d(s_d)$, then there exists an infinite family of bijections of the form $f: \mathcal{S} \rightarrow \mathcal{S}$ such that

$$\int p_\theta(x | s)p(s)ds = \int p_{\theta'}(x | f(s))p(f(s))ds$$

for some alternative generative model with parameters θ' .

$$p_\theta(x)$$

VAE Non-identifiability

Theorem. Let $p(s) = \prod p_d(s_d)$, then there exists an infinite family of bijections of the form $f: \mathcal{S} \rightarrow \mathcal{S}$ such that

$$\int p_\theta(x | s)p(s)ds = \int p_{\theta'}(x | f(s))p(f(s))ds$$

for some alternative generative model with parameters θ' .

$$p_\theta(x) \approx p_{\theta^*}(x)$$

VAE Non-identifiability

Theorem. Let $p(s) = \prod p_d(s_d)$, then there exists an infinite family of bijections of the form $f: \mathcal{S} \rightarrow \mathcal{S}$ such that

$$\int p_\theta(x | s)p(s)ds = \int p_{\theta'}(x | f(s))p(f(s))ds$$

for some alternative generative model with parameters θ' .

$$p_\theta(x) \approx p_{\theta^*}(x)$$

$$p_\theta(s, x)$$

VAE Non-identifiability

Theorem. Let $p(s) = \prod p_d(s_d)$, then there exists an infinite family of bijections of the form $f: \mathcal{S} \rightarrow \mathcal{S}$ such that

$$\int p_\theta(x | s)p(s)ds = \int p_{\theta'}(x | f(s))p(f(s))ds$$

for some alternative generative model with parameters θ' .

