The background of the slide is a faded, light-colored photograph of a university campus. A prominent feature is a tall, white, square tower with a pointed top, likely the Sather Tower at UC Berkeley. To the right, there is a large, classical-style building with a portico supported by columns. The sky is a pale, overcast blue.

# ICML 2021 Tutorial

## Unsupervised Learning for RL

Aravind Srinivas, Pieter Abbeel  
UC Berkeley

Part 1: Representation Learning in RL

# Structure

- Part 1:  
Representation Learning in RL (Aravind):  
How can you use Self(Un)-supervised Learning to improve RL
- Part 2:  
Reward-Free RL (Pieter):  
How can you train a completely unsupervised  
(or un-reinforced) agent



# Why Unsupervised Learning

# Why Unsupervised Learning

Please also check out Alex Graves and Marc Ranzato's Deep Unsupervised Learning (NeurIPS 2018) tutorial

# Why Unsupervised Learning

Please also check out Alex Graves and Marc Ranzato's Deep Unsupervised Learning (NeurIPS 2018) tutorial

Objective:

Introduce some of the recent UL/SSL methods and some of their RL applications to folks doing RL.

# LeCake



LeCake

“If intelligence is a cake,  
bulk of the cake is un(self-)supervised learning,  
the icing on the cake is supervised learning, and,  
the cherry on the cake is reinforcement learning.”

— *Yann LeCun*



# LeCake



LeCake

“If intelligence is a cake,  
bulk of the cake is un(self-)supervised learning,  
the icing on the cake is supervised learning, and,  
the cherry on the cake is reinforcement learning.”

— *Yann LeCun*

Definitions:

1. Unsupervised Learning: Learn from *unannotated* data
2. Supervised Learning: Learn from *annotations* (human labels)
3. Reinforcement Learning: Learn from *reward* signals

# Forms of Learning



LeCake

	<b>With teacher</b>	<b>Without teacher</b>
<b>Active</b>	Reinforcement Learning	Intrinsic Motivation (Exploration)
<b>Passive</b>	Supervised Learning	Self-(un)supervised Learning



# Forms of Learning



LeCake

	With teacher	Without teacher
Active	Reinforcement Learning	Intrinsic Motivation (Exploration)
Passive	Supervised Learning	Self-(un)supervised Learning

Learning is purely based on *extrinsic* reward optimization

# Forms of Learning



LeCake

	With teacher	Without teacher
Active	Reinforcement Learning	Intrinsic Motivation (Exploration)
Passive	Supervised Learning	Self-(un)supervised Learning

**Examples:** DQN, AlphaGo, Alphazero, OpenAI Five (Dota), AlphaStar

# Forms of Learning



LeCake

	With teacher	Without teacher
Active	Reinforcement Learning	Intrinsic Motivation (Exploration)
Passive	Supervised Learning	Self-(un)supervised Learning

Learning is purely based on *intrinsic* reward optimization

# Forms of Learning



LeCake

	With teacher	Without teacher
Active	Reinforcement Learning	Intrinsic Motivation (Exploration)
Passive	Supervised Learning	Self-(un)supervised Learning

**Examples:** Count-Based Exploration (DeepMind, Bellemare), Curiosity (Pathak, Efros), Random Network Distillation (OpenAI), Go-Explore (Uber AI)

# Forms of Learning



LeCake

	With teacher	Without teacher
Active	Reinforcement Learning	Intrinsic Motivation (Exploration)
Passive	Supervised Learning	Self-(un)supervised Learning

Learning is purely based on predicting the output from an input using *labeled (annotated)* data

# Forms of Learning



LeCake

	With teacher	Without teacher
Active	Reinforcement Learning	Intrinsic Motivation (Exploration)
Passive	Supervised Learning	Self-(un)supervised Learning

**Examples:** AlexNet, VGG, ResNet, Neural Image Captioning, seq2seq (Transformers), DeepSpeech, WaveNet, AlphaFold, DALL-E, .....

# Forms of Learning



LeCake

	<b>With teacher</b>	<b>Without teacher</b>
<b>Active</b>	Reinforcement Learning	Intrinsic Motivation (Exploration)
<b>Passive</b>	Supervised Learning	Self-(un)supervised Learning

Learn representations, or world models, or generative models, from unlabeled (un-annotated or rewarded) data.

# Forms of Learning



LeCake

	With teacher	Without teacher
Active	Reinforcement Learning	Intrinsic Motivation (Exploration)
Passive	Supervised Learning	Self-(un)supervised Learning

**Examples:** PixelCNN, GANs, VAEs, GPT-1,2,3; BERT, T5, Electra, iGPT, Contrastive Predictive Coding (CPC), Momentum Contrast (MoCo), AMDIM, CMC, SimCLR, BYOL, SimSiam, DINO, Barlow Twins, ....



# LeCake



LeCake

“If intelligence is a cake,  
bulk of the cake is un(self-)supervised learning,  
the icing on the cake is supervised learning, and,  
the cherry on the cake is reinforcement learning.”

— *Yann LeCun*

Bits Argument:

1. Unsupervised Learning: Predict missing from given (*millions* of bits)
2. Supervised Learning: Predict human-annotations (*thousands* of bits)
3. Reinforcement Learning: Predict scalar rewards (*fewer* bits)

# Practical arguments in favor of unsupervised learning



LeCake

# Practical arguments in favor of unsupervised learning



LeCake

1. Label creation, annotation and maintenance is time-consuming and a challenging discipline in itself.

# Practical arguments in favor of unsupervised learning



LeCake

1. Label creation, annotation and maintenance is time-consuming and a challenging discipline in itself.
2. Some domains require significant expertise in annotation (medicine, law, ethics, science, biology, etc.)

# Practical arguments in favor of unsupervised learning



LeCake

1. Label creation, annotation and maintenance is time-consuming and a challenging discipline in itself.
2. Some domains require significant expertise in annotation (medicine, law, ethics, science, biology, etc.)
3. Reward annotation is even trickier - sparse or dense, what's the right UI for scalable annotation of rewards by humans, should you use a continuous-valued reward or categorical (good/bad/neural), etc.

# Practical arguments in favor of unsupervised learning



LeCake

1. Label creation, annotation and maintenance is time-consuming and a challenging discipline in itself.
2. Some domains require significant expertise in annotation (medicine, law, ethics, science, biology, etc.)
3. Reward annotation is even trickier - sparse or dense, what's the right UI for scalable annotation of rewards by humans, should you use a continuous-valued reward or categorical (good/bad/neural), etc.
4. Collecting human demos for behavior cloning is challenging at scale

# Practical arguments in favor of unsupervised learning



LeCake

1. Label creation, annotation and maintenance is time-consuming and a challenging discipline in itself.
2. Some domains require significant expertise in annotation (medicine, law, ethics, science, biology, etc.)
3. Reward annotation is even trickier - sparse or dense, what's the right UI for scalable annotation of rewards by humans, should you use a continuous-valued reward or categorical (good/bad/neural), etc.
4. Collecting human demos for behavior cloning is challenging at scale
5. Behavior cloning for real world problems prone to distribution mismatch and compounding errors without sufficient data

# Practical arguments in favor of unsupervised learning



LeCake

1. Label creation, annotation and maintenance is time-consuming and a challenging discipline in itself.
2. Some domains require significant expertise in annotation (medicine, law, ethics, science, biology, etc.)
3. Reward annotation is even trickier - sparse or dense, what's the right UI for scalable annotation of rewards by humans, should you use a continuous-valued reward or categorical (good/bad/neural), etc.
4. Collecting human demos for behavior cloning is challenging at scale
5. Behavior cloning for real world problems prone to distribution mismatch and compounding errors without sufficient data
6. Good (useful) behavioral data *even without* annotations is quite challenging to collect.



# Philosophical arguments in favor of unsupervised learning



LeCake

# Philosophical arguments in favor of unsupervised learning



LeCake

1. Inspired by how human infants learn

# Philosophical arguments in favor of unsupervised learning



LeCake

1. Inspired by how human infants learn
2. Feels more *human-like*

# Philosophical arguments in favor of unsupervised learning



LeCake

1. Inspired by how human infants learn
2. Feels more *human-like*
3. Build mental models of the world (Kenneth Craik) - “*mind is a predictive modelling engine*”, “*organism carries a ‘small-scale model’ of external reality and of its own possible actions within its head*”

# Philosophical arguments in favor of unsupervised learning



LeCake

1. Inspired by how human infants learn
2. Feels more *human-like*
3. Build mental models of the world (Kenneth Craik) - "*mind is a predictive modelling engine*", "*organism carries a 'small-scale model' of external reality and of its own possible actions within its head*"
4. Learn skills, not tasks (Satinder Singh)

# Philosophical arguments in favor of unsupervised learning



LeCake

1. Inspired by how human infants learn
2. Feels more *human-like*
3. Build mental models of the world (Kenneth Craik) - “*mind is a predictive modelling engine*”, “*organism carries a ‘small-scale model’ of external reality and of its own possible actions within its head*”
4. Learn skills, not tasks (Satinder Singh)
5. Build a *general-purpose* understanding of the world (representations, world models, common sense, etc etc) to be able to **generalize** well to new scenarios and learn fast (**transfer**)

# Bits Argument



LeCake

# Bits Argument

1. Consider ImageNet (1.28M images) with 1000 categories, 256x256.



LeCake



# Bits Argument



LeCake

1. Consider ImageNet (1.28M images) with 1000 categories, 256x256.
2. Supervised:  $\log_2(1000) * 1.28M \sim 12.8$  Mbits

# Bits Argument



LeCake

1. Consider ImageNet (1.28M images) with 1000 categories, 256x256.
2. Supervised:  $\log_2(1000) * 1.28M \sim 12.8$  Mbits
3. Unsupervised:  $256x256x3x8 \sim 2000$  Gbits, 150000x more.

# Bits Argument



LeCake

1. Consider ImageNet (1.28M images) with 1000 categories, 256x256.
2. Supervised:  $\log_2(1000) * 1.28M \sim 12.8$  Mbits
3. Unsupervised:  $256x256x3x8 \sim 2000$  Gbits, 150000x more.
4. Flaw in the argument: Not all bits are equal, a 8-bit color-channel information conveys much less than a human-language aligned categorical classification bit.

# Bits Argument



LeCake

1. Consider ImageNet (1.28M images) with 1000 categories, 256x256.
2. Supervised:  $\log_2(1000) * 1.28M \sim 12.8$  Mbits
3. Unsupervised:  $256x256x3x8 \sim 2000$  Gbits, 150000x more.
4. Flaw in the argument: Not all bits are equal, a 8-bit color-channel information conveys much less than a human-language aligned categorical classification bit.
5. Somehow, *intrinsically capture* the bits *that matter* for *downstream tasks* and *rely less* on human-provided bits

# Bits Argument



LeCake

1. Consider ImageNet (1.28M images) with 1000 categories, 256x256.
2. Supervised:  $\log_2(1000) * 1.28M \sim 12.8$  Mbits
3. Unsupervised:  $256x256x3x8 \sim 2000$  Gbits, 150000x more.
4. Flaw in the argument: Not all bits are equal, a 8-bit color-channel information conveys much less than a human-language aligned categorical classification bit.
5. Somehow, *intrinsically capture* the bits *that matter* for *downstream tasks* and *rely less* on human-provided bits
6. Use self-supervised learning to improve the label and reward efficiency of supervised and reinforcement learning systems

# Bits Argument



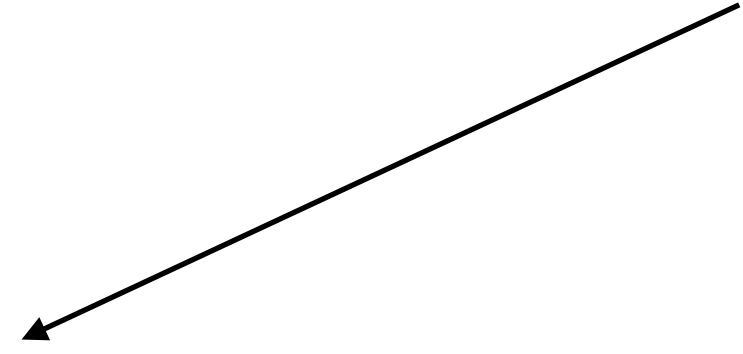
LeCake

1. Consider ImageNet (1.28M images) with 1000 categories, 256x256.
2. Supervised:  $\log_2(1000) * 1.28M \sim 12.8$  Mbits
3. Unsupervised:  $256x256x3x8 \sim 2000$  Gbits, 150000x more.
4. Flaw in the argument: Not all bits are equal, a 8-bit color-channel information conveys much less than a human-language aligned categorical classification bit.
5. Somehow, *intrinsically capture* the bits *that matter* for *downstream tasks* and *rely less* on human-provided bits
6. Use self-supervised learning to improve the label and reward efficiency of supervised and reinforcement learning systems

↓  
**Focus of this tutorial**

# Forms of Unsupervised Learning

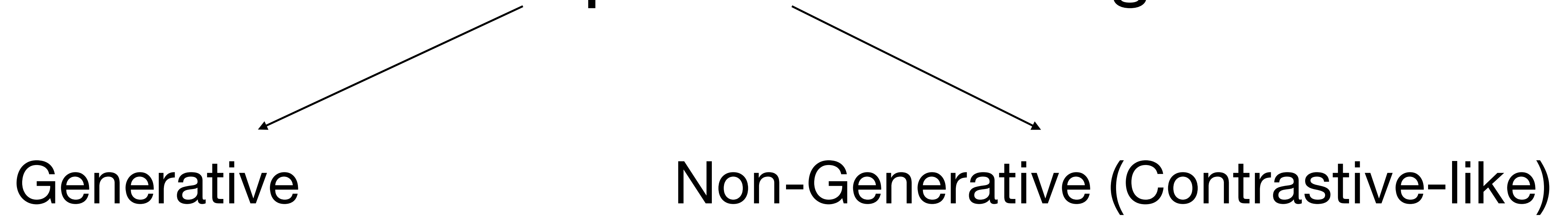
# Forms of Unsupervised Learning



Generative



# Forms of Unsupervised Learning

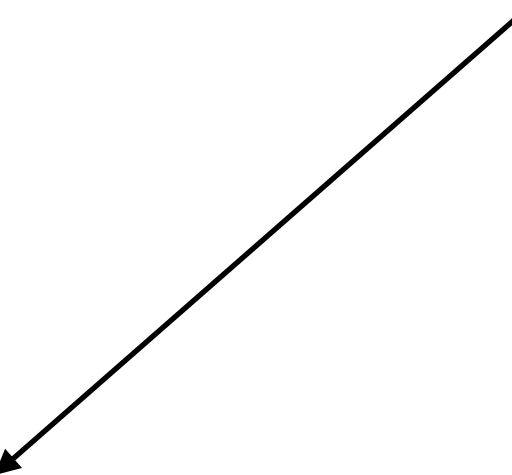


# Forms of Unsupervised Learning



Generative

Non-Generative (Contrastive-like)



Density Modeling

# Forms of Unsupervised Learning

```
graph TD; A[Forms of Unsupervised Learning] --> B[Generative]; A --> C[Non-Generative (Contrastive-like)]; B --> D["Density Modeling (PixelCNN, GPT-x, Bigan (implicit))"];
```

Generative

Non-Generative (Contrastive-like)

Density Modeling  
(PixelCNN, GPT-x,  
Bigan (implicit))

# Forms of Unsupervised Learning

```
graph TD; A[Forms of Unsupervised Learning] --> B[Generative]; A --> C[Non-Generative (Contrastive-like)]; B --> D["Density Modeling (PixelCNN, GPT-x, Bigan (implicit))"]; B --> E[Masked Auto-Encoding];
```

Generative

Non-Generative (Contrastive-like)

Density Modeling

Masked Auto-Encoding

(PixelCNN, GPT-x,  
Bigan (implicit))

# Forms of Unsupervised Learning

```
graph TD; A[Forms of Unsupervised Learning] --> B[Generative]; A --> C[Non-Generative (Contrastive-like)]; B --> D["Density Modeling  
(PixelCNN, GPT-x,  
Bigan (implicit))"]; B --> E["Masked Auto-Encoding  
(BERT, Electra)"];
```

Generative

Non-Generative (Contrastive-like)

Density Modeling  
(PixelCNN, GPT-x,  
Bigan (implicit))

Masked Auto-Encoding  
(BERT, Electra)

# Forms of Unsupervised Learning

```
graph TD; A[Forms of Unsupervised Learning] --> B[Generative]; A --> C[Non-Generative (Contrastive-like)]; B --> D["Density Modeling (PixelCNN, GPT-x, Bigan (implicit))"]; B --> E["Masked Auto-Encoding (BERT, Electra)"]; C --> F[Siamese Networks]
```

Generative

Non-Generative (Contrastive-like)

Density Modeling  
(PixelCNN, GPT-x,  
Bigan (implicit))

Masked Auto-Encoding  
(BERT, Electra)

Siamese Networks

# Forms of Unsupervised Learning

## Generative

Density Modeling  
(PixelCNN, GPT-x,  
Bigan (implicit))

Masked Auto-Encoding  
(BERT, Electra)

## Non-Generative (Contrastive-like)

Siamese Networks  
(SimCLR, MoCo, AMDIM,  
BYOL, SimSiam, DINO,  
Barlow Twins ...)

# Forms of Unsupervised Learning

```
graph TD; A[Forms of Unsupervised Learning] --> B[Generative]; A --> C[Non-Generative (Contrastive-like)]; B --> D[Density Modeling]; B --> E[Masked Auto-Encoding]; C --> F[Siamese Networks]; C --> G[Contrastive Prediction];
```

## Generative

## Non-Generative (Contrastive-like)

Density Modeling  
(PixelCNN, GPT-x,  
Bigan (implicit))

Masked Auto-Encoding  
(BERT, Electra)

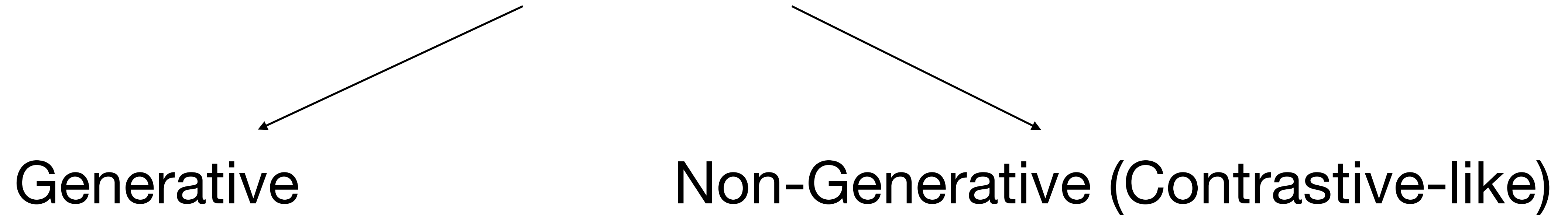
Siamese Networks  
(SimCLR, MoCo, AMDIM,  
BYOL, SimSiam, DINO,  
Barlow Twins ...)

Contrastive Prediction  
(CPC)



**How best to use UL for RL?**

# How best to use UL for RL?



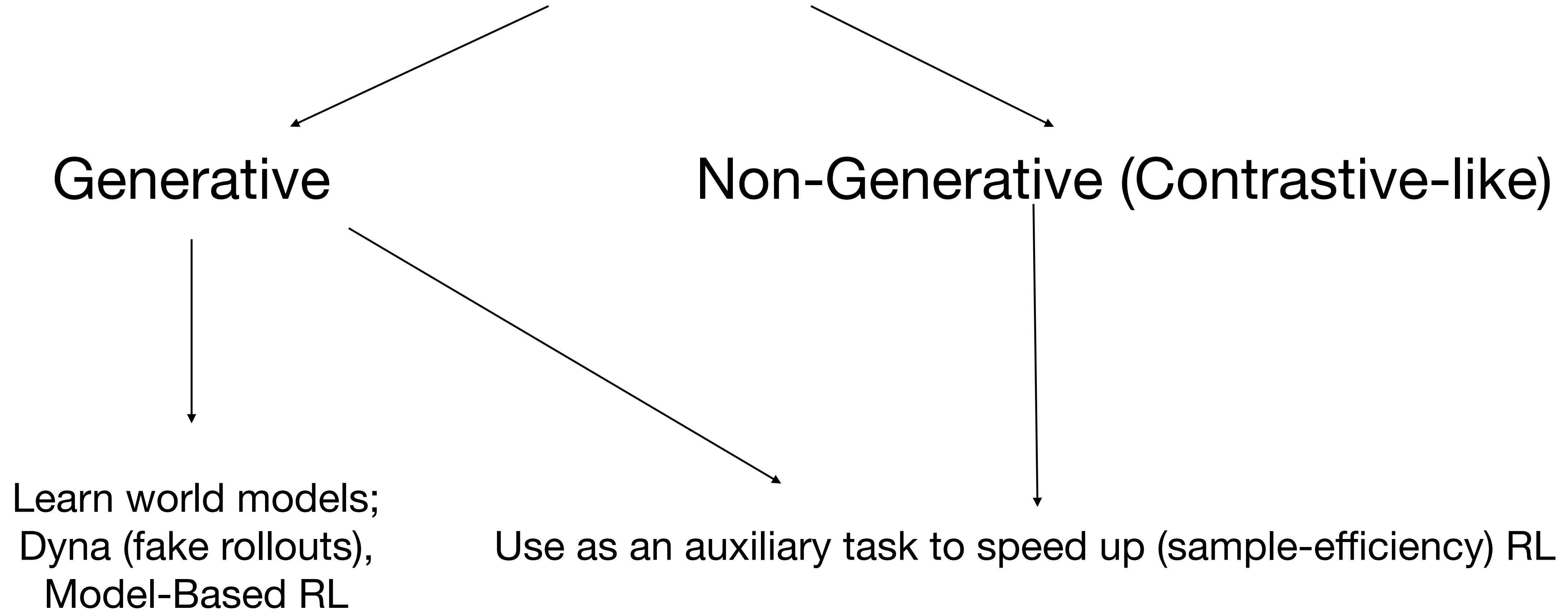
# How best to use UL for RL?

Generative

Non-Generative (Contrastive-like)

Learn world models;  
Dyna (fake rollouts),  
Model-Based RL

# How best to use UL for RL?



# How best to use UL for RL?

Generative

Non-Generative (Contrastive-like)

Learn world models;  
Dyna (fake rollouts),  
Model-Based RL

Use as an auxiliary task to speed up (sample-efficiency) RL

Beyond the scope of this tutorial;  
but will highlight some representative work,  
Refer to Mordatch/Hamrick Tutorial (ICML 20')

# How best to use UL for RL?

Generative

Non-Generative (Contrastive-like)

Learn world models;  
Dyna (fake rollouts),  
Model-Based RL

Use as an auxiliary task to speed up (sample-efficiency) RL

**Focus of this tutorial**

Beyond the scope of this tutorial;  
but will highlight some representative work,  
Refer to Mordatch/Hamrick Tutorial (ICML 20')

Why use UL as an auxiliary task to speed up RL?

# Why use UL as an auxiliary task to speed up RL?

1. **Simplest** way to combine UL/SSL and RL



# Why use UL as an auxiliary task to speed up RL?

1. **Simplest** way to combine UL/SSL and RL
2. Maximize reward, and learn about the world, using the **same shared network** - design choices: how many parameters for UL, how many for RL, etc. but simplest way is to just use all for both.

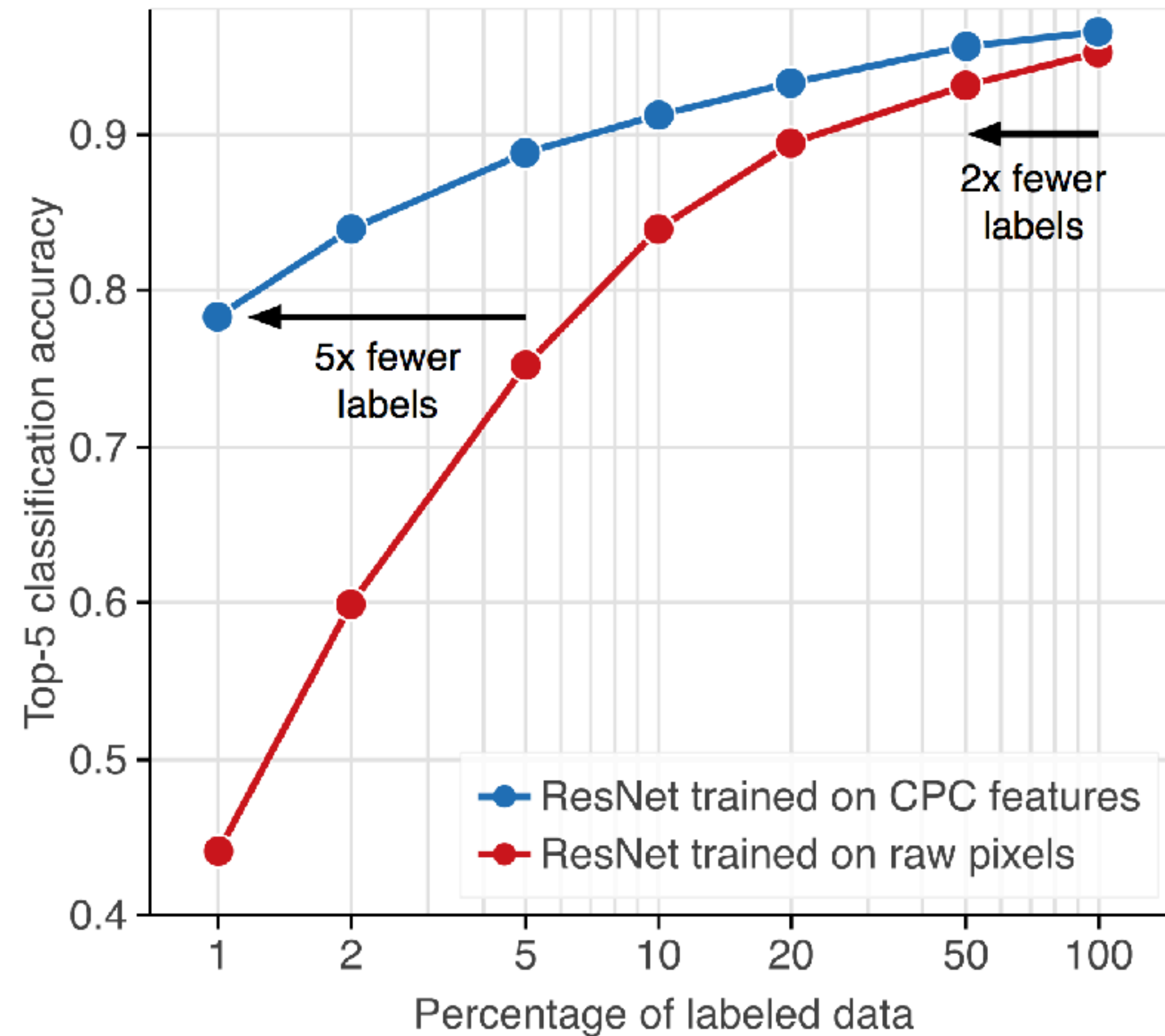
# Why use UL as an auxiliary task to speed up RL?

1. **Simplest** way to combine UL/SSL and RL
2. Maximize reward, and learn about the world, using the **same shared network** - design choices: how many parameters for UL, how many for RL, etc. but simplest way is to just use all for both.
3. Hope: Representations learned with UL/SSL **help** for the RL task(s)

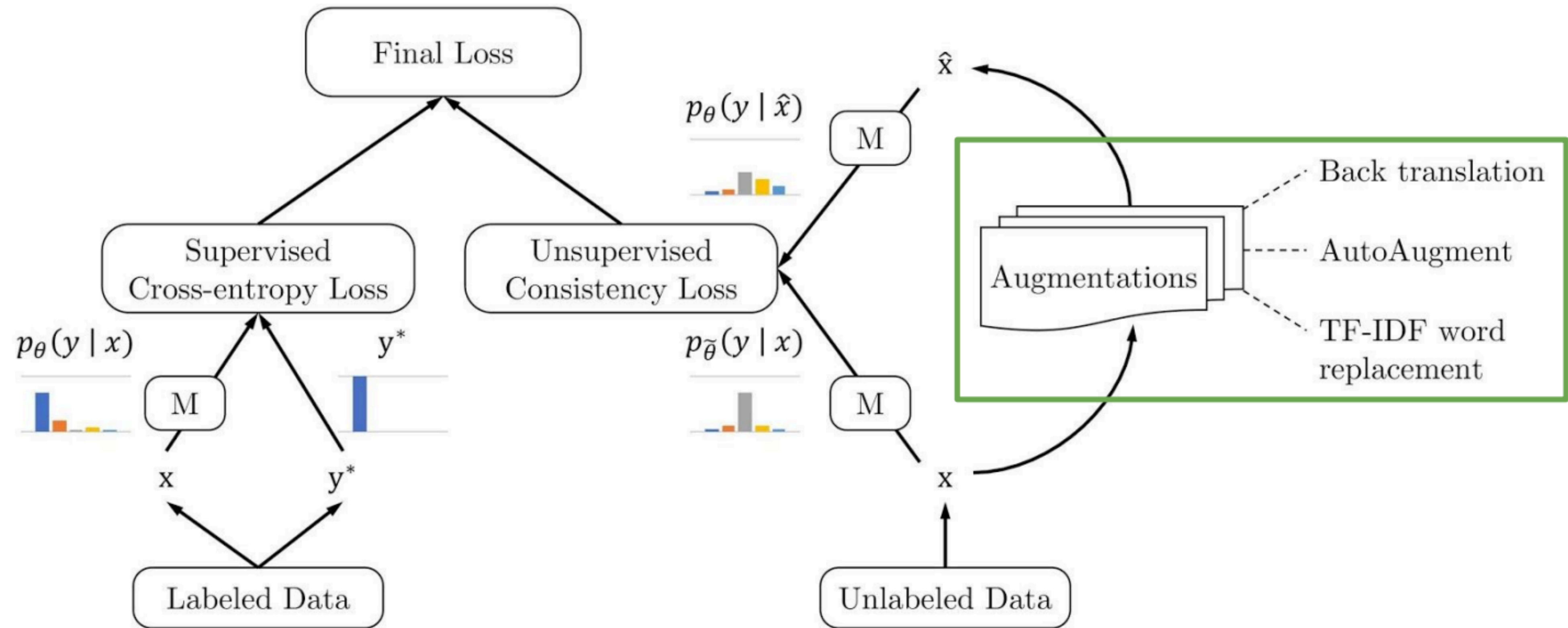
# Why use UL as an auxiliary task to speed up RL?