$$p_\theta(x) \approx p_{\theta^*}(x)$$

$$p_\theta(s, x) \not\approx p_{\theta^*}(s, x)$$

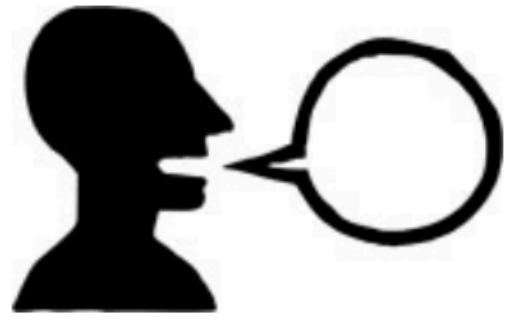
The Cocktail Party Problem

The Cocktail Party Problem

music



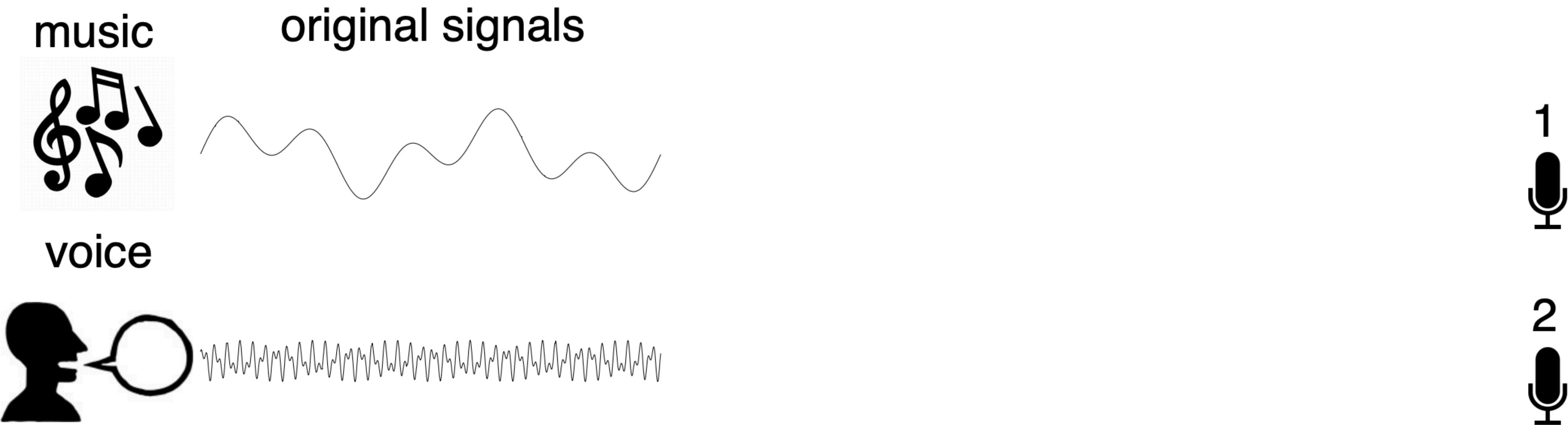
voice



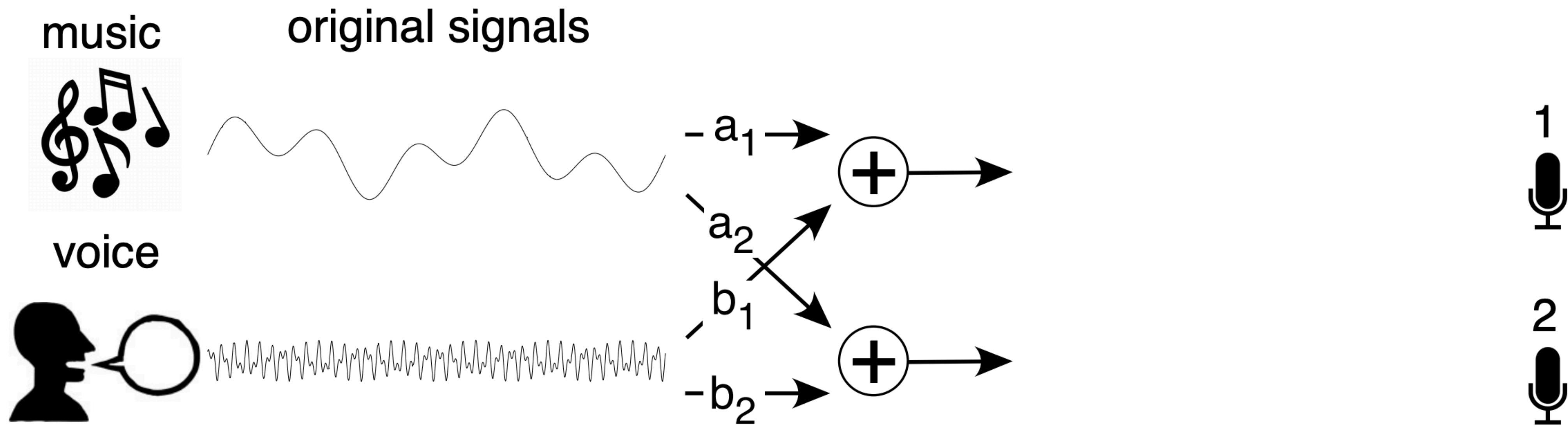
1

2

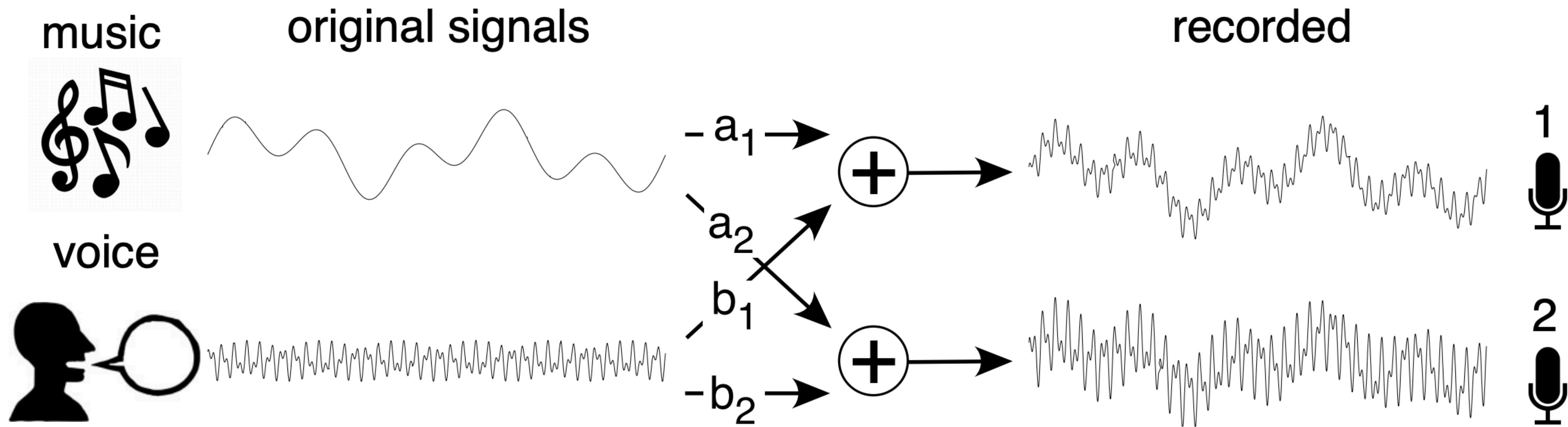
The Cocktail Party Problem



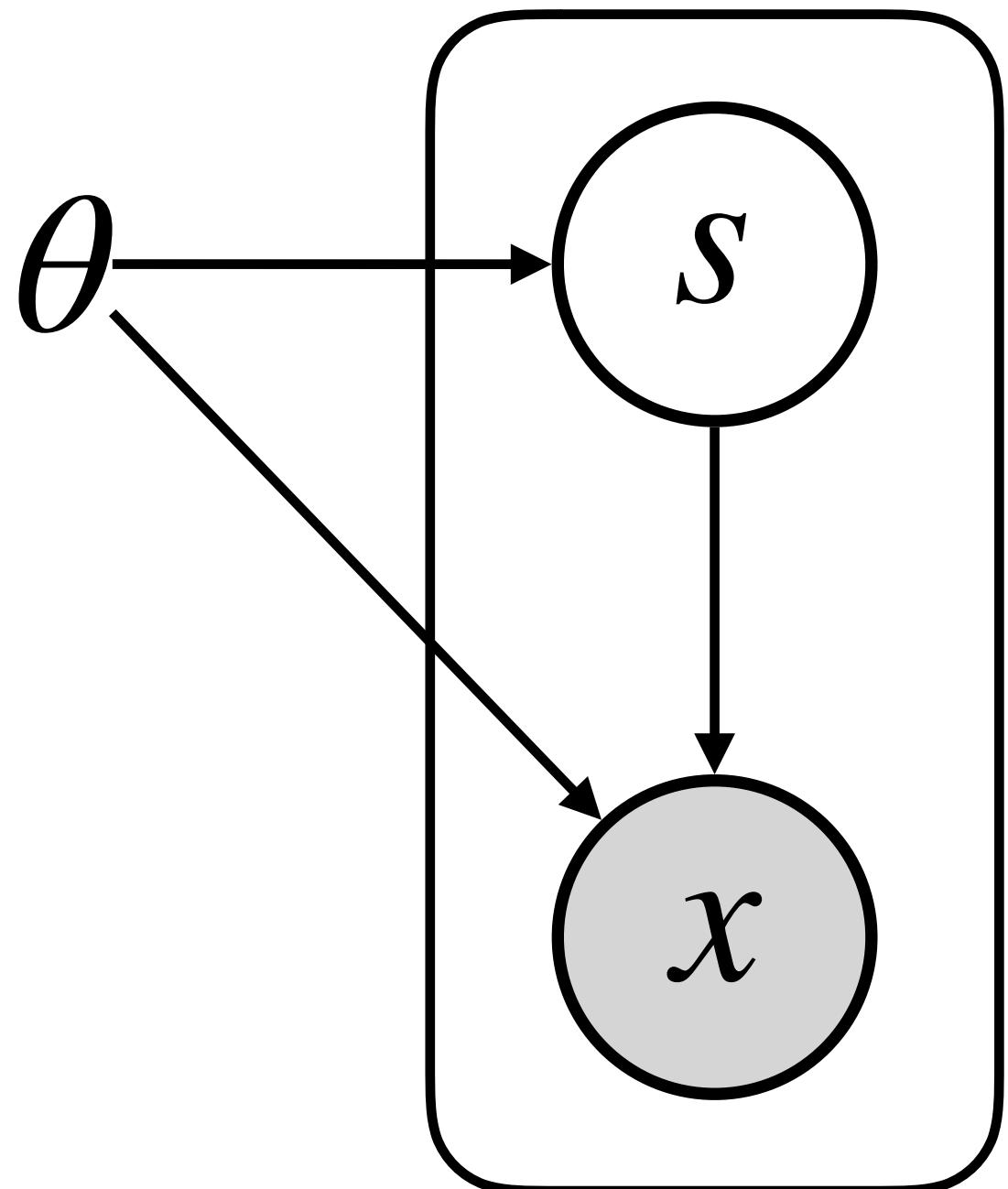
The Cocktail Party Problem



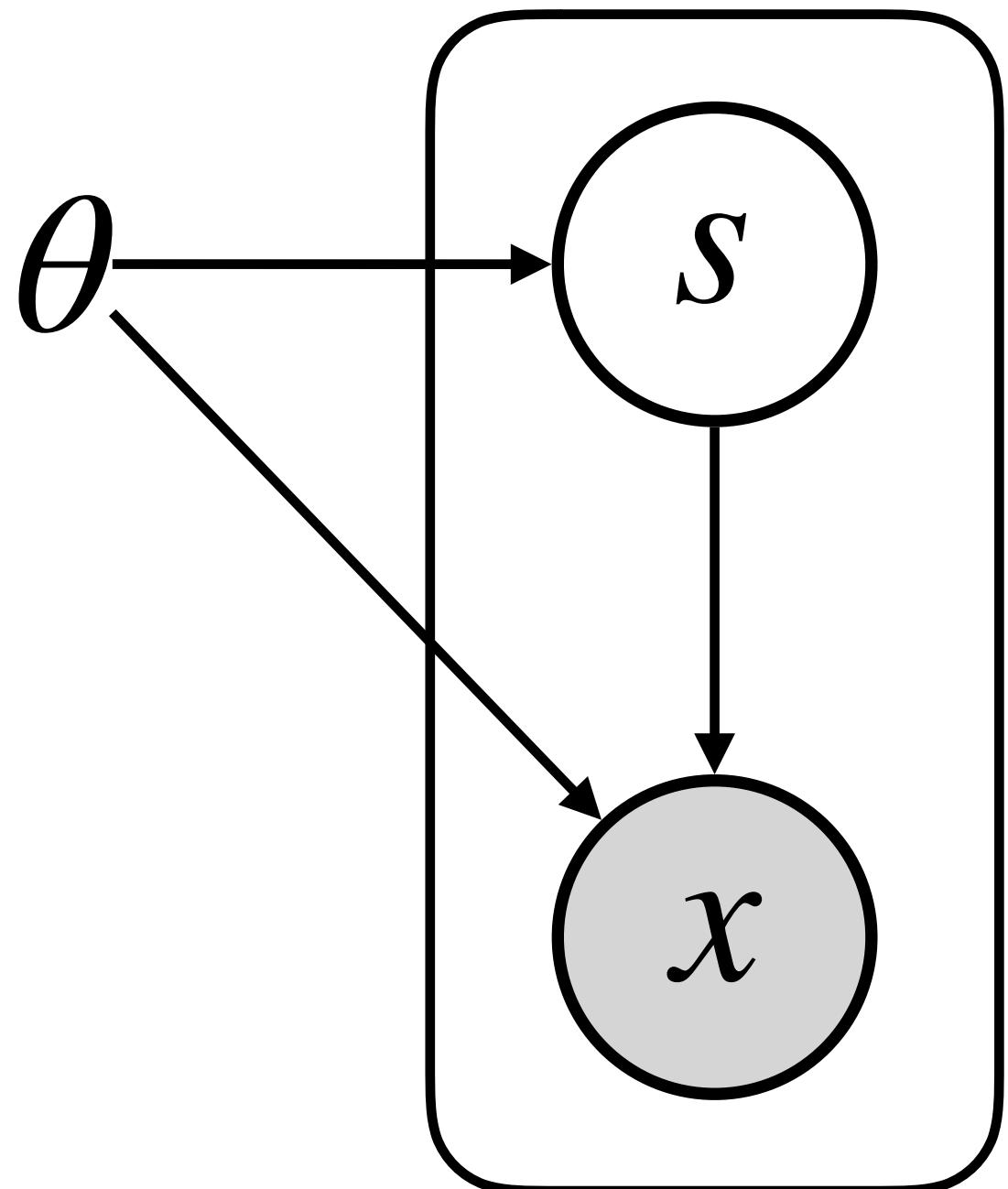
The Cocktail Party Problem



Linear ICA

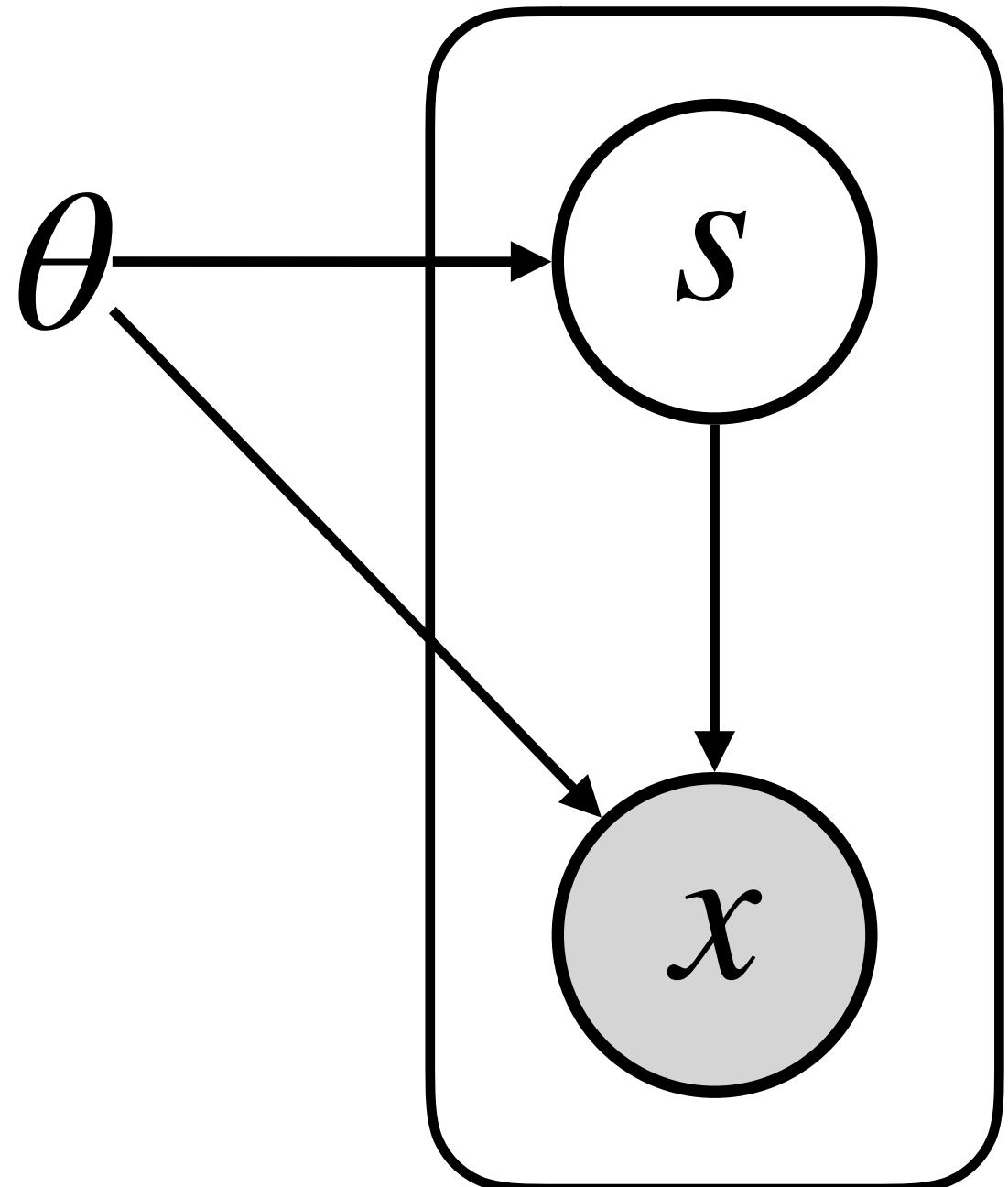


Linear ICA



$$p_{\theta}(s) = \prod_{i=1}^d p_i(s_i; \theta)$$

Linear ICA

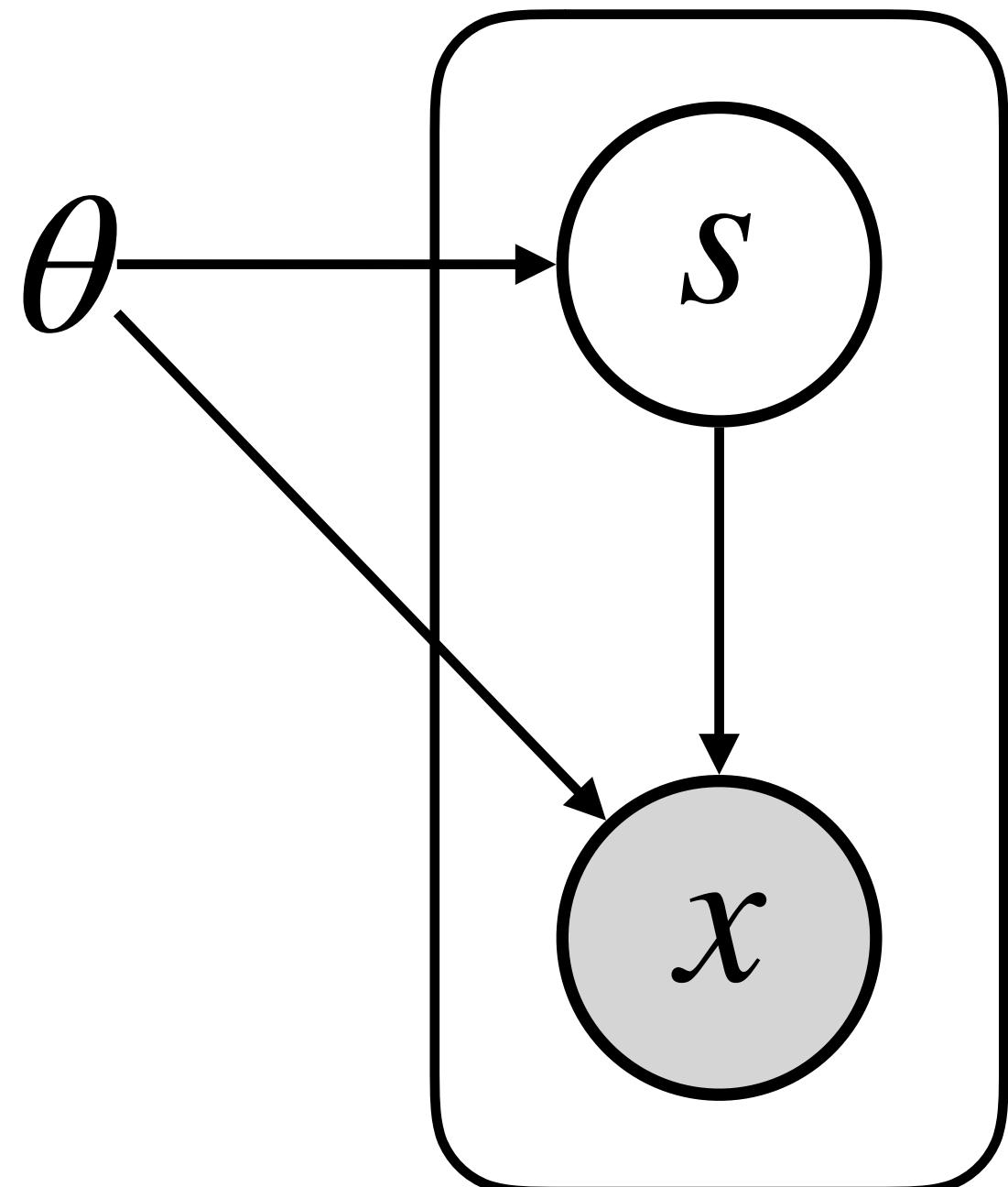


$$p_{\theta}(s) = \prod_{i=1}^d p_i(s_i; \theta)$$

$$x = As; \quad A \in \mathbb{R}^{d \times d}$$

Linear ICA

Goal: recover A^{-1} $\implies s = A^{-1}x$



$$p_{\theta}(s) = \prod_{i=1}^d p_i(s_i; \theta)$$

$$x = As; \quad A \in \mathbb{R}^{d \times d}$$

Maximum Likelihood Solution

D. Mackay. 2003.

Maximum Likelihood Solution

$$W = A^{-1};$$

Maximum Likelihood Solution

$$W = A^{-1};$$

$$a_i = \sum_j W_{ij} x_j$$

Maximum Likelihood Solution

$$W = A^{-1};$$

$$a_i = \sum_j W_{ij} x_j$$

$$\phi(a_i) = \frac{d}{da_i} \log p_i(a_i)$$

Maximum Likelihood Solution

$$W = A^{-1}; \quad a_i = \sum_j W_{ij} x_j \quad \phi(a_i) = \frac{d}{da_i} \log p_i(a_i)$$

$$\log p(x^{(n)} | A) = \log |\det W| + \sum_i \log p_i(a_i^{(n)})$$

Maximum Likelihood Solution

$$W = A^{-1};$$

$$a_i = \sum_j W_{ij} x_j$$

$$\phi(a_i) = \frac{d}{da_i} \log p_i(a_i)$$

$$\log p(x^{(n)} | A) = \log |\det W| + \sum_i \log p_i(a_i^{(n)})$$

$$\frac{\partial}{\partial W_{ij}} \log p(x^{(n)} | A) = [W^T]_{ij}^{-1} + x_j^{(n)} \phi(a_i^{(n)})$$

Maximum Likelihood Solution

Initialize W

for $x \in X$; do

$$a \leftarrow Wx$$

$$s \leftarrow \phi(a)$$

$$W \leftarrow W + \eta([W^T]^{-1} + sx^T)$$

Maximum Likelihood Solution

Initialize W

for $x \in X$; do

$$a \leftarrow Wx$$

$$s \leftarrow \phi(a)$$

$$W \leftarrow W + \eta([W^T]^{-1} + sx^T)$$

Choice of $p \iff \phi$

Maximum Likelihood Solution

Initialize W

for $x \in X$; do

$$a \leftarrow Wx$$

$$s \leftarrow \phi(a)$$

$$W \leftarrow W + \eta([W^T]^{-1} + sx^T)$$

Choice of $p \iff \phi$

Identifiable when:

1. p_i is not Gaussian (except perhaps one).
2. $p_i \perp p_j$ for all i, j

Non-linear ICA

A. Hyvarinen and P. Pajunen. 1998.

Non-linear ICA