1. **Simplest** way to combine UL/SSL and RL
2. Maximize reward, and learn about the world, using the **same shared network** - design choices: how many parameters for UL, how many for RL, etc. but simplest way is to just use all for both.
3. Hope: Representations learned with UL/SSL **help** for the RL task(s)
4. Inspired by: **Success** of **Un/Self/Semi**-Supervised Learning for **label-efficient** Supervised Learning

# Inspiration from success in supervised learning

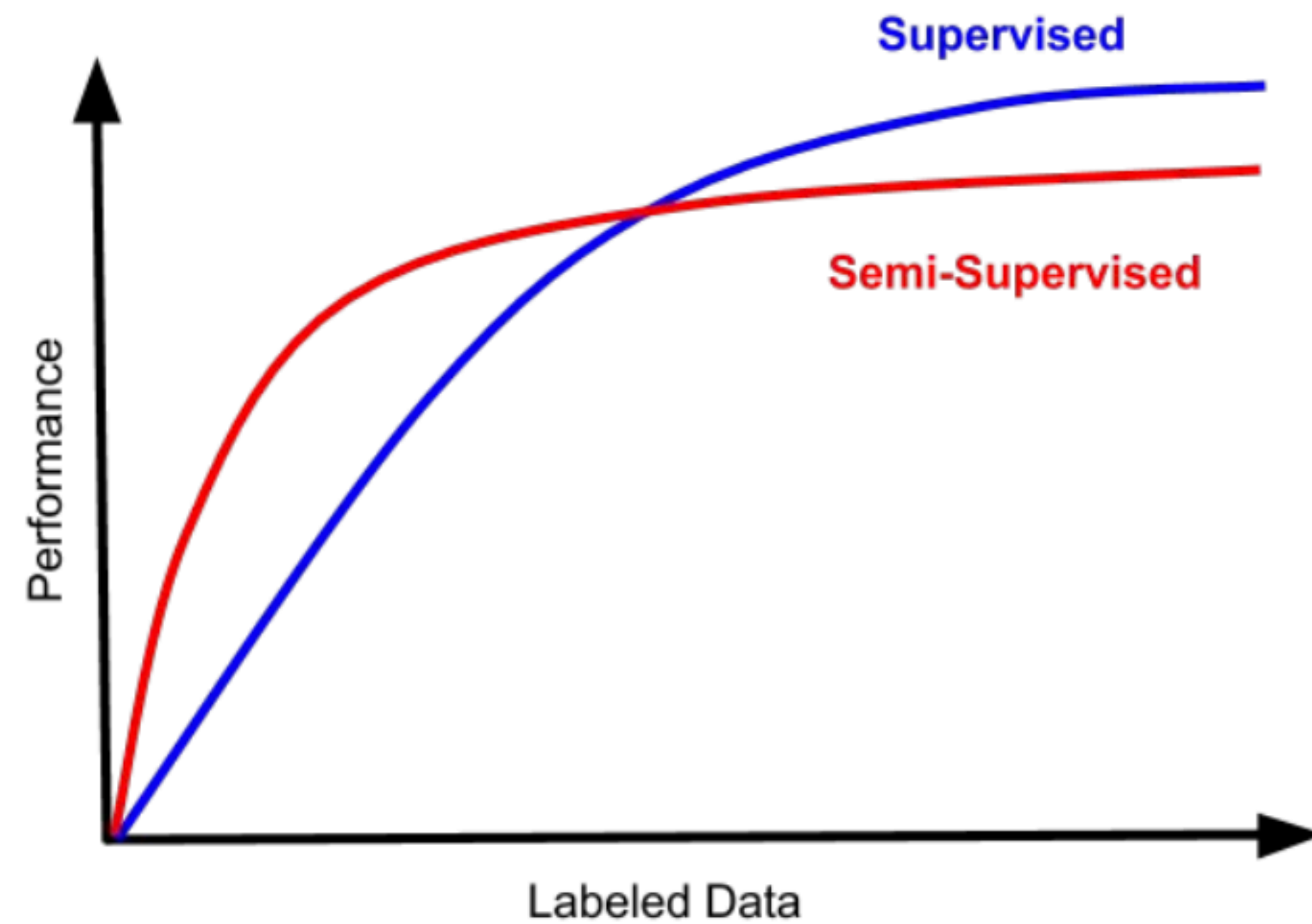


# Inspiration from success in supervised learning

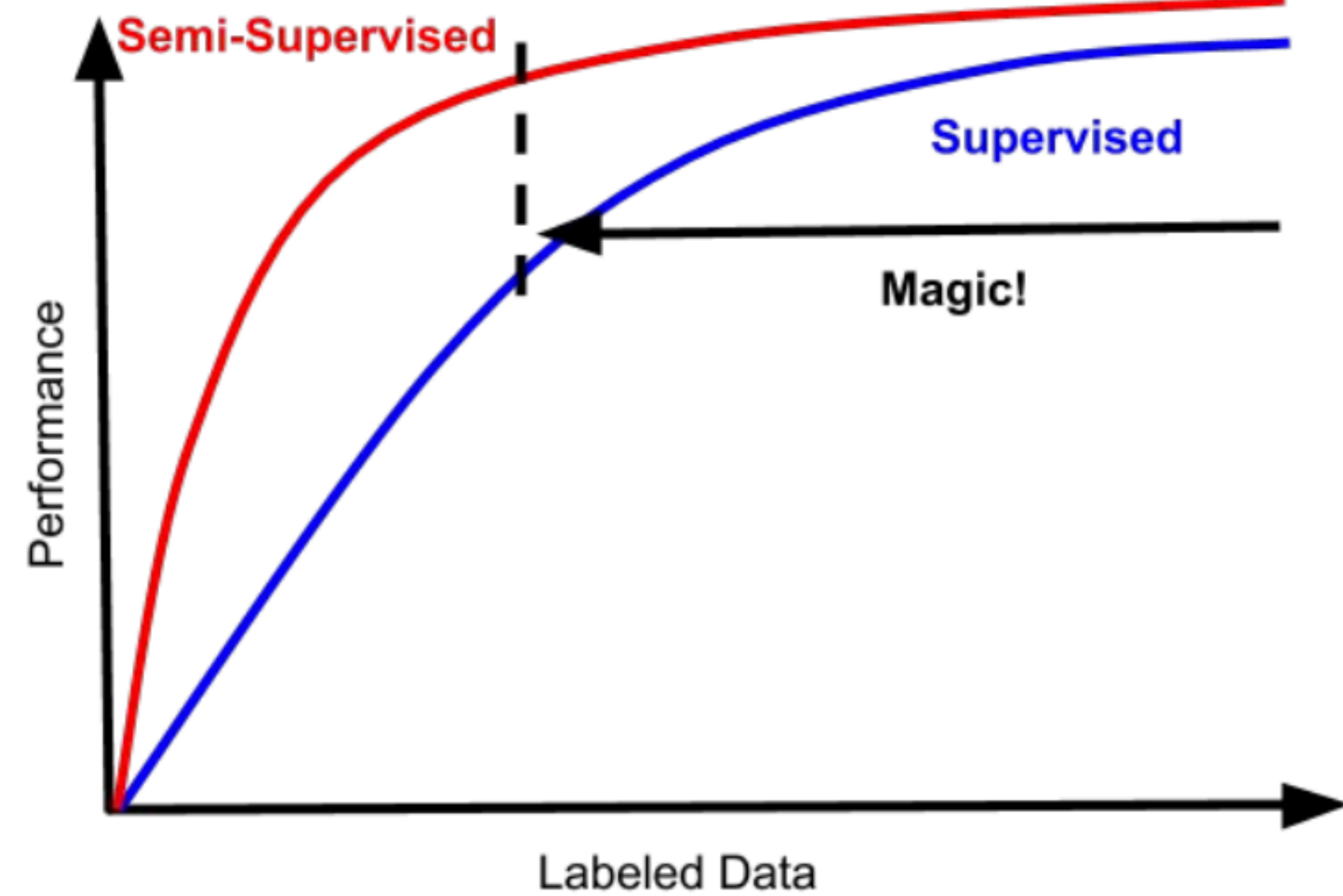


Unsupervised Data Augmentation (Xie et al 2019), figure from Thang Luong.

# Inspiration from success in supervised learning



Belief of many ML practitioners



Belief of many SSL researchers

Figure credit: Vincent Vanhoucke

# Why use UL as an auxiliary task to speed up RL?

1. **Simplest** way to combine UL/SSL and RL
2. Maximize reward, and learn about the world, using the **same shared network** - design choices: how many parameters for UL, how many for RL, etc. but simplest way is to just use all for both.
3. Hope: Representations learned with UL/SSL **help** for the RL task(s)
4. Inspired by: **Success** of **Un/Self/Semi**-Supervised Learning for **label-efficient** Supervised Learning

# Why use UL as an auxiliary task to speed up RL?

1. **Simplest** way to combine UL/SSL and RL
2. Maximize reward, and learn about the world, using the **same shared network** - design choices: how many parameters for UL, how many for RL, etc. but simplest way is to just use all for both.
3. Hope: Representations learned with UL/SSL **help** for the RL task(s)
4. Inspired by: **Success** of **Un/Self/Semi**-Supervised Learning for **label-efficient** Supervised Learning
5. Challenge: What is the right kind of UL objective that will work well in tandem with RL? How to **capture** the **useful** aspects of a high dimensional sensory stream?



# Why use UL as an auxiliary task to speed up RL?

1. **Simplest** way to combine UL/SSL and RL
2. Maximize reward, and learn about the world, using the **same shared network** - design choices: how many parameters for UL, how many for RL, etc. but simplest way is to just use all for both.
3. Hope: Representations learned with UL/SSL **help** for the RL task(s)
4. Inspired by: **Success** of **Un/Self/Semi**-Supervised Learning for **label-efficient** Supervised Learning
5. Challenge: What is the right kind of UL objective that will work well in tandem with RL? How to **capture** the **useful** aspects of a high dimensional sensory stream?

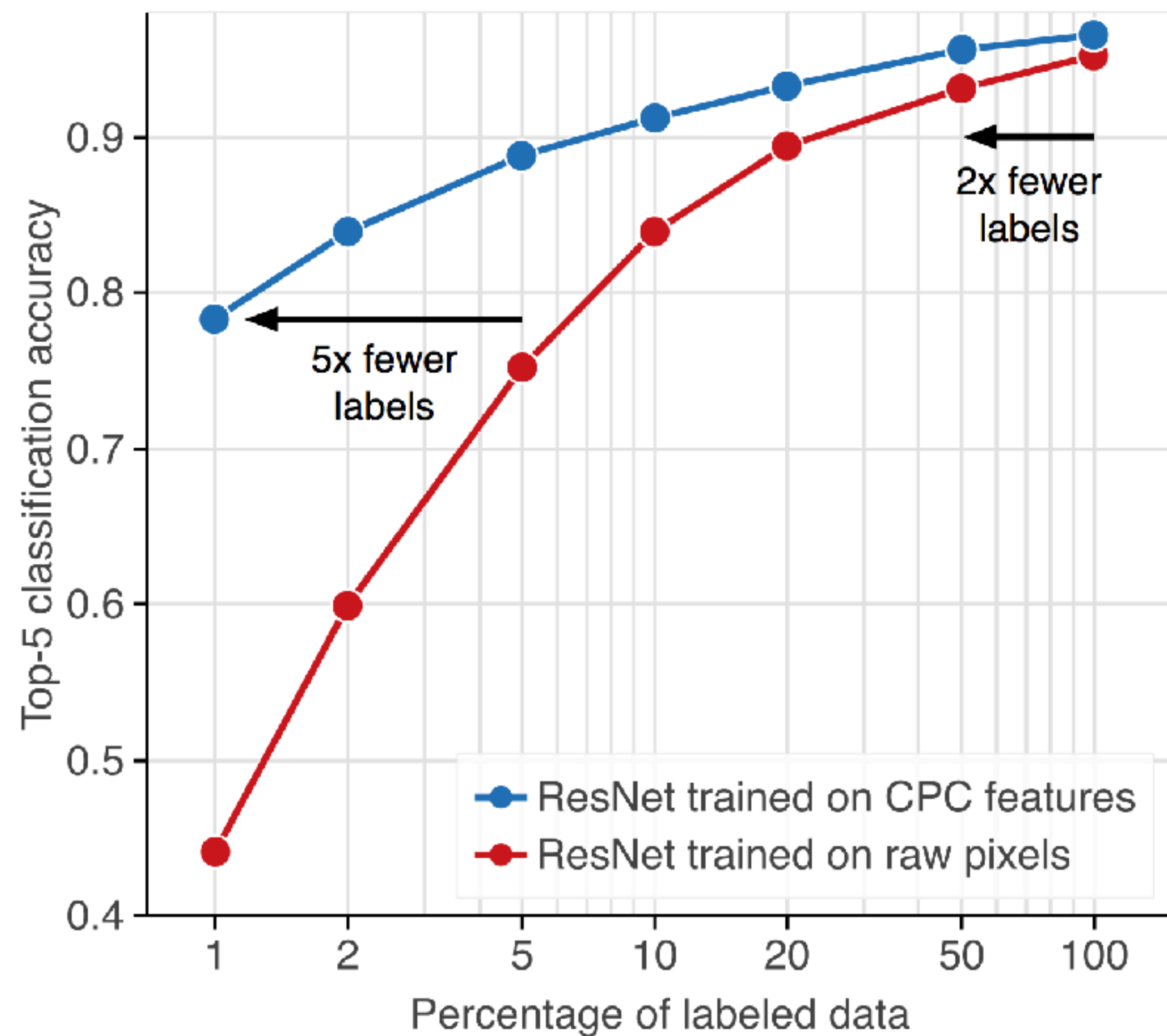
*... we lived our lives under the constantly changing sky without sparing it a glance or thought.  
And why indeed should we? If the various formations had some meaning, if, for example, there had  
been some concealed signs and messages for us which it was important to decode correctly, unceasing  
attention to what was happening would have been inescapable*