$$p_{\theta}(s) = \prod_{i=1}^d p_i(s_i; \theta)$$

Non-linear ICA

$$p_{\theta}(s) = \prod_{i=1}^d p_i(s_i; \theta)$$

$$x = f(s; \theta); \quad f: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

Non-linear ICA

$$p_{\theta}(s) = \prod_{i=1}^d p_i(s_i; \theta)$$

$$x = f(s; \theta); \quad f: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

Goal: recover f^{-1}

Non-linear ICA

Existence: *for any random variable $x \in \mathbb{R}^d$, there exists a function $g: \mathcal{X} \rightarrow \mathcal{S}$ such that s_1, \dots, s_d have density $p_\theta(s)$ (constructive).*

$$p_\theta(s) = \prod_{i=1}^d p_i(s_i; \theta)$$

$$x = f(s; \theta); \quad f: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

Goal: recover f^{-1}

Non-linear ICA

Existence: *for any random variable $x \in \mathbb{R}^d$, there exists a function $g: \mathcal{X} \rightarrow \mathcal{S}$ such that s_1, \dots, s_d have density $p_\theta(s)$ (constructive).*

Non-uniqueness: *the number of solutions g is at least as large as the class of measure-preserving functions $h: [0,1]^n \rightarrow [0,1]^n$.*

$$p_\theta(s) = \prod_{i=1}^d p_i(s_i; \theta)$$
$$x = f(s; \theta); \quad f: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

Goal: recover f^{-1}

Nonstationary Nonlinear ICA

A. Hyvarinen and H. Morioka. 2016.

Nonstationary Nonlinear ICA

$$s_t = (s_1(t), \dots, s_n(t))$$

Nonstationary Nonlinear ICA

$$s_t = (s_1(t), \dots, s_n(t))$$

$$p_{i,\tau}(s_i) \propto \lambda_i(\tau) q_i(s_i)$$

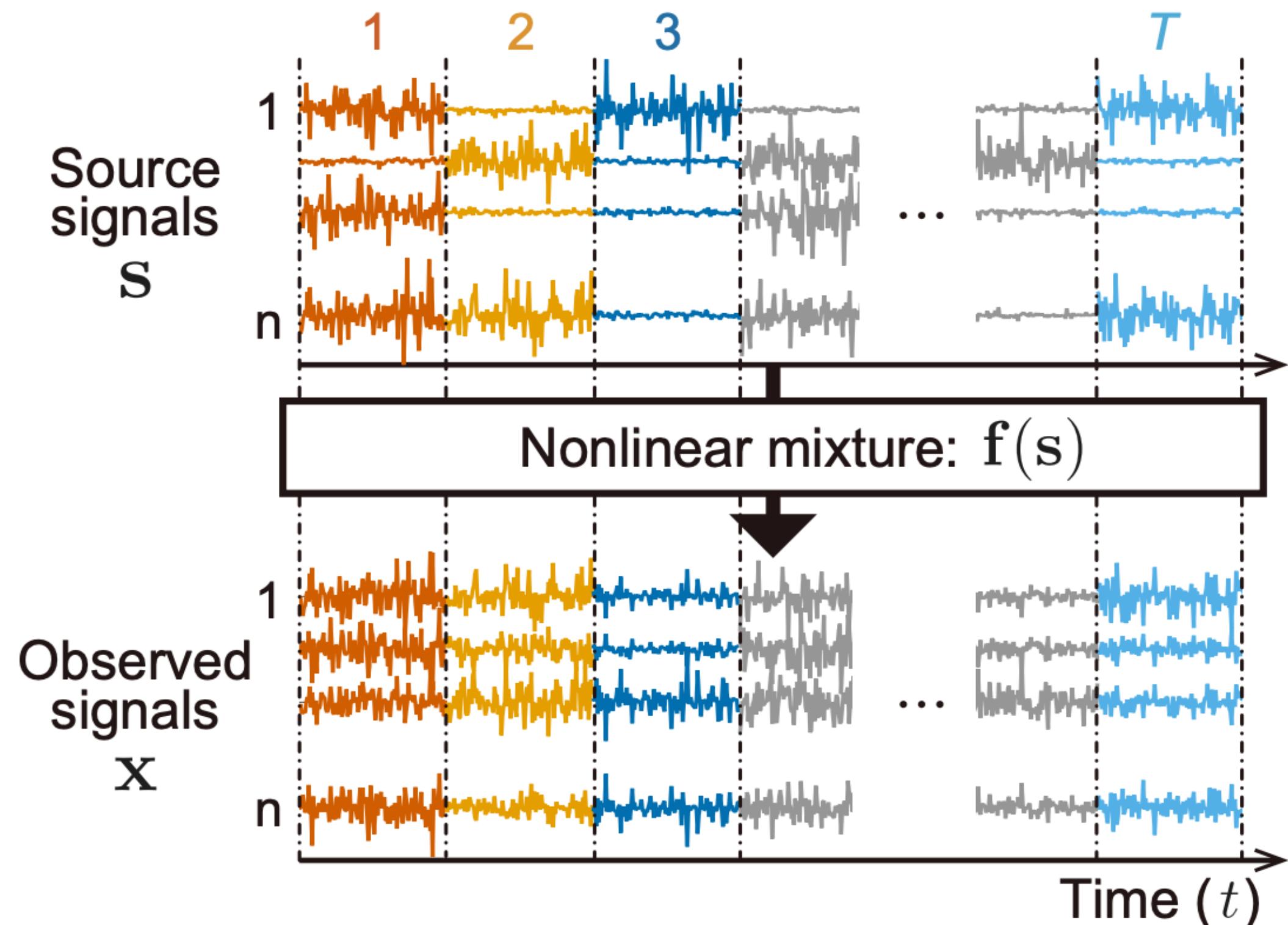
Nonstationary Nonlinear ICA

$$s_t = (s_1(t), \dots, s_n(t))$$

$$p_{i,\tau}(s_i) \propto \lambda_i(\tau) q_i(s_i)$$

$$x_t = f(s_t)$$

Nonstationary Nonlinear ICA



$$s_t = (s_1(t), \dots, s_n(t))$$

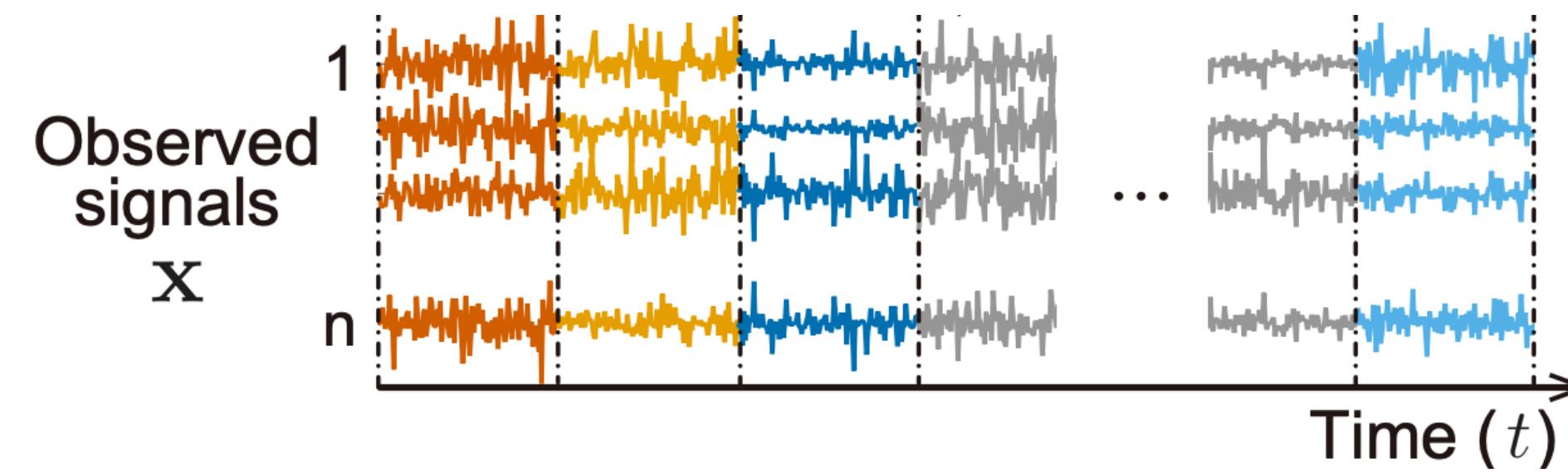
$$p_{i,\tau}(s_i) \propto \lambda_i(\tau) q_i(s_i)$$