- Karl Ove Knausgaard, A Death in the Family

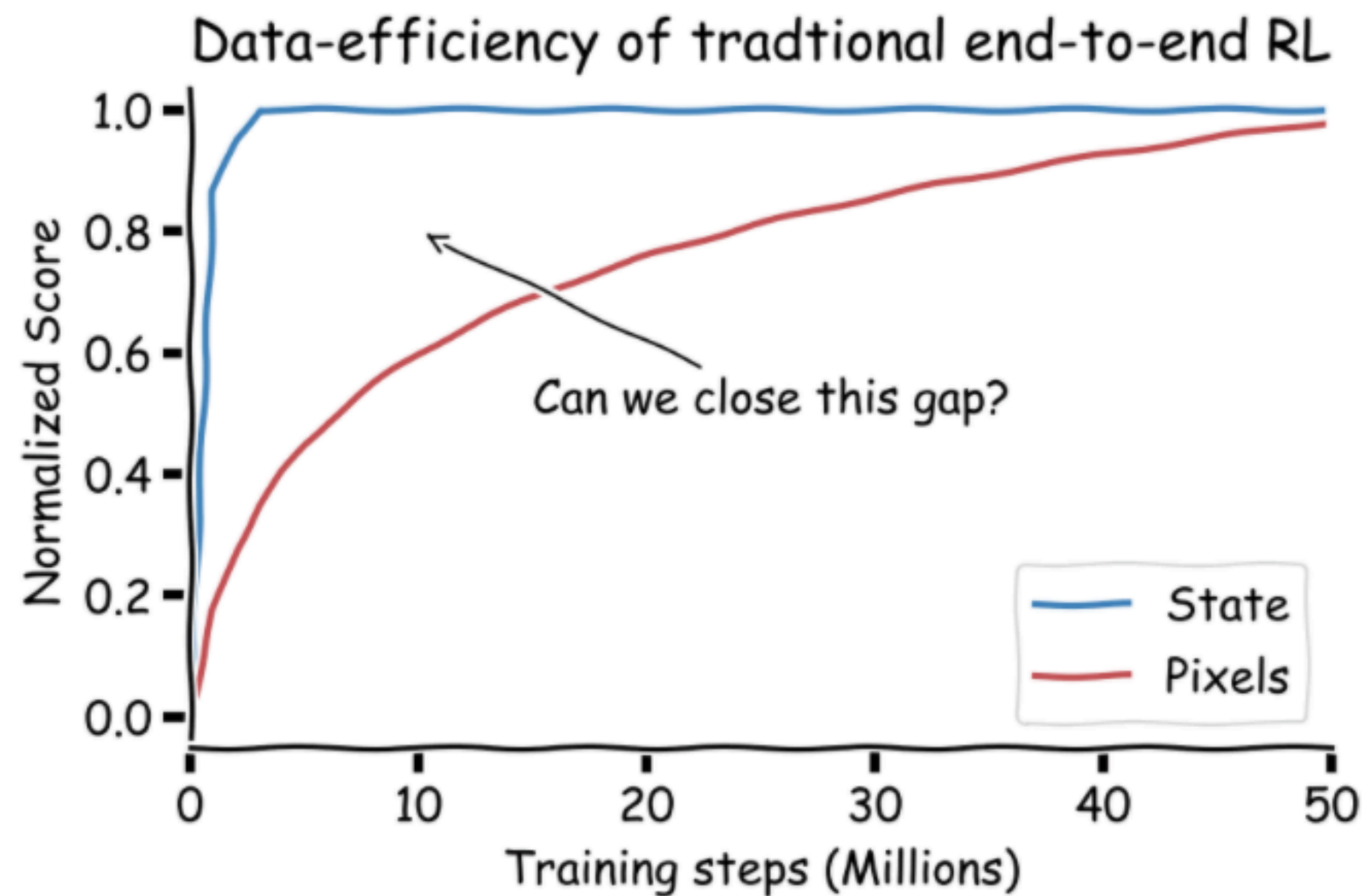
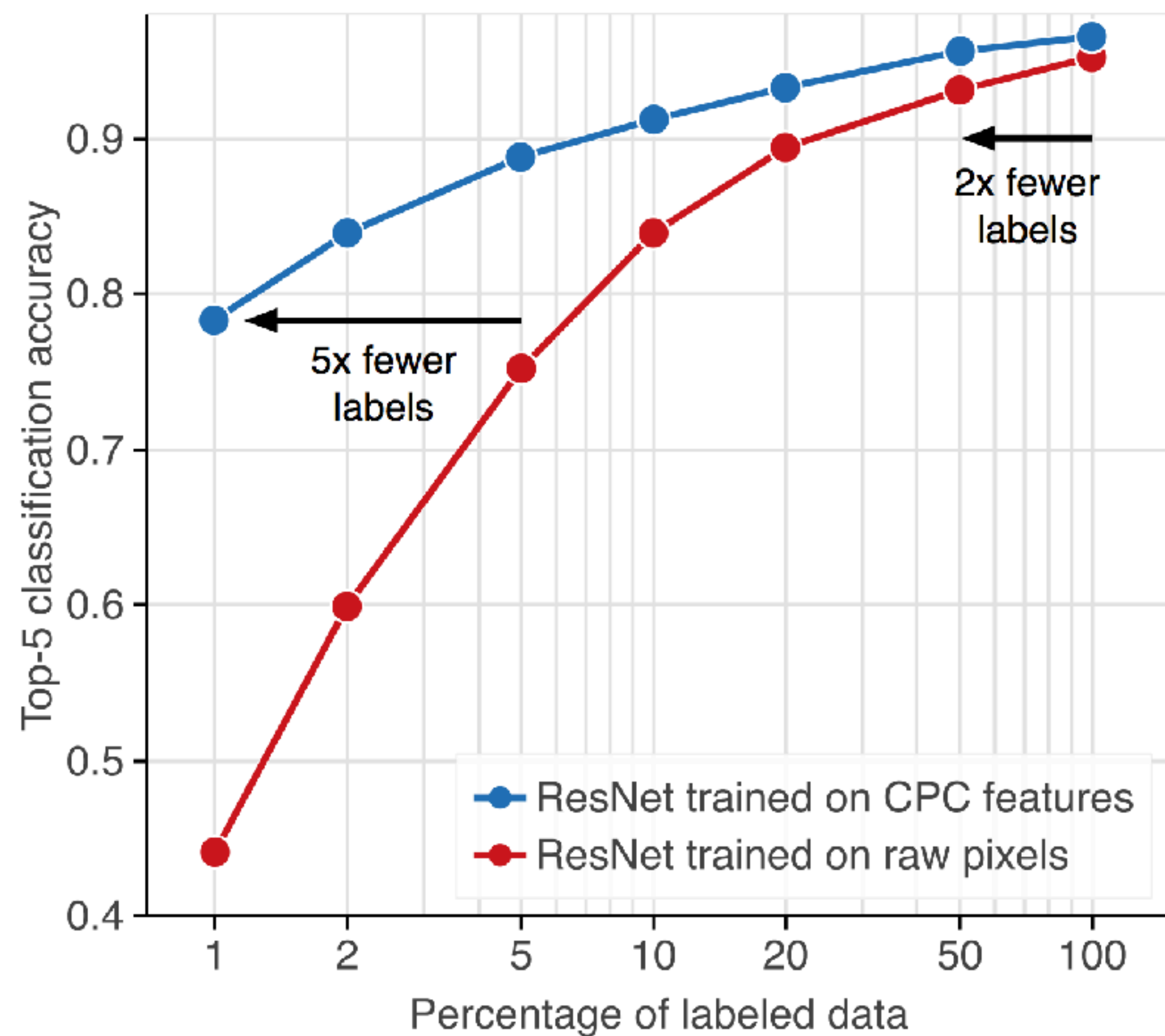
(From Alex Graves' tutorial)

Capture useful aspects of high dimensional sensory stream

# Capture useful aspects of high dimensional sensory stream



# Capture useful aspects of high dimensional sensory stream



# Quick background

1. Autoencoder
2. Variational Autoencoder
3. Contrastive Learning
4. Siamese Networks
5. Data-Augmentations

# Quick background

1. Autoencoder

2. Variational Autoencoder

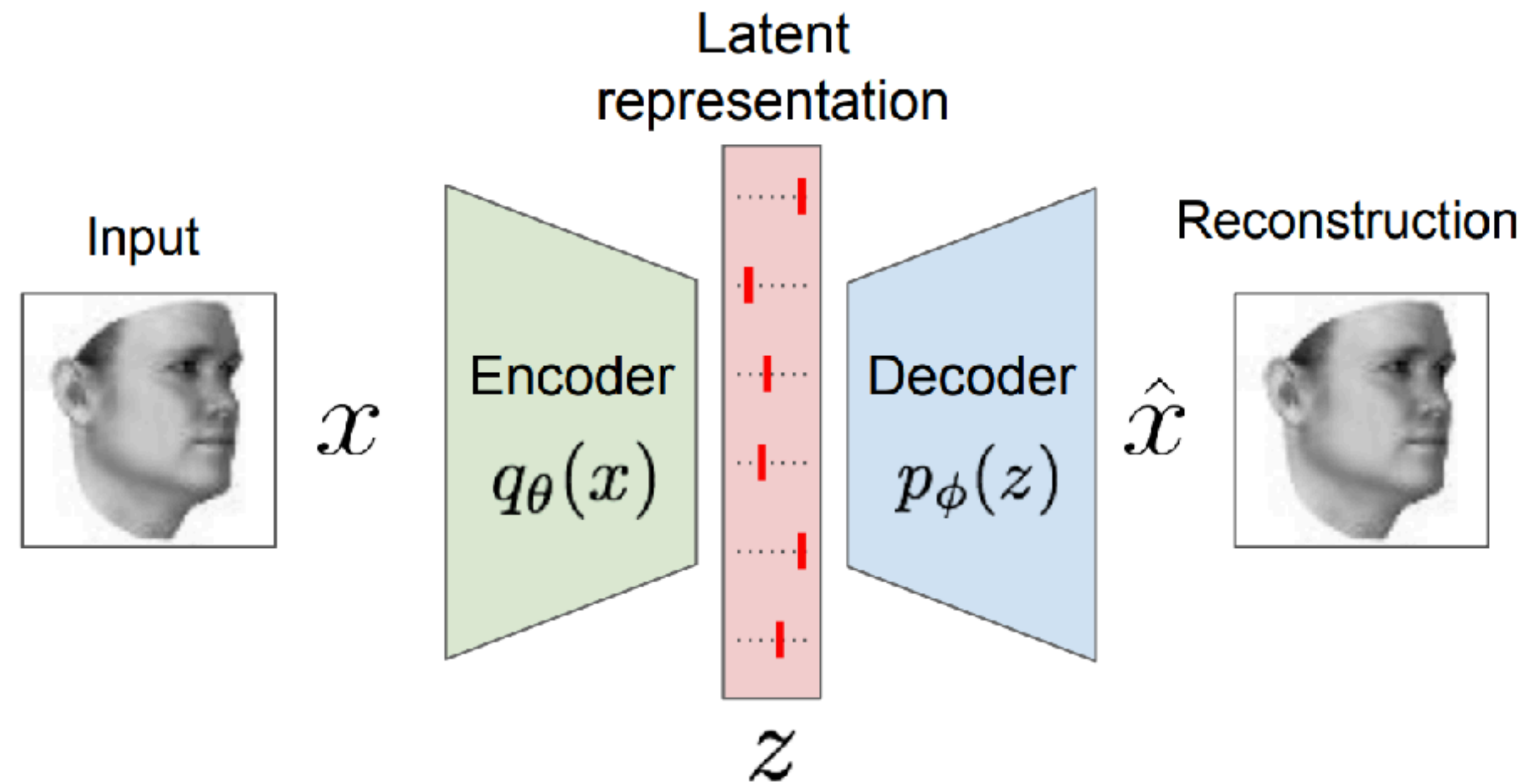
3. Contrastive Learning

4. Siamese Networks

5. Data-Augmentations

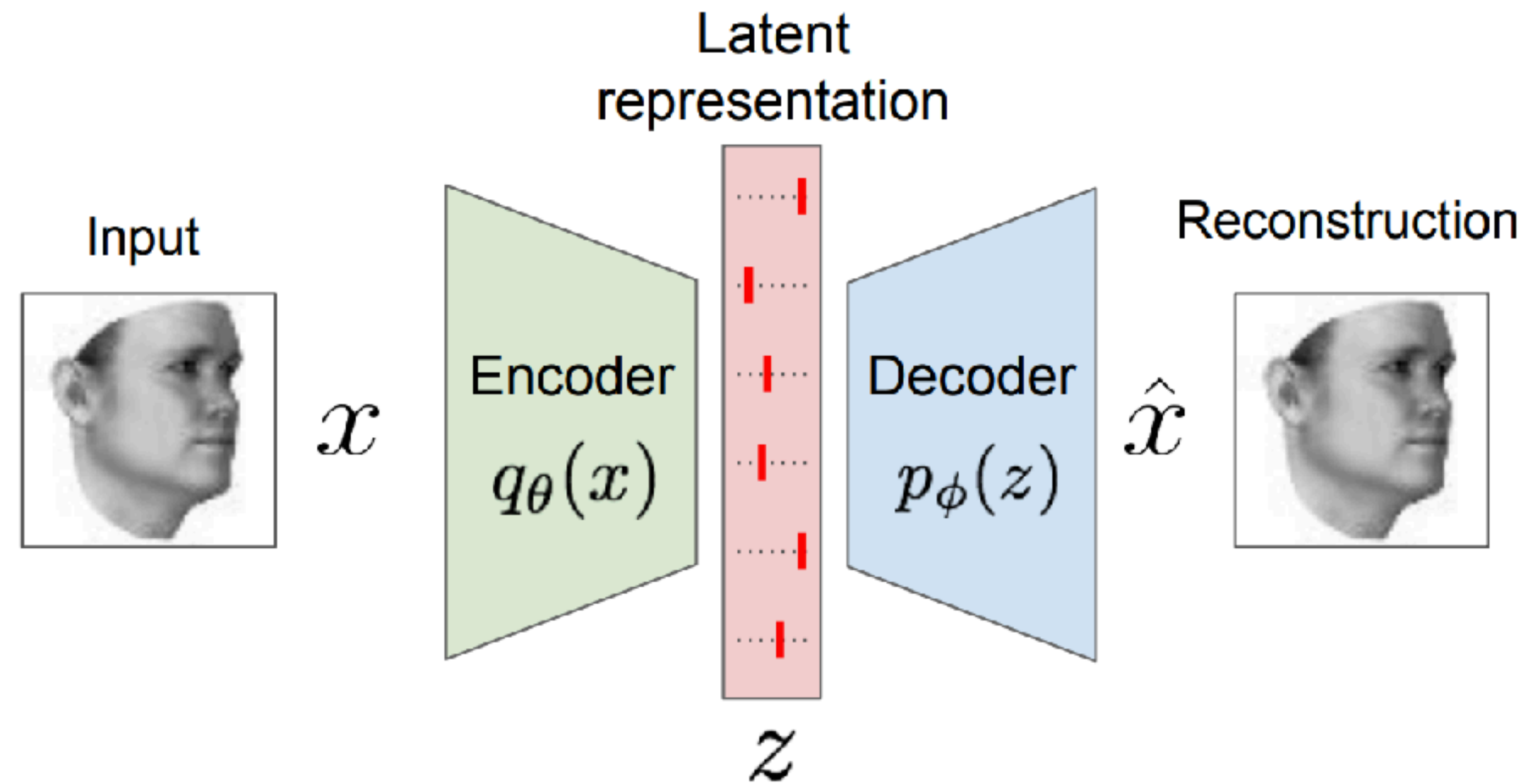
**Generative UL**

# AutoEncoder



$$\mathcal{L}^{AE}(\mathbf{x}; \theta, \phi) = \frac{[\mathbf{x} - p_{\theta}(q_{\phi}(\mathbf{x}))]^2}{\text{Reconstruction cost}}$$

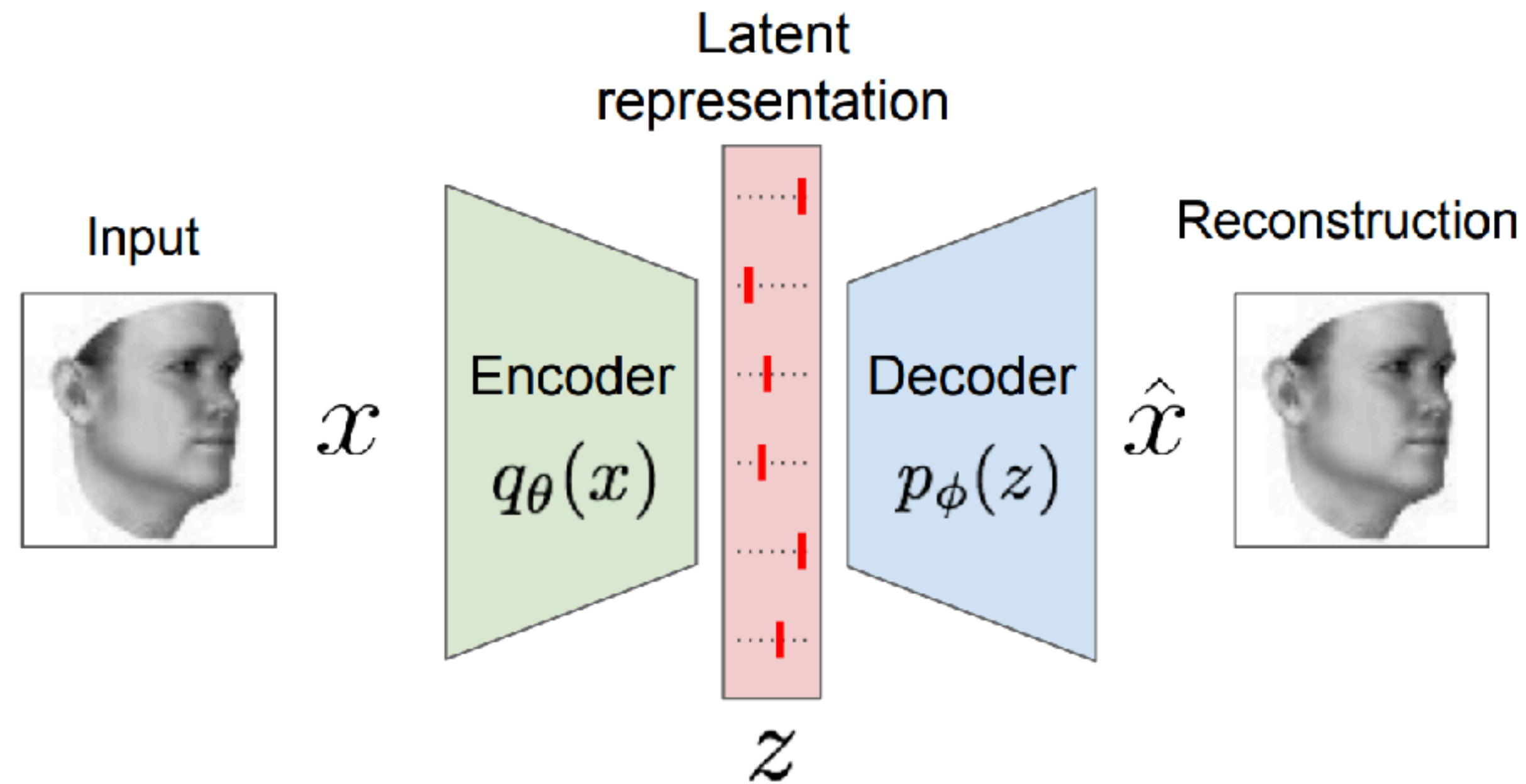
# AutoEncoder



$$\mathcal{L}^{AE}(\mathbf{x}; \theta, \phi) = \underbrace{\left[ \mathbf{x} - p_{\theta}(q_{\phi}(\mathbf{x})) \right]^2}_{\text{Reconstruction cost}} - \log p_{\theta}(q_{\phi}(\mathbf{x}))$$