$$x_t = f(s_t)$$

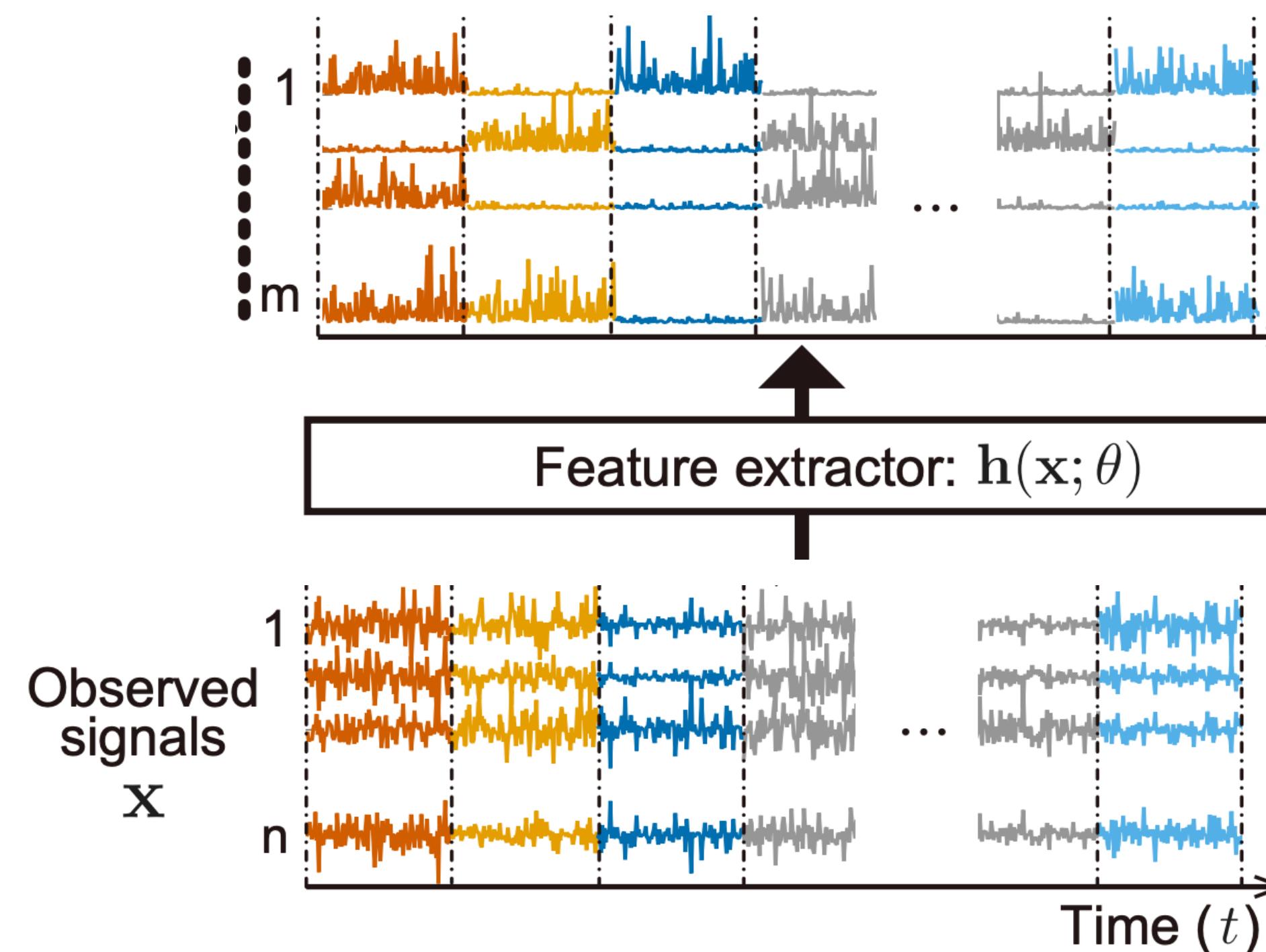
Time-Contrastive Learning

A. Hyvarinen and H. Morioka. 2016.

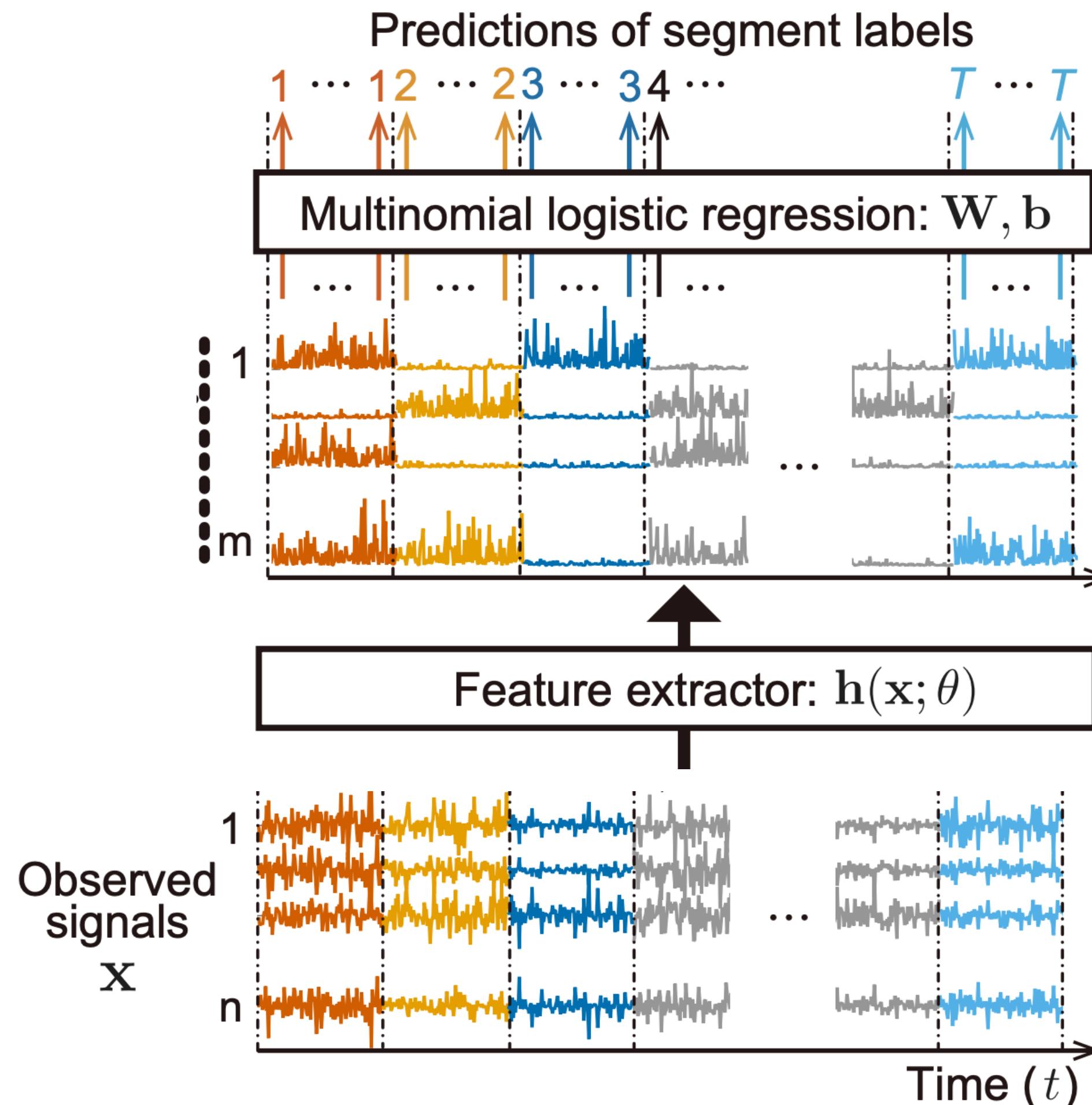
Time-Contrastive Learning



Time-Contrastive Learning



Time-Contrastive Learning



Identifiability via TCL

A. Hyvarinen and H. Morioka. 2016.

Identifiability via TCL

Theorem. Assume that

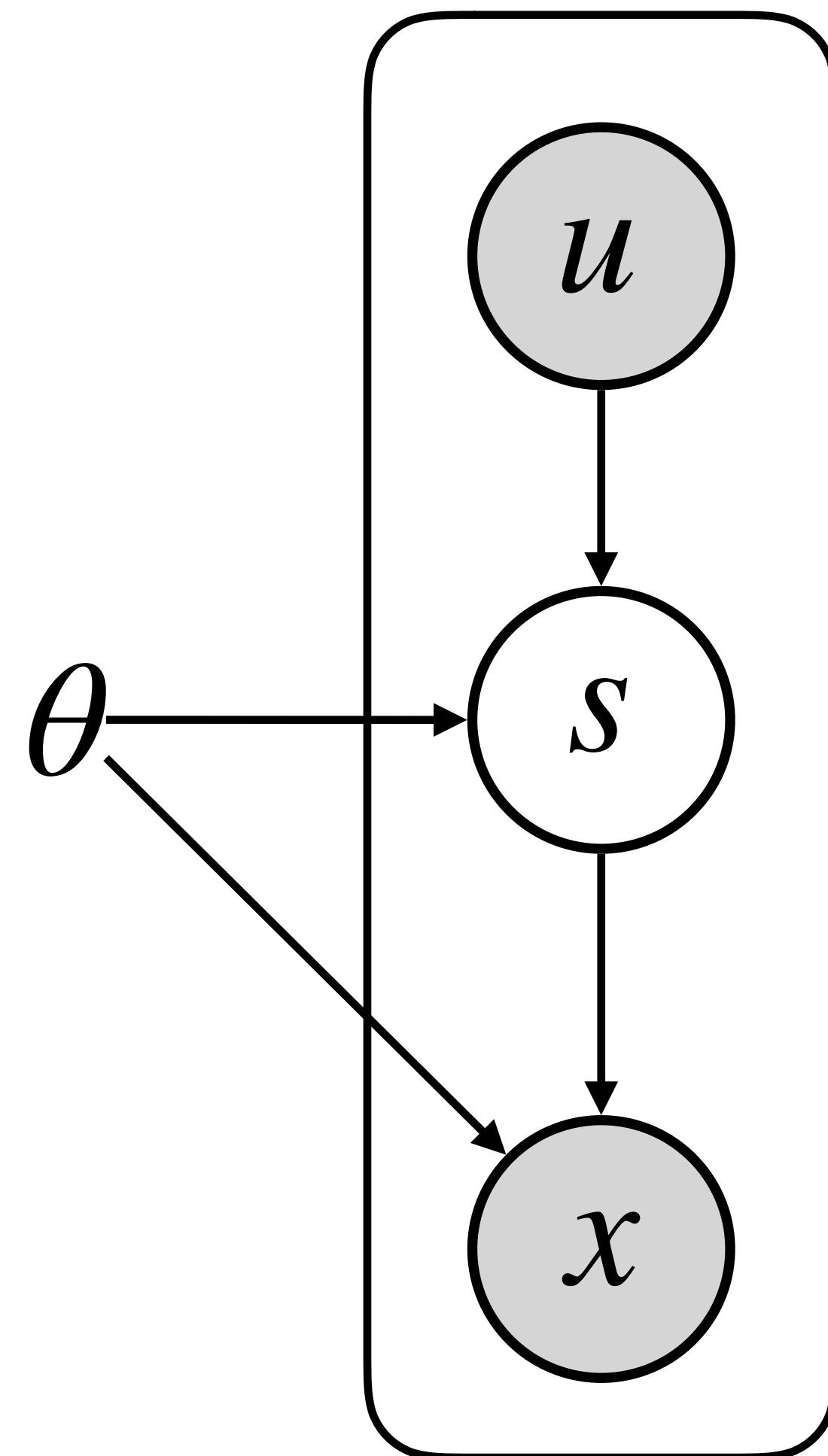
- (i) the observed data are generated from the detailed model,
- (ii) TCL is applied to learn a feature extractor $h(x_t; \theta)$, and
- (iii) the parameters $\lambda_{i,v}$ are “well-behaved”.