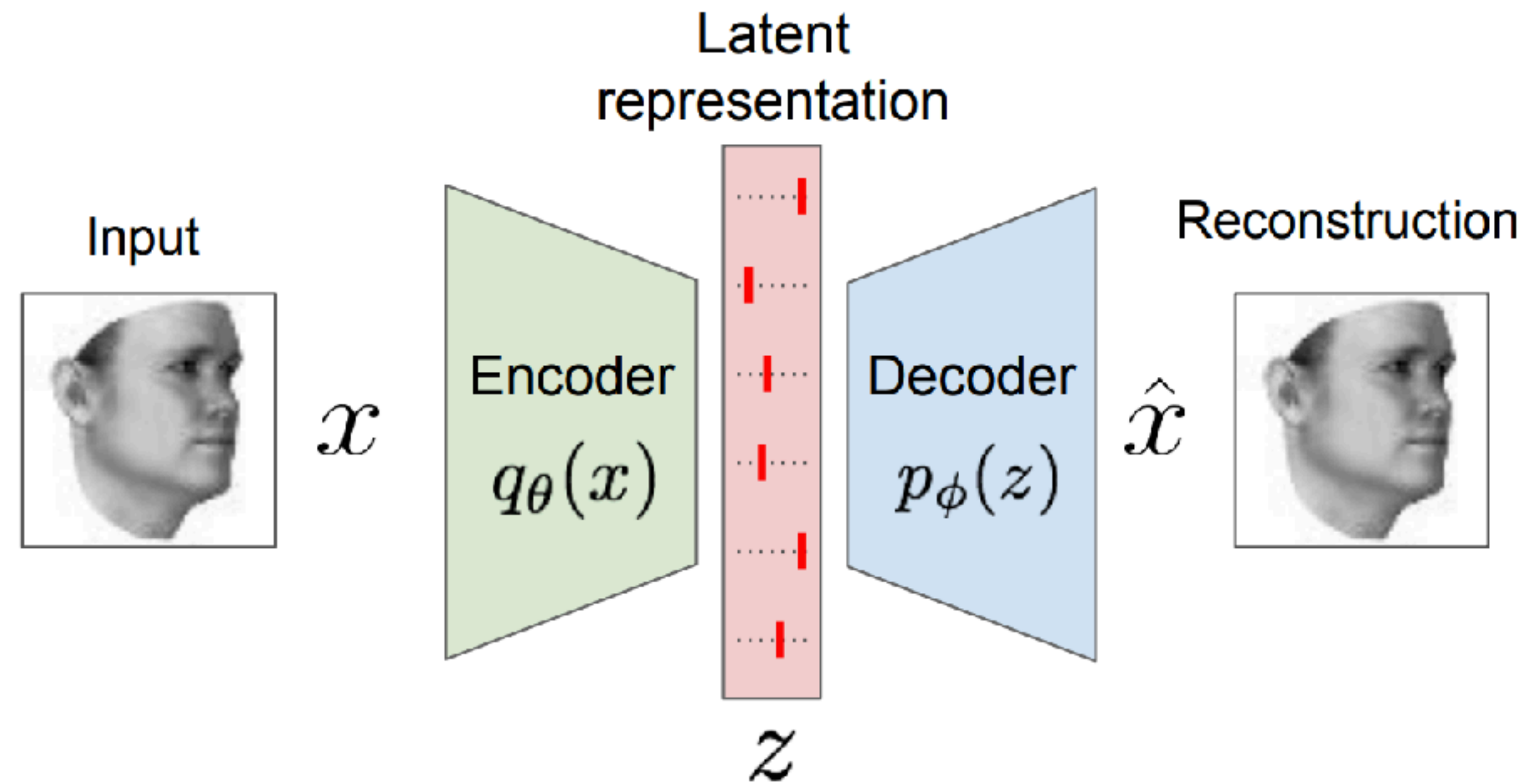
# AutoEncoder



$$\mathcal{L}^{AE}(\mathbf{x}; \theta, \phi) = \underbrace{\left[ \mathbf{x} - p_{\theta}(q_{\phi}(\mathbf{x})) \right]^2}_{\text{Reconstruction cost}} - \log p_{\theta}(q_{\phi}(\mathbf{x}))$$

A lossy (bottleneck) representation of the input - but packs as much information as possible

# AutoEncoder



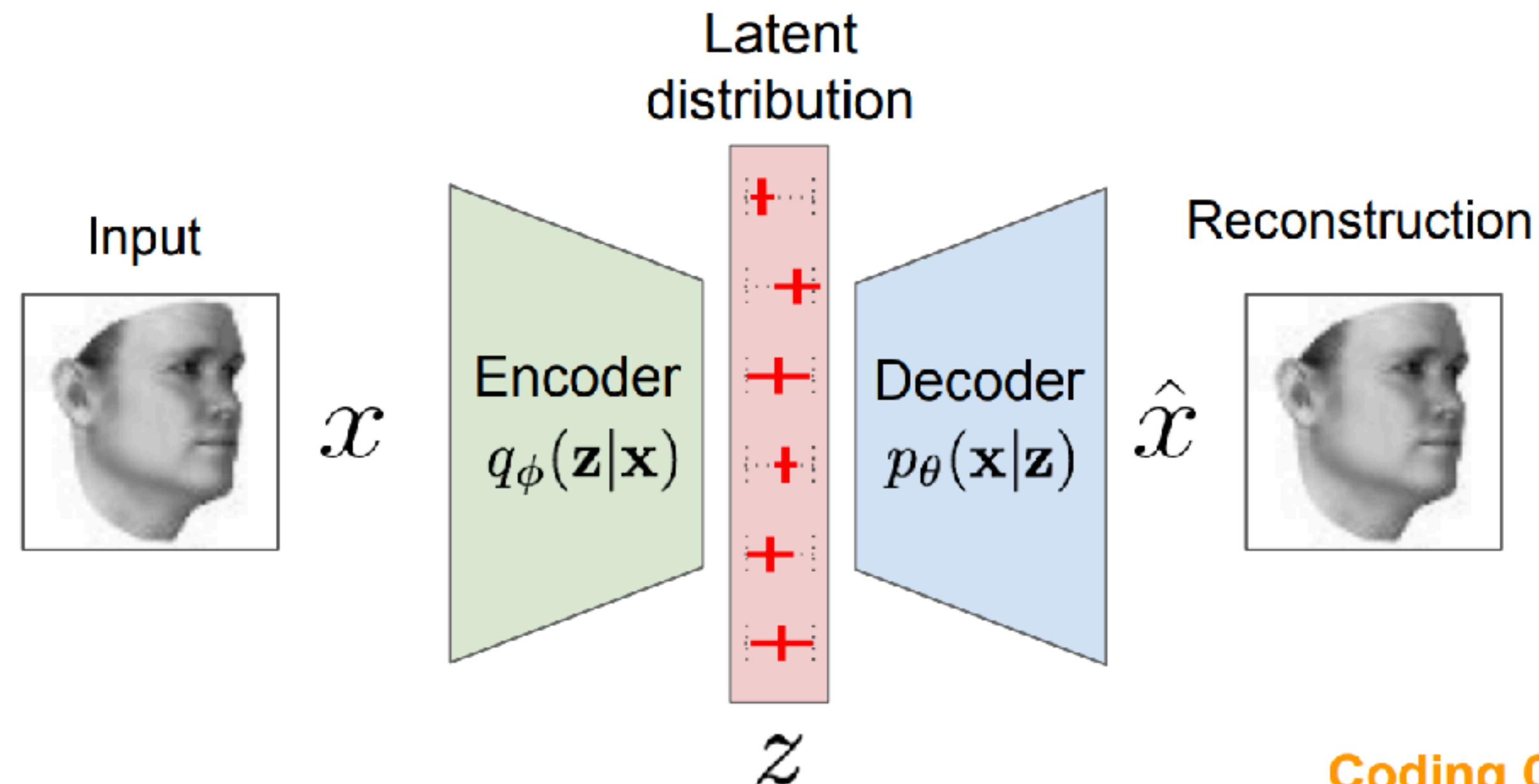
$$\mathcal{L}^{AE}(\mathbf{x}; \theta, \phi) = \underbrace{\left[ \mathbf{x} - p_{\theta}(q_{\phi}(\mathbf{x})) \right]^2}_{\text{Reconstruction cost}} - \log p_{\theta}(q_{\phi}(\mathbf{x}))$$

A lossy (bottleneck) representation of the input - but packs as much information as possible

Adding some structure to the latent (bottleneck) can help add more semantic meaning

# Variational AutoEncoder (VAE)

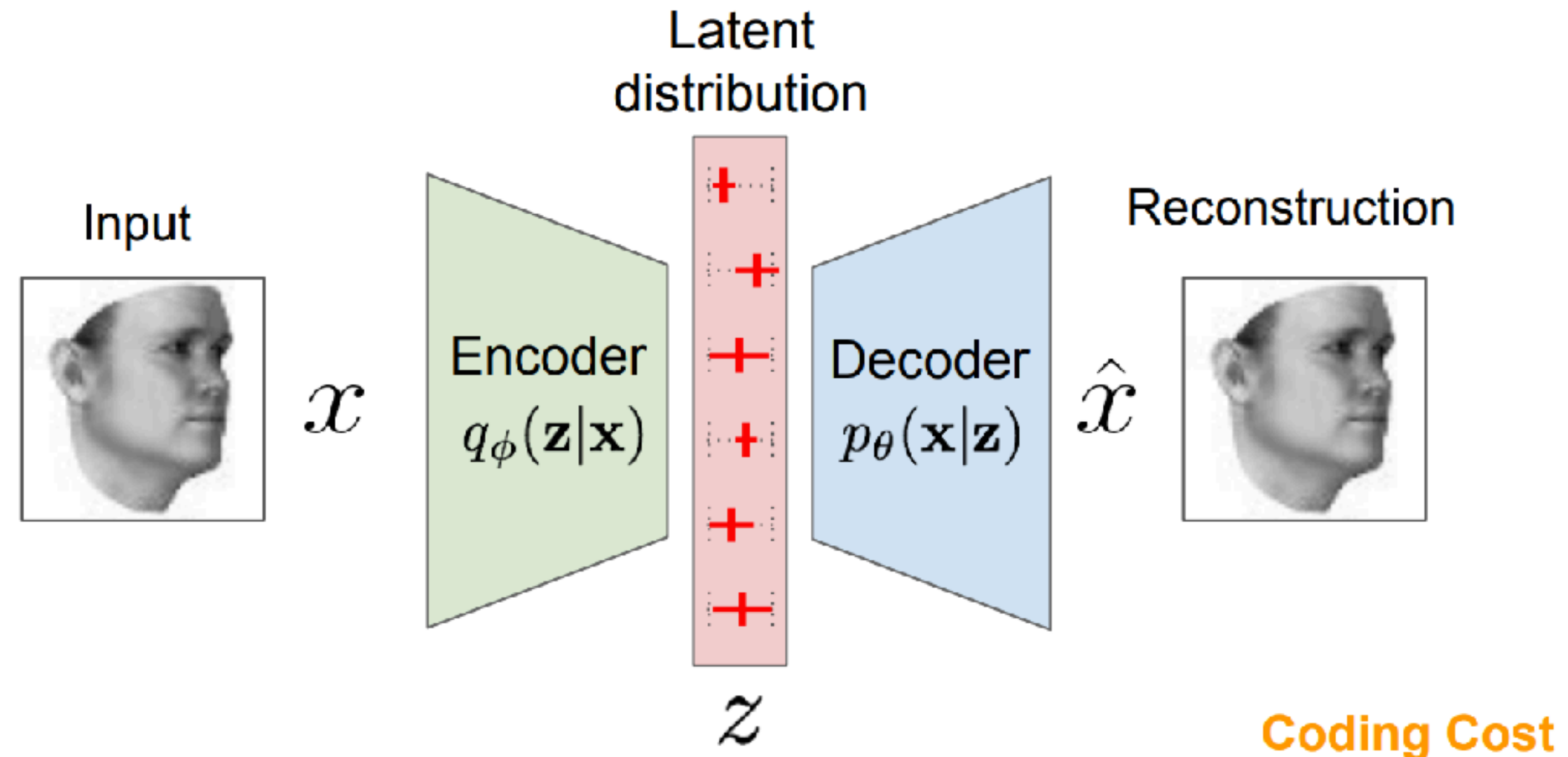
Kingma et al, 2014  
Rezende et al, 2014



$$\mathcal{L}_{VAE}(\mathbf{x}; \theta, \phi) = \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [-\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction cost}} + \underbrace{KL(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{Coding Cost}}$$

# Variational AutoEncoder (VAE)

Kingma et al, 2014  
Rezende et al, 2014

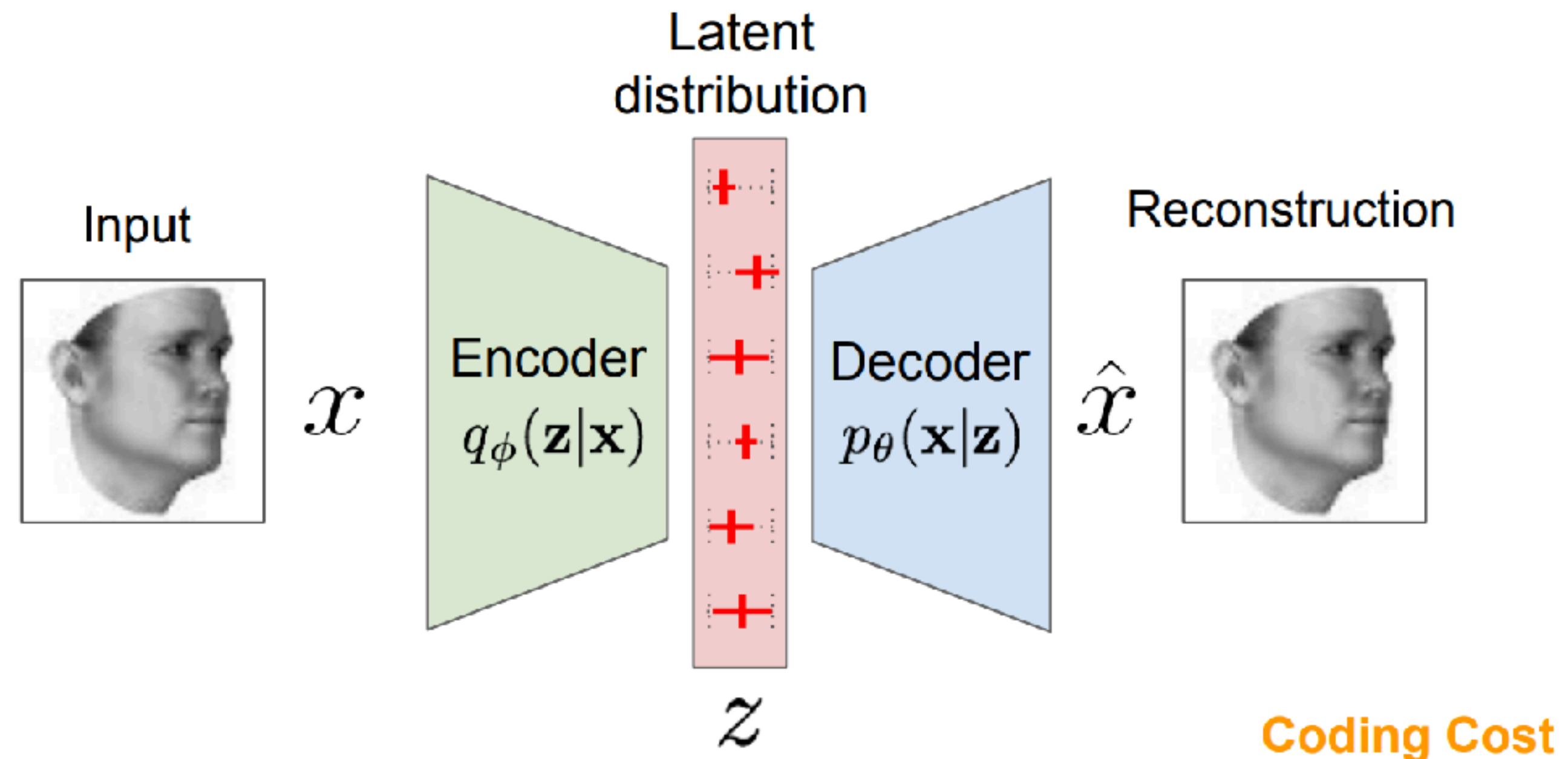


$$\mathcal{L}_{VAE}(\mathbf{x}; \theta, \phi) = \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [-\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction cost}} + \underbrace{KL(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{Coding Cost}}$$

Make sure there's an additional penalty for the latents (posterior) to match a prior

# Variational AutoEncoder (VAE)

Kingma et al, 2014  
Rezende et al, 2014

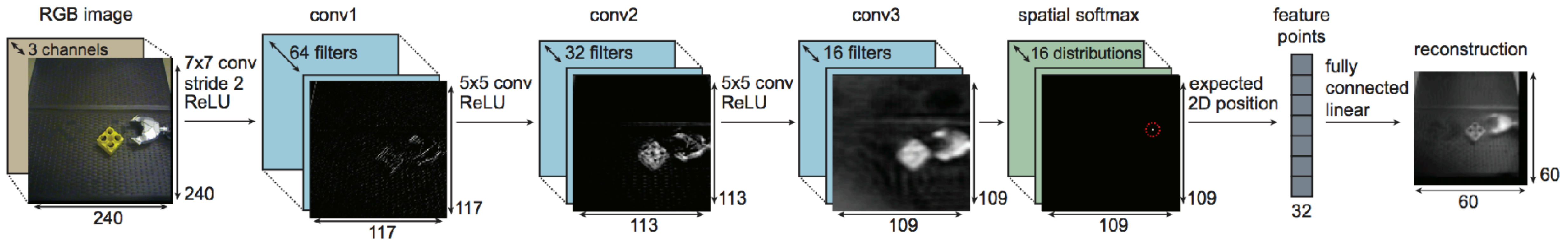


$$\mathcal{L}_{VAE}(\mathbf{x}; \theta, \phi) = \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [-\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction cost}} + \underbrace{KL(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{Coding Cost}}$$

Make sure there's an additional penalty for the latents (posterior) to match a prior

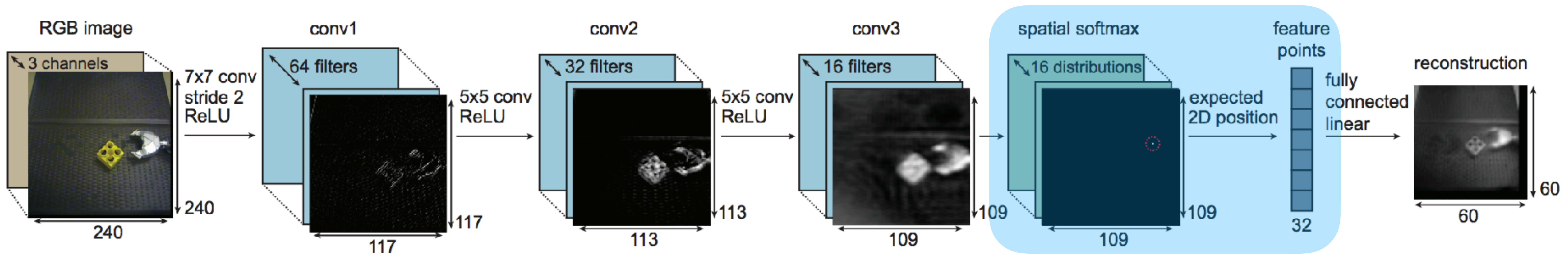
reconstruction pathway goes through a cost for how much information it can pack in; also has to match the prior

# Spatial AutoEncoder



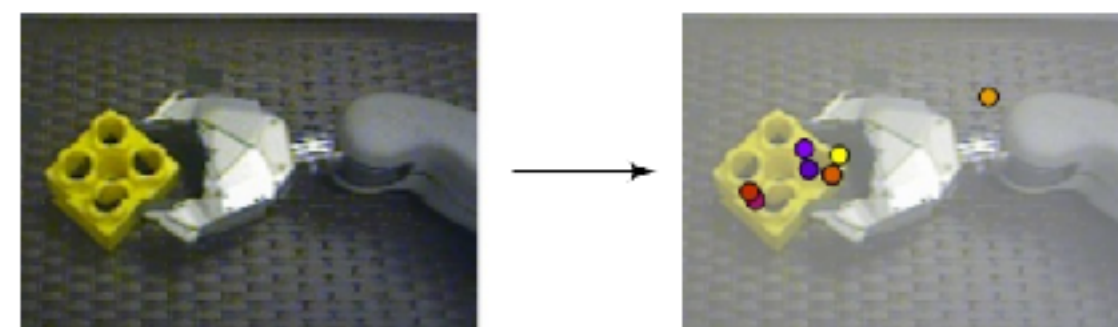
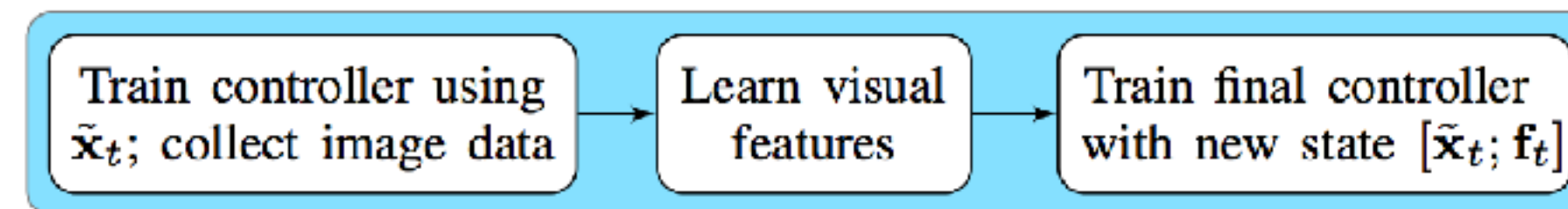
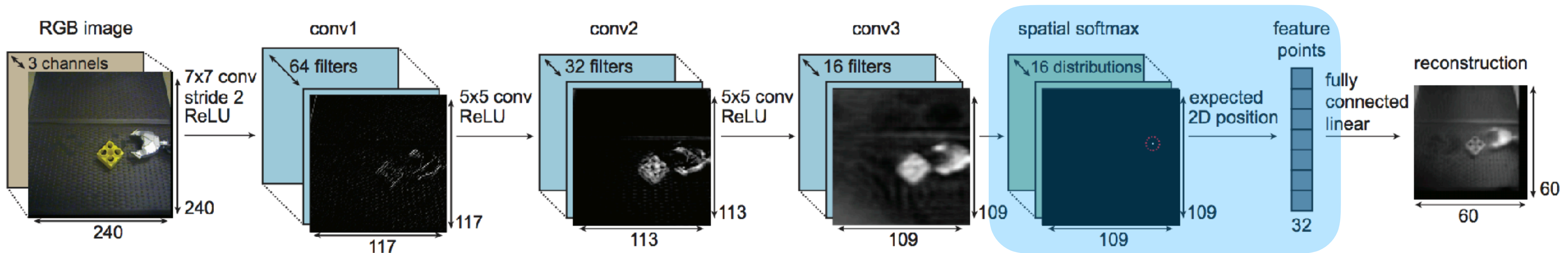
Deep Spatial Autoencoders for Visuomotor Learning (Finn et al 2015)

# Spatial AutoEncoder



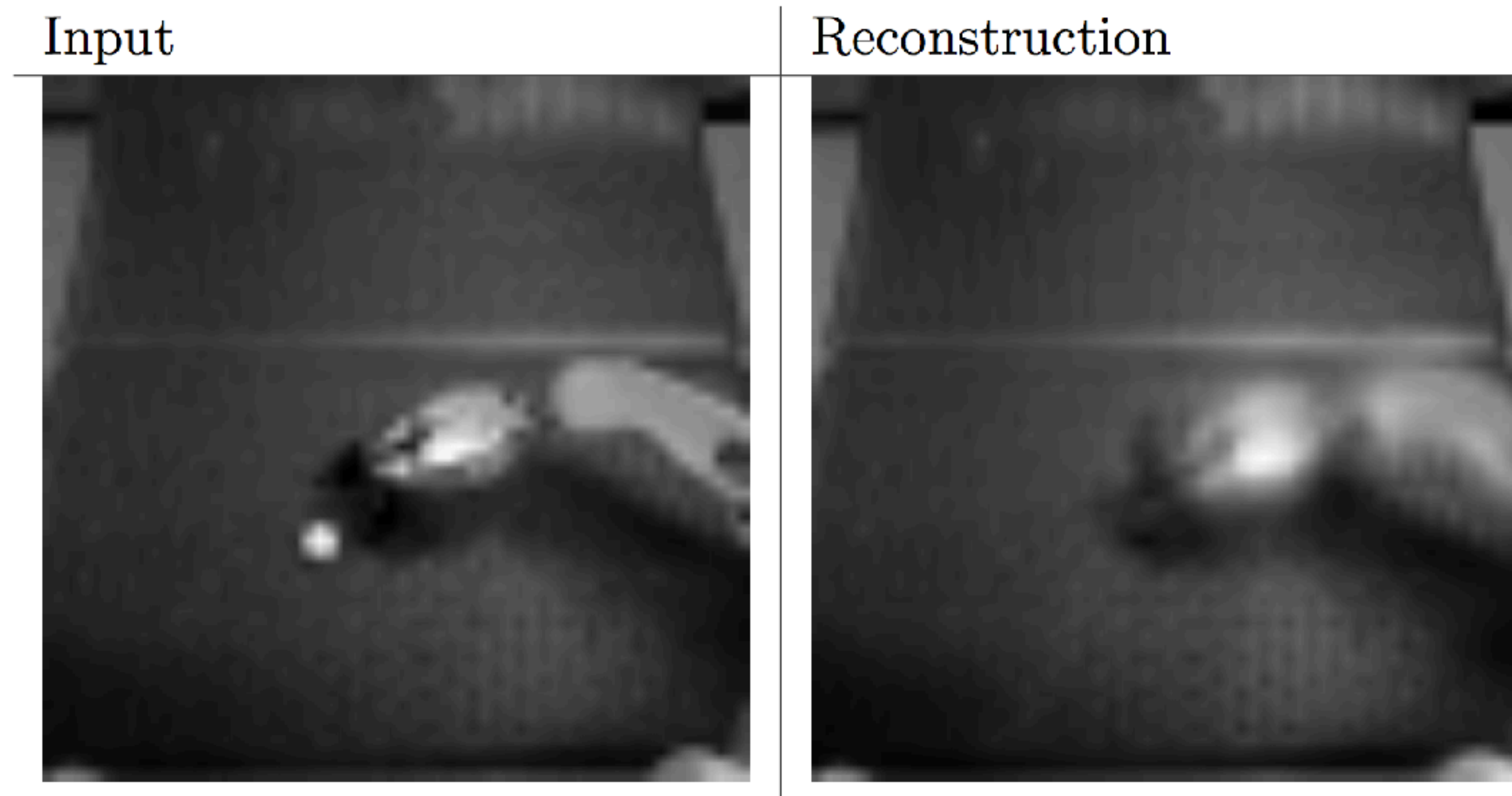
Deep Spatial Autoencoders for Visuomotor Learning (Finn et al 2015)

# Spatial AutoEncoder





# AutoEncoders can ignore relevant features for a task



Reconstruction error isn't significantly altered by presence or absence of the ping-pong ball.

# AutoEncoders can ignore relevant features for a task



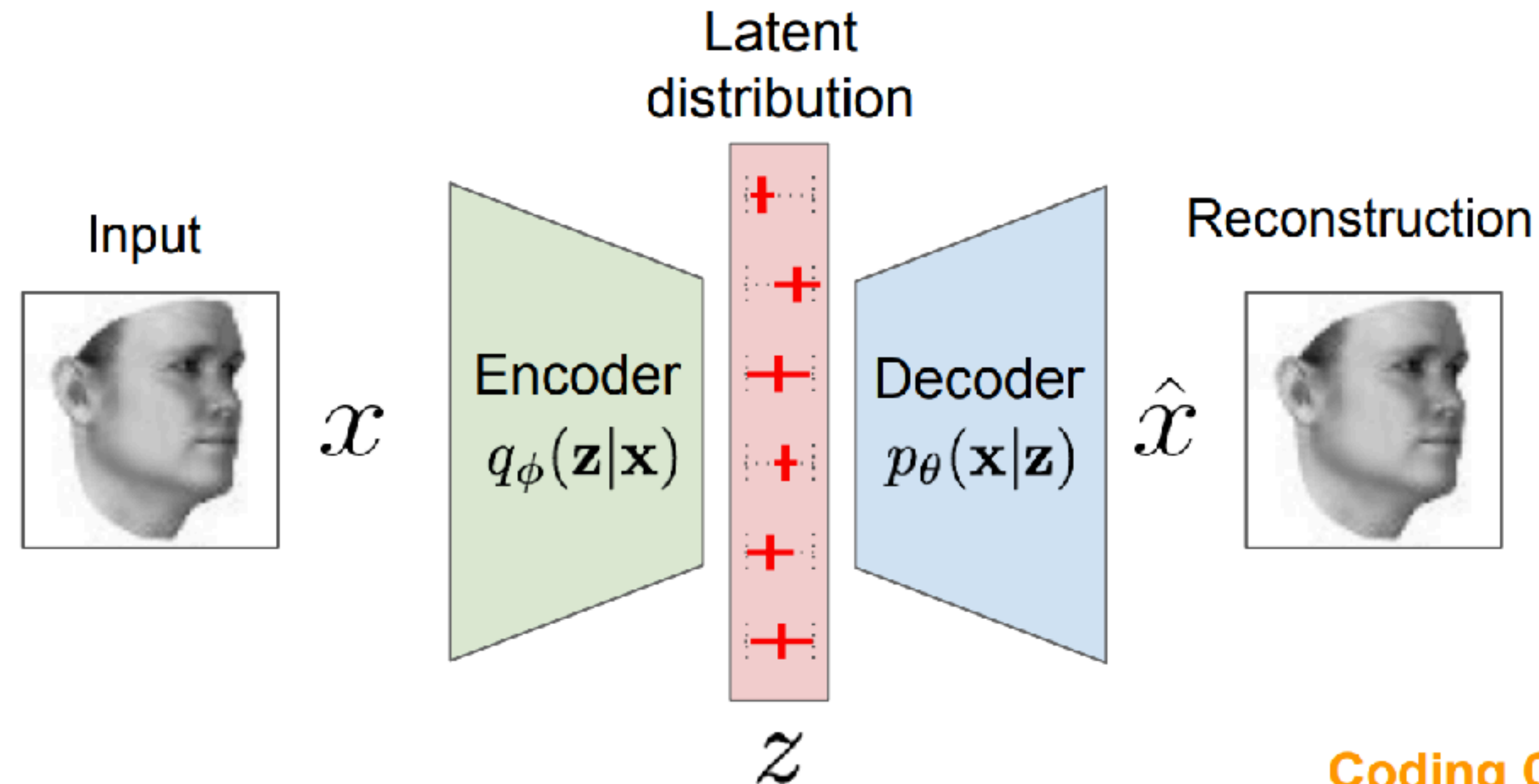
# AutoEncoders can ignore relevant features for a task



If the task is to accurately identify the pastry place, some of the details are blurry. However, future incarnations could potentially address these issues. It all depends on how much you downsample (8x in this case).

# Beta-Variational AutoEncoder (beta-VAE)

Kingma et al, 2014  
Rezende et al, 2014

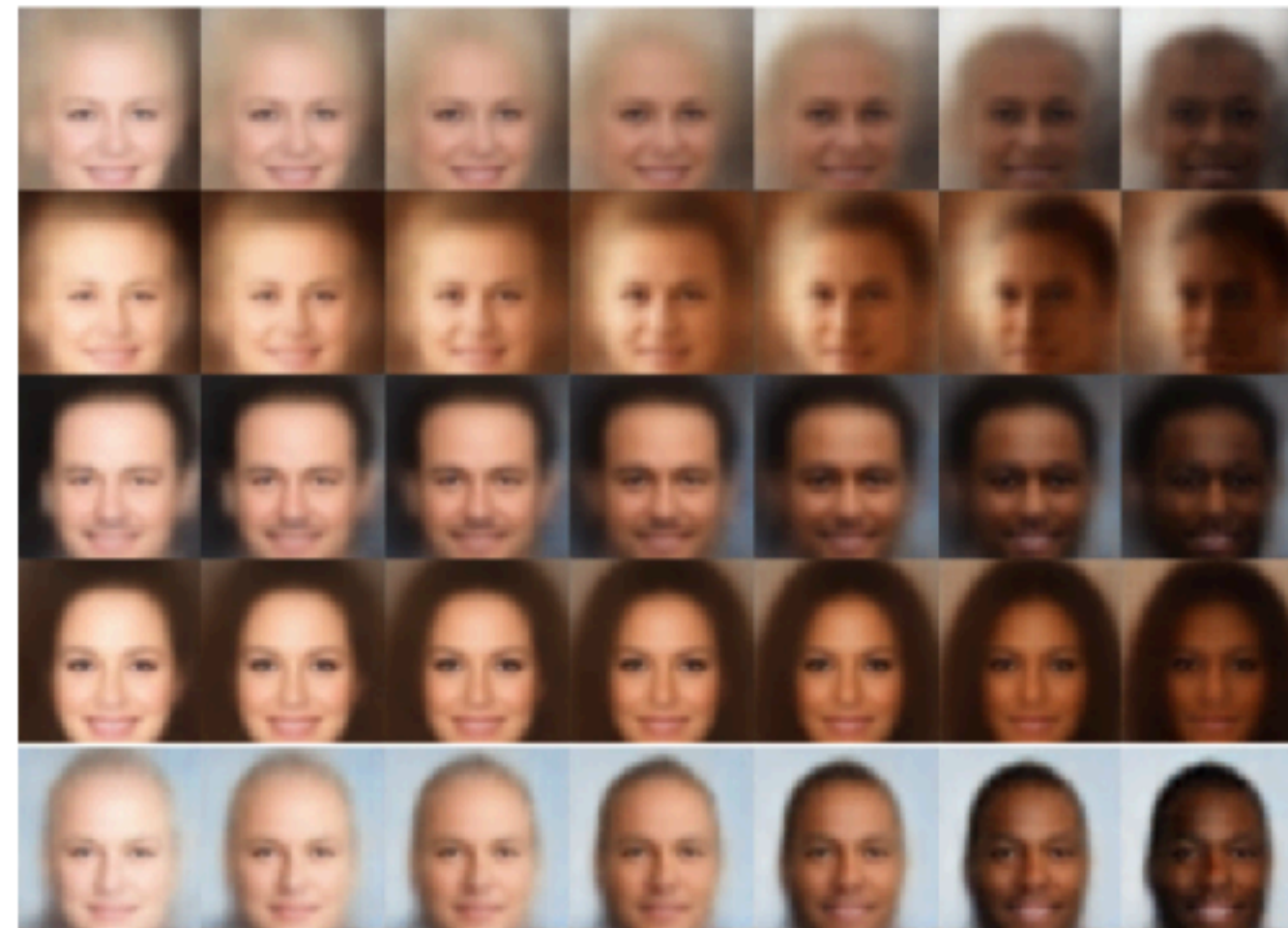


$$\mathcal{L}_{VAE}(\mathbf{x}; \theta, \phi) = \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [-\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction cost}} + \underbrace{KL(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))}_{\text{Coding Cost}} \cdot \beta$$

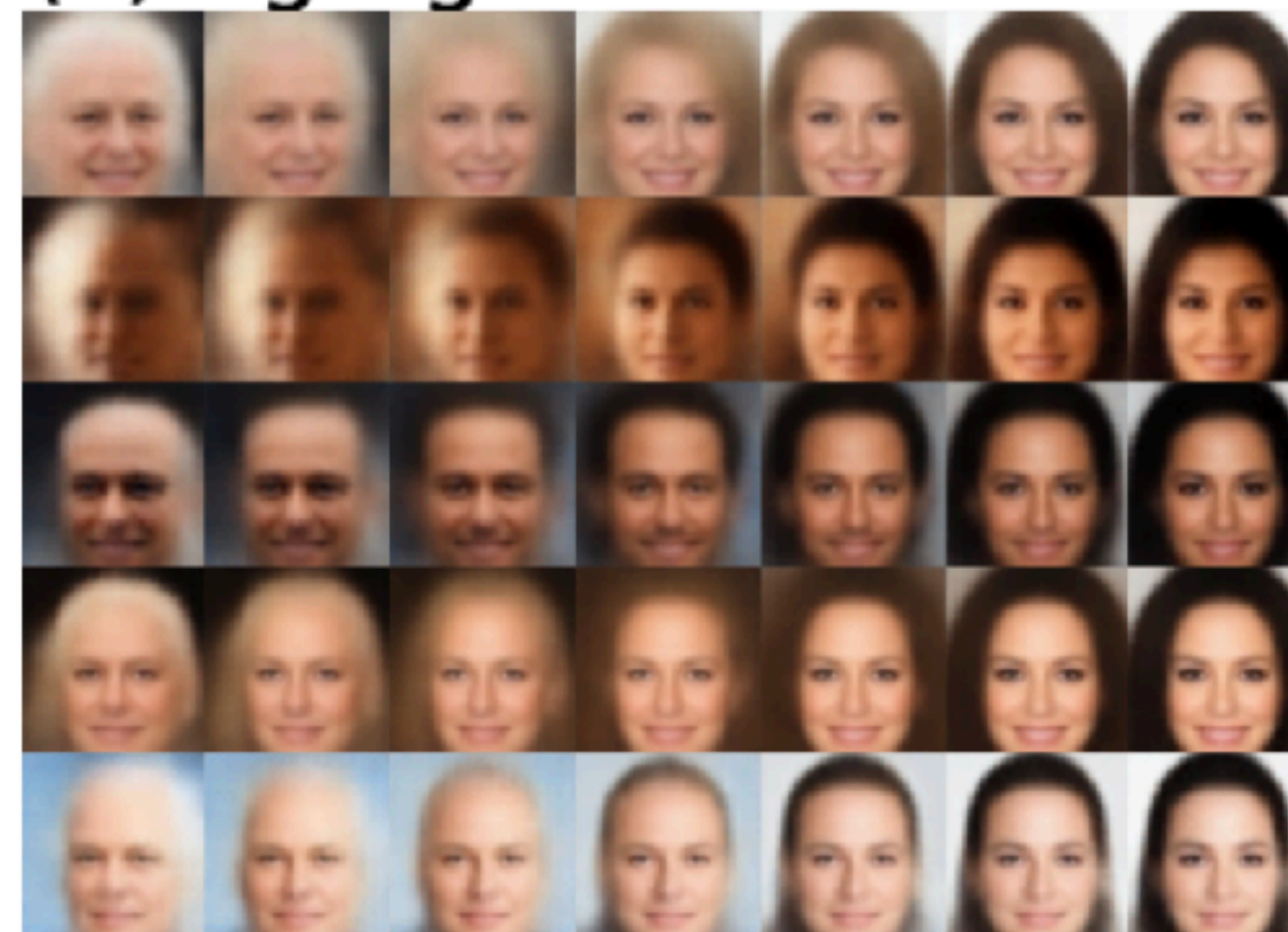
Increase the KL cost on the latent. Leads to more disentanglement

# Beta-Variational AutoEncoder (beta-VAE)

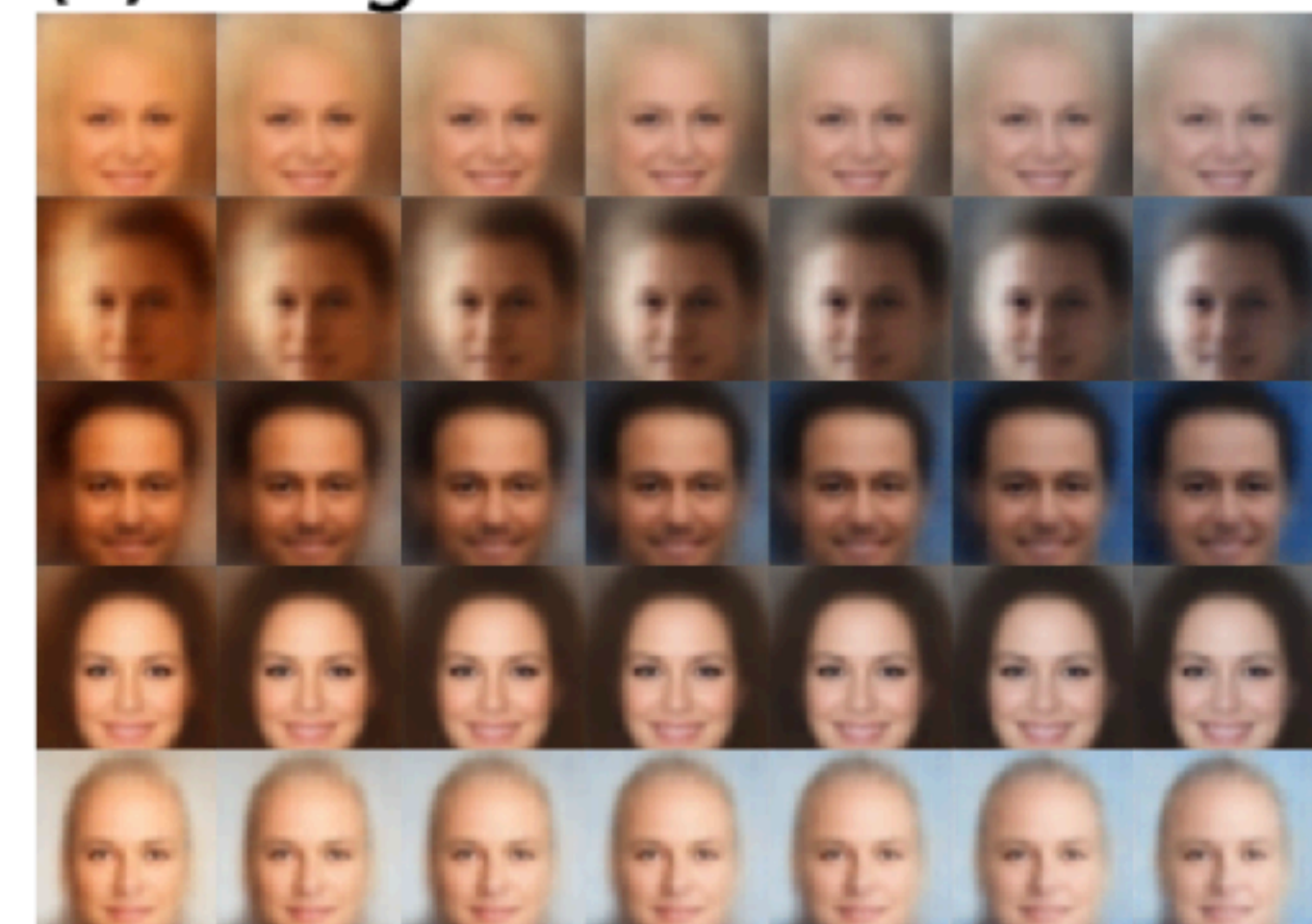
(a) Skin colour



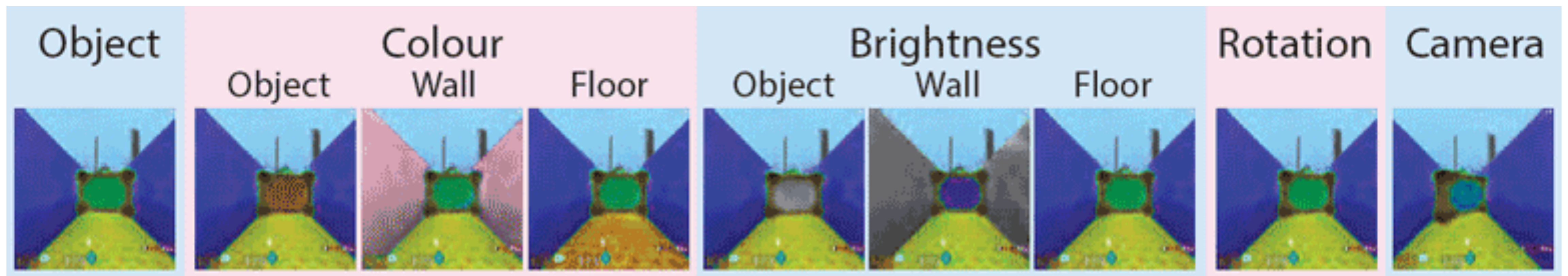
(b) Age/gender



(c) Image saturation

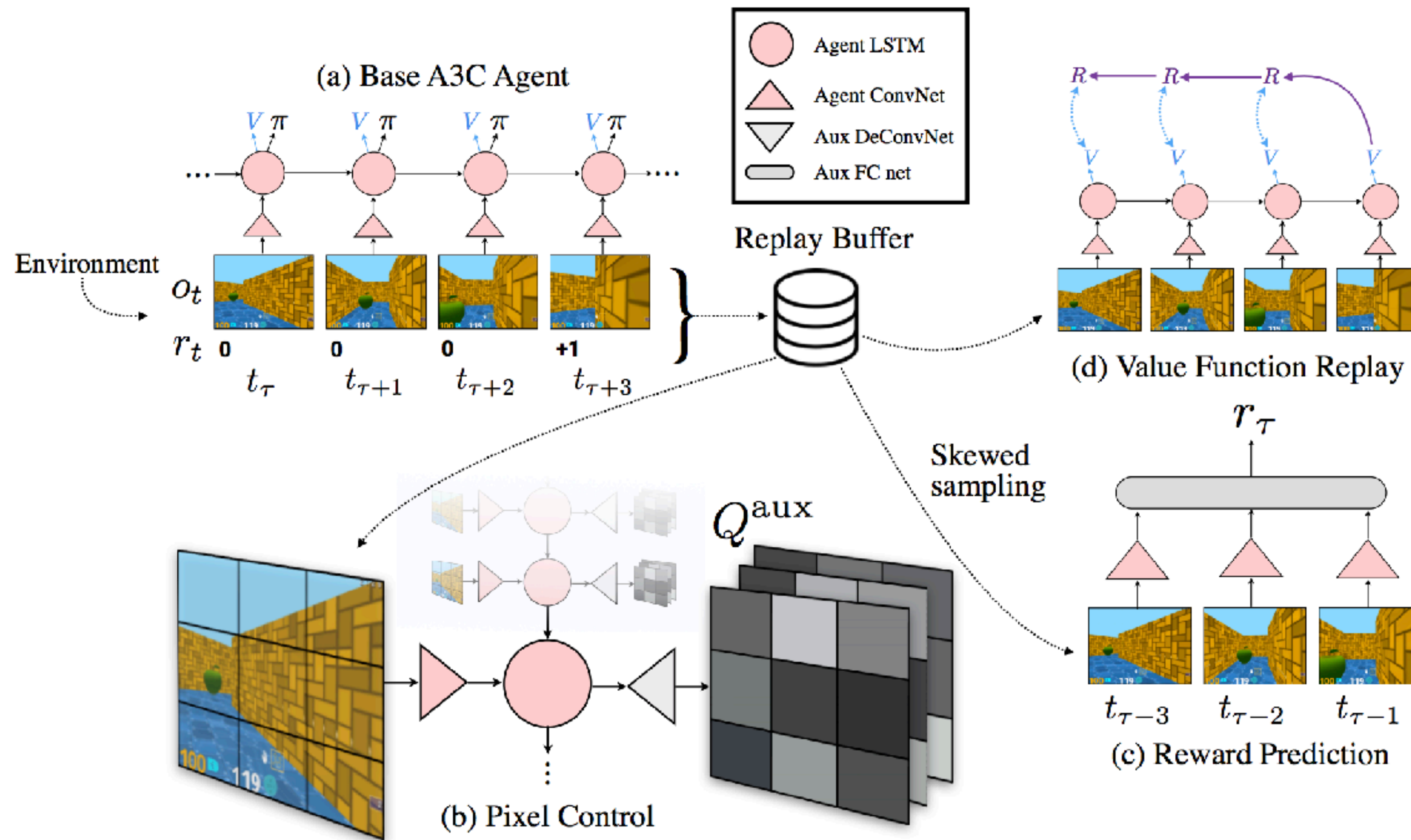


# Beta-Variational AutoEncoder (beta-VAE)

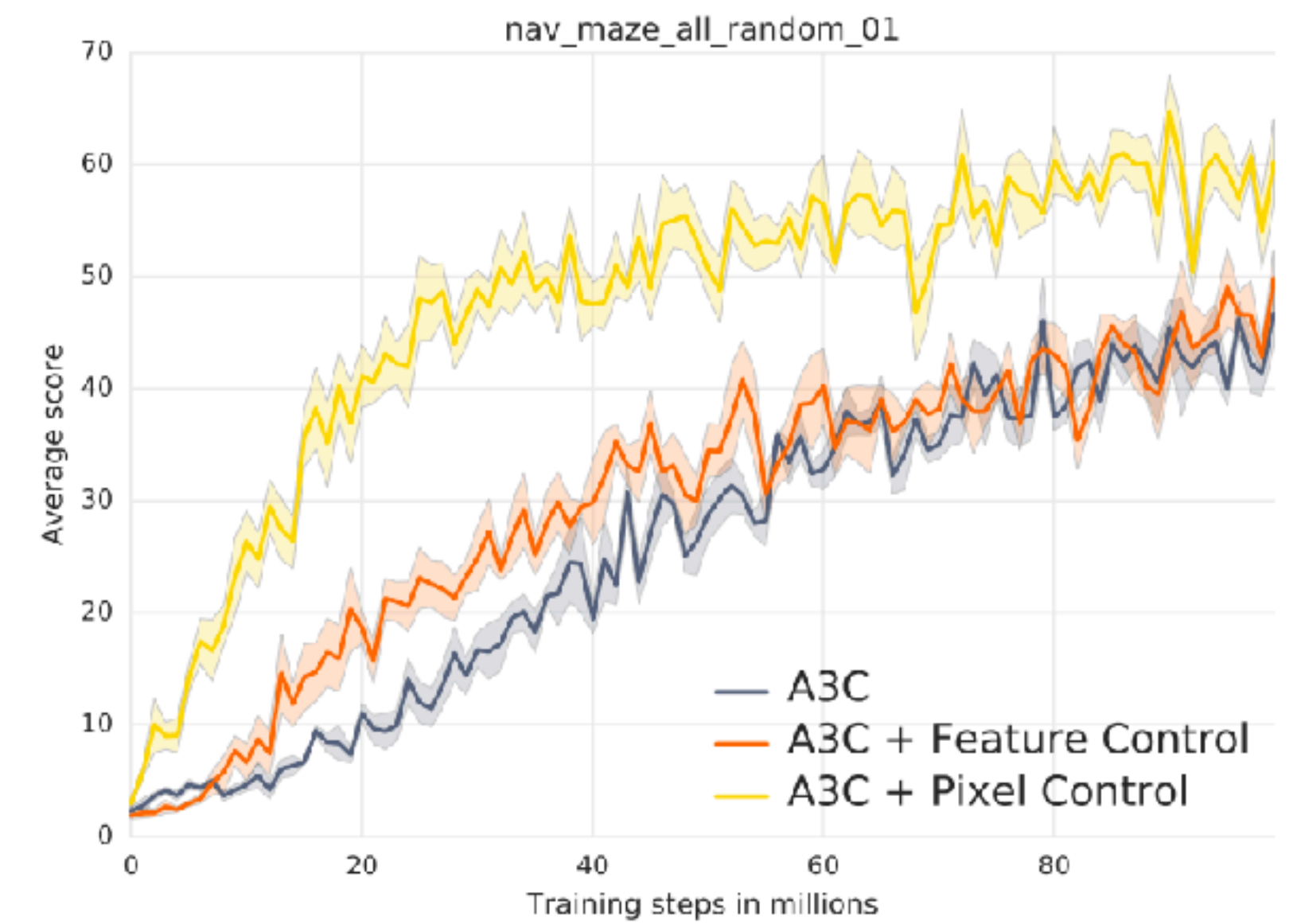
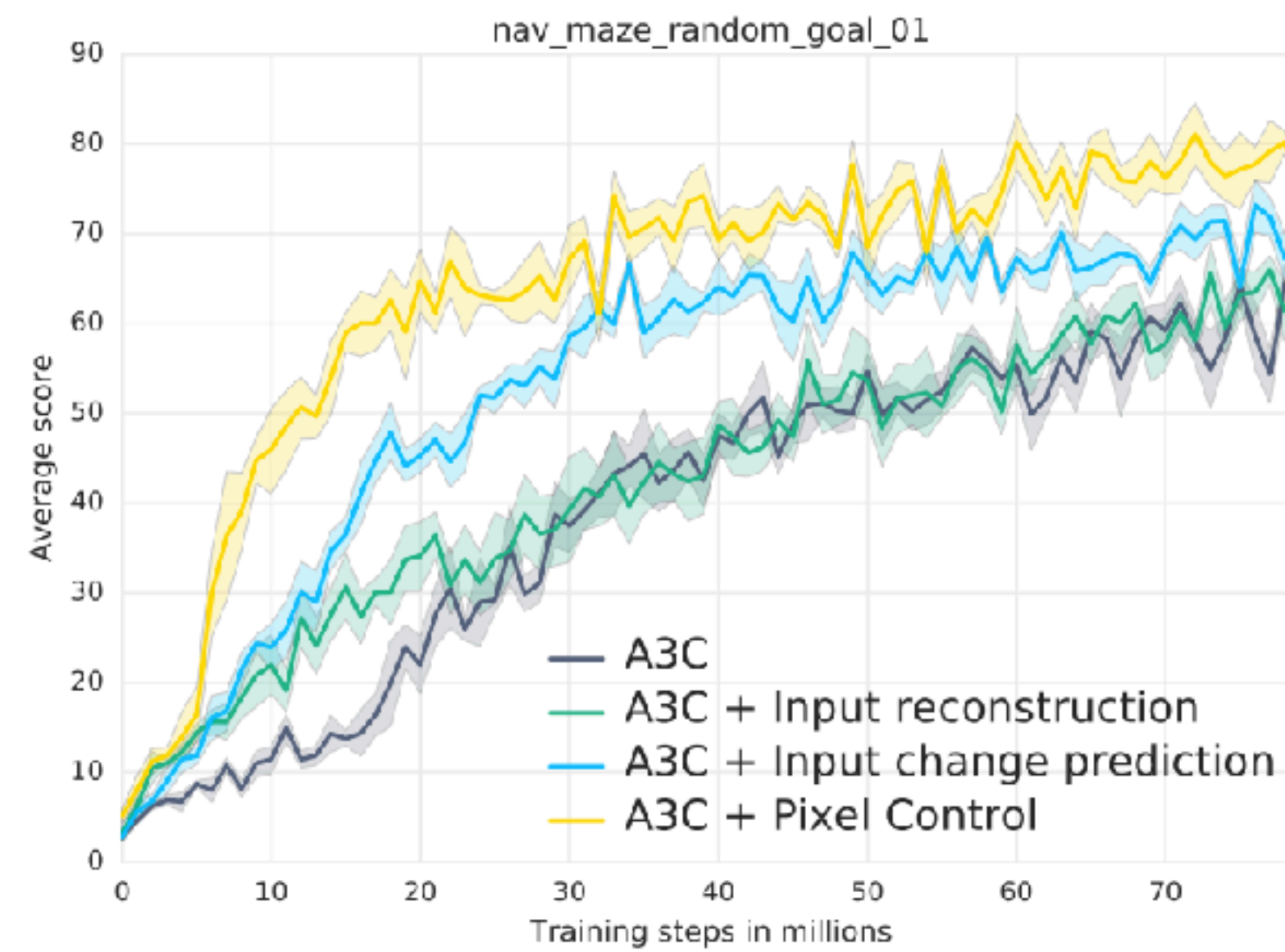
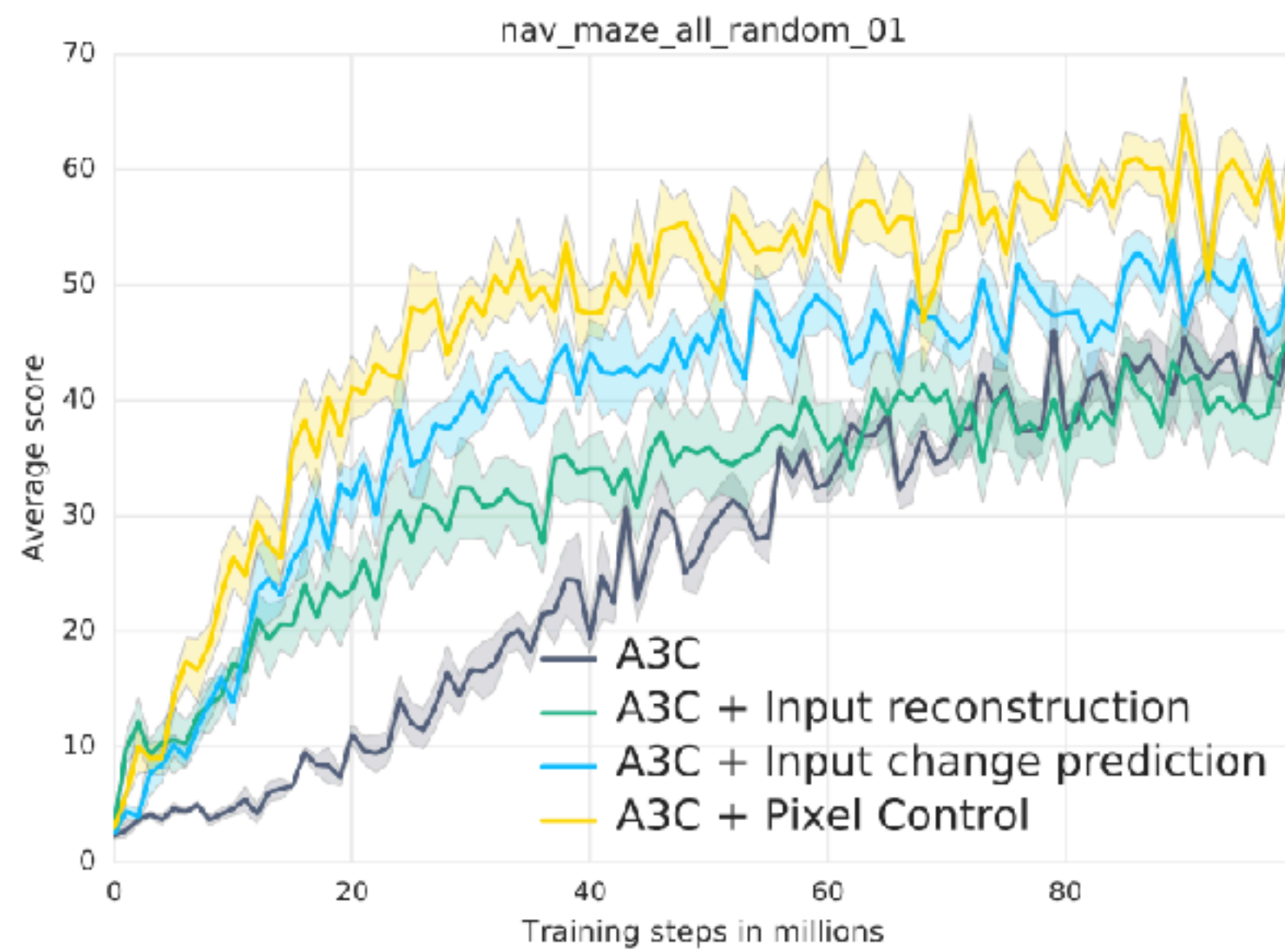


SCAN: Learning Abstract Hierarchical Compositional Visual Concepts (DeepMind Blog)

# UNREAL Architecture



# UNREAL Architecture





# DARLA: Improving Zero-Shot Transfer in RL

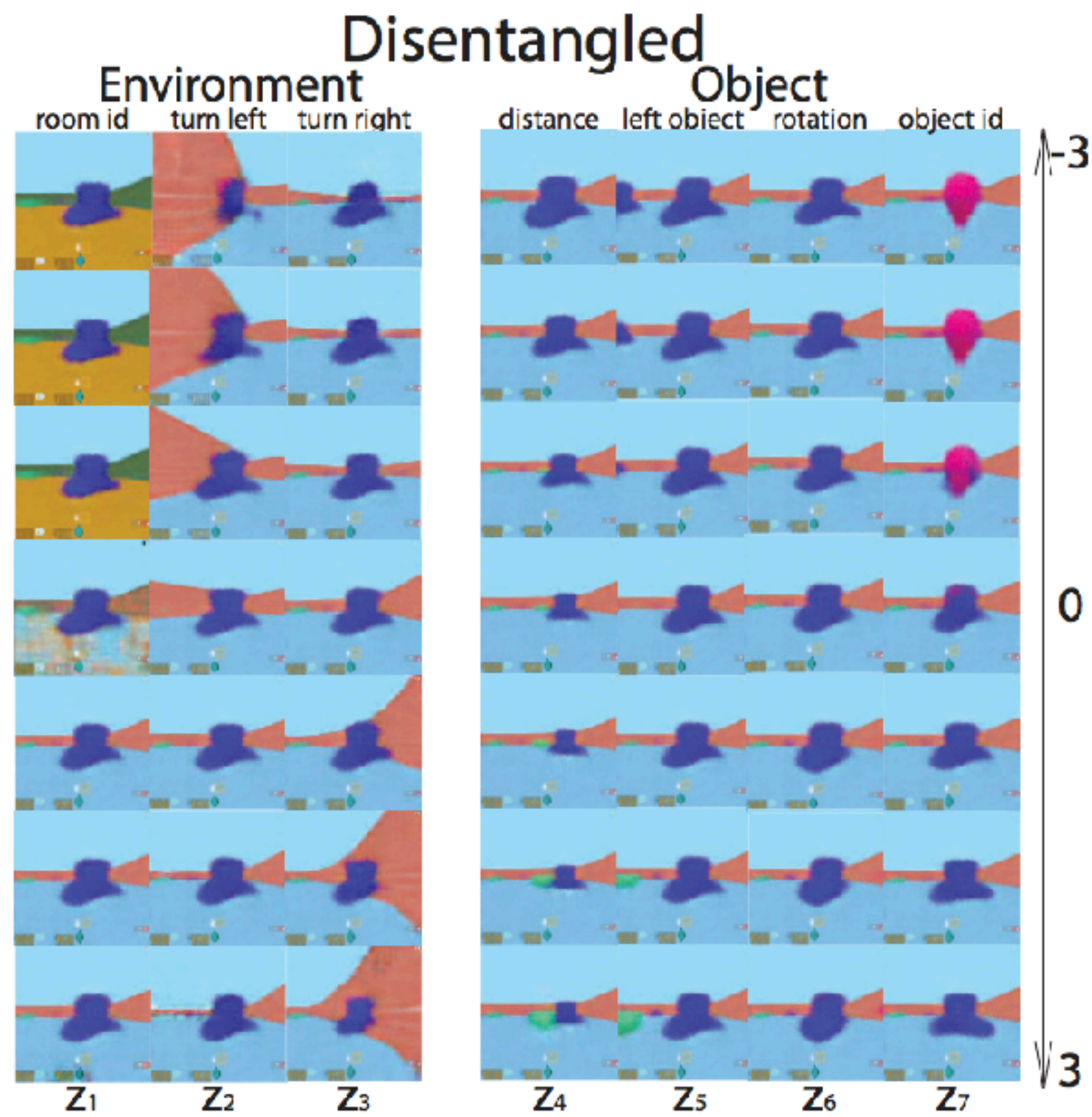