Then, with infinite data we have that

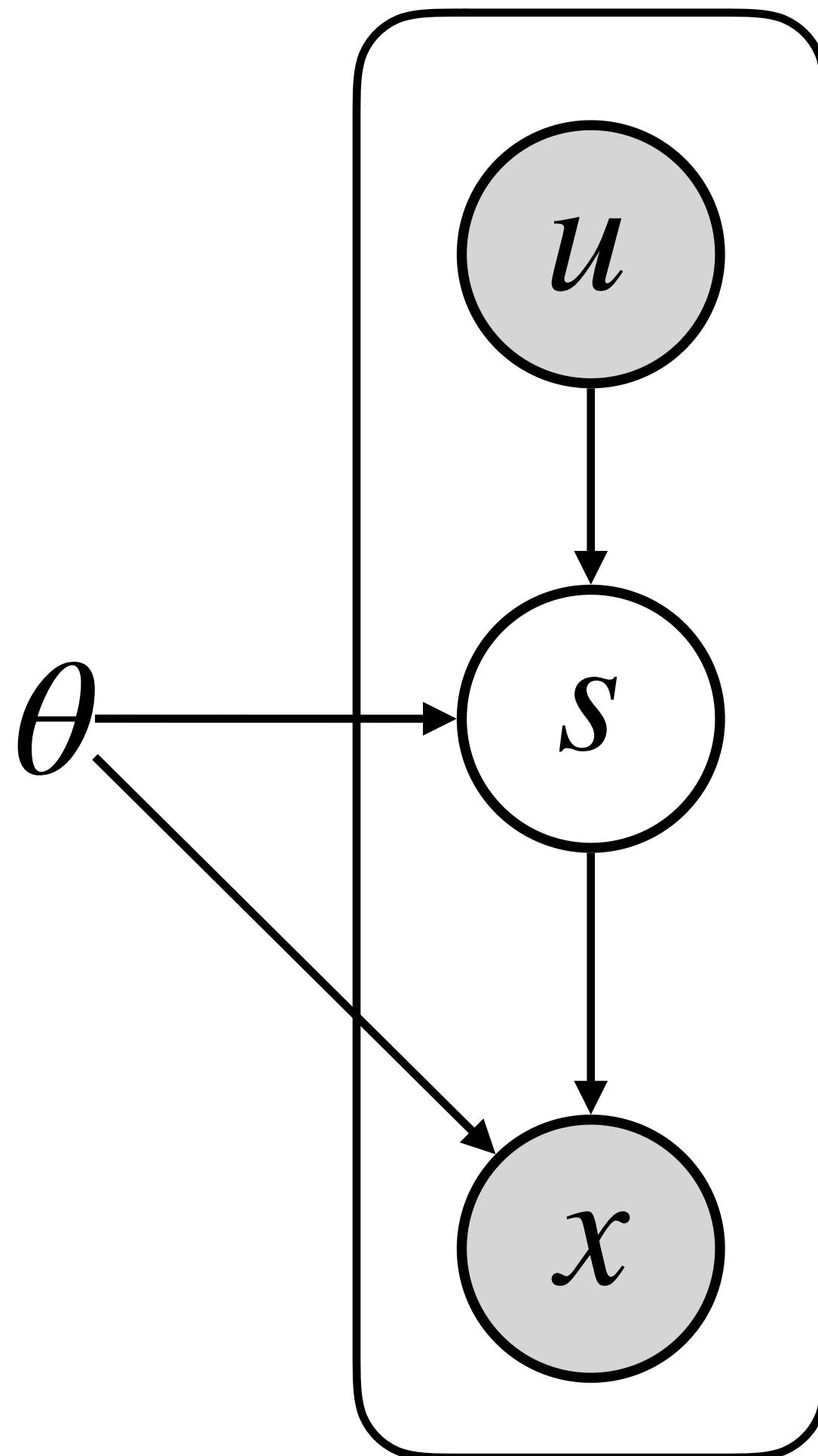
$$q(s_t) = Ah(x; \theta) + b .$$

Generalization with Auxiliary Variables

Generalization with Auxiliary Variables

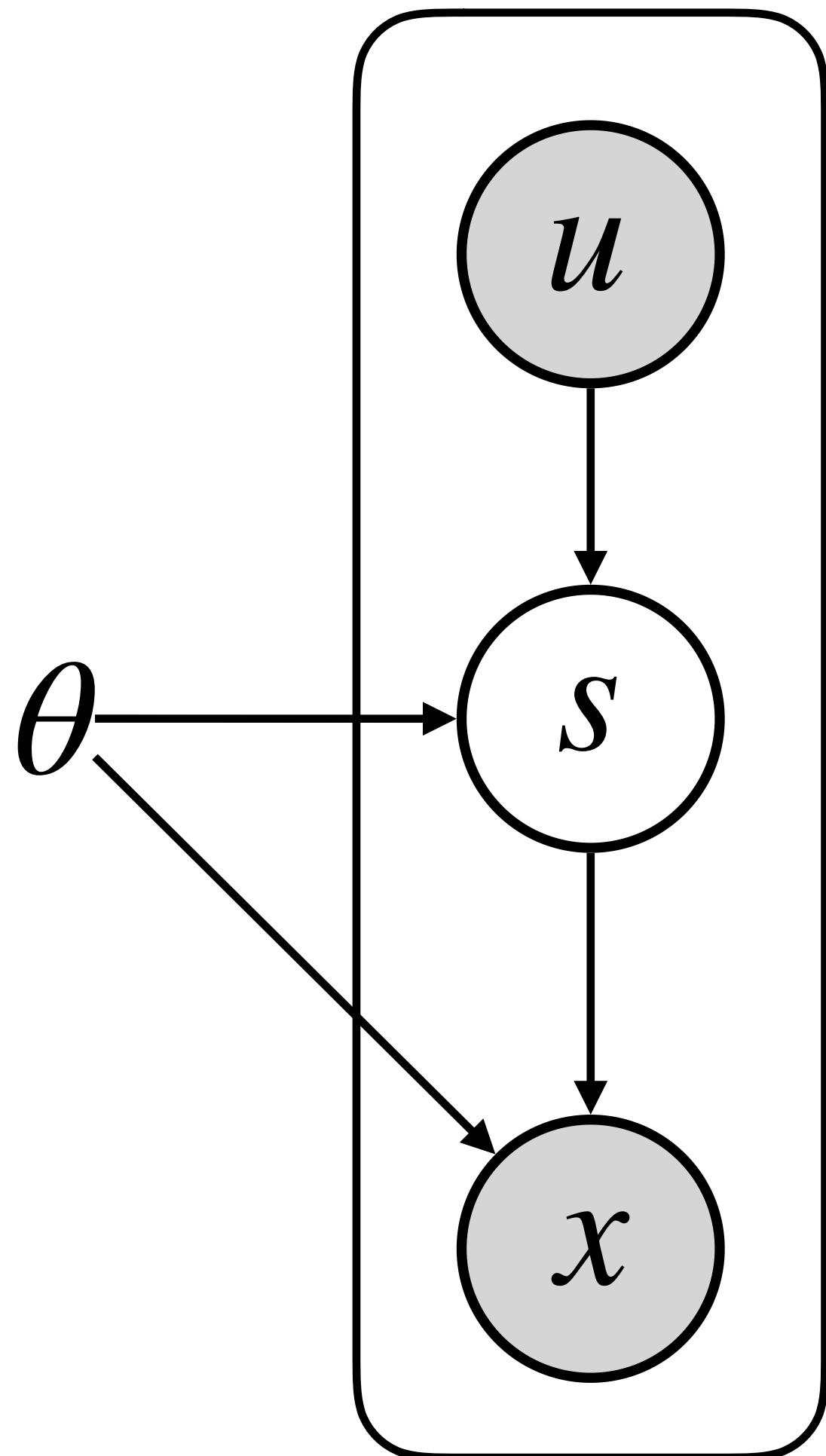


Generalization with Auxiliary Variables



$$p_{\theta}(s) = \prod_{i=1}^d p_i(s_i | u; \theta)$$
$$x = f(s; \theta); \quad f: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

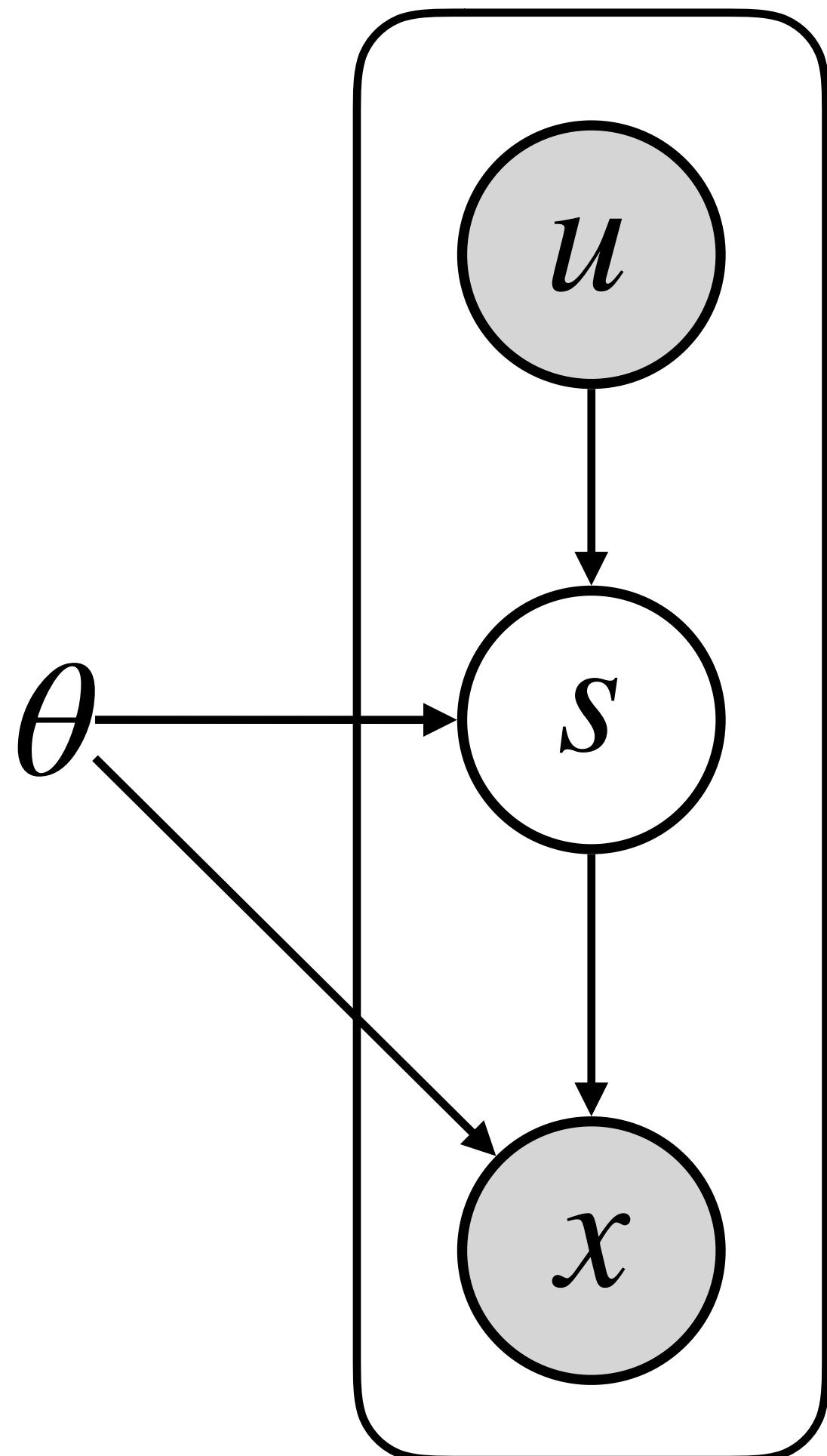
Generalization with Auxiliary Variables



$$p_{\theta}(s) = \prod_{i=1}^d p_i(s_i | u; \theta)$$
$$x = f(s; \theta); \quad f: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$\tilde{x} = (x, u); \quad \tilde{x}^* = (x, u^*)$$

Generalization with Auxiliary Variables



$$p_{\theta}(s) = \prod_{i=1}^d p_i(s_i | u; \theta)$$
$$x = f(s; \theta); \quad f: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$r(x, u) = \sum_{i=1}^n \psi_i(h(x; \theta), u)$$

$$\tilde{x} = (x, u); \quad \tilde{x}^* = (x, u^*)$$

Summary

Summary

- *Contrastive learning* –> Self-supervised “heuristics”

Summary

- *Contrastive learning –> Self-supervised “heuristics”*
- *SSL methods seem to lead to useful representations*

Summary

- *Contrastive learning –> Self-supervised “heuristics”*
- *SSL methods seem to lead to useful representations*
- *Methods largely motivated (heuristically) by mutual information*

Summary

- *Contrastive learning* → *Self-supervised “heuristics”*
- *SSL methods seem to lead to useful representations*
- *Methods largely motivated (heuristically) by mutual information*
- *Self-supervised learning <→ identifiability in generative models*

References

- F. Locatello et al. **Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations.** 2019.
- J. Shlens. **A Tutorial on Independent Component Analysis.** 2014.
- D. Mackay. **Information Theory, Inference, and Learning Algorithms.** 2003.
- A. Hyvarinen and P. Pajunen. **Nonlinear ICA: Existence and Uniqueness Results.** 1998.
- A. Hyvarinen and H. Morioka. **Unsupervised Feature Extraction by Time-Contrastive Learning and Nonlinear ICA.** 2016.
- A. Hyvarinen et al. **Nonlinear ICA Using Auxiliary Variables and Generalized Contrastive Learning.** 2019.
- I. Khemakhem et al. **VAEs and Nonlinear ICA: a Unifying Framework.** 2019.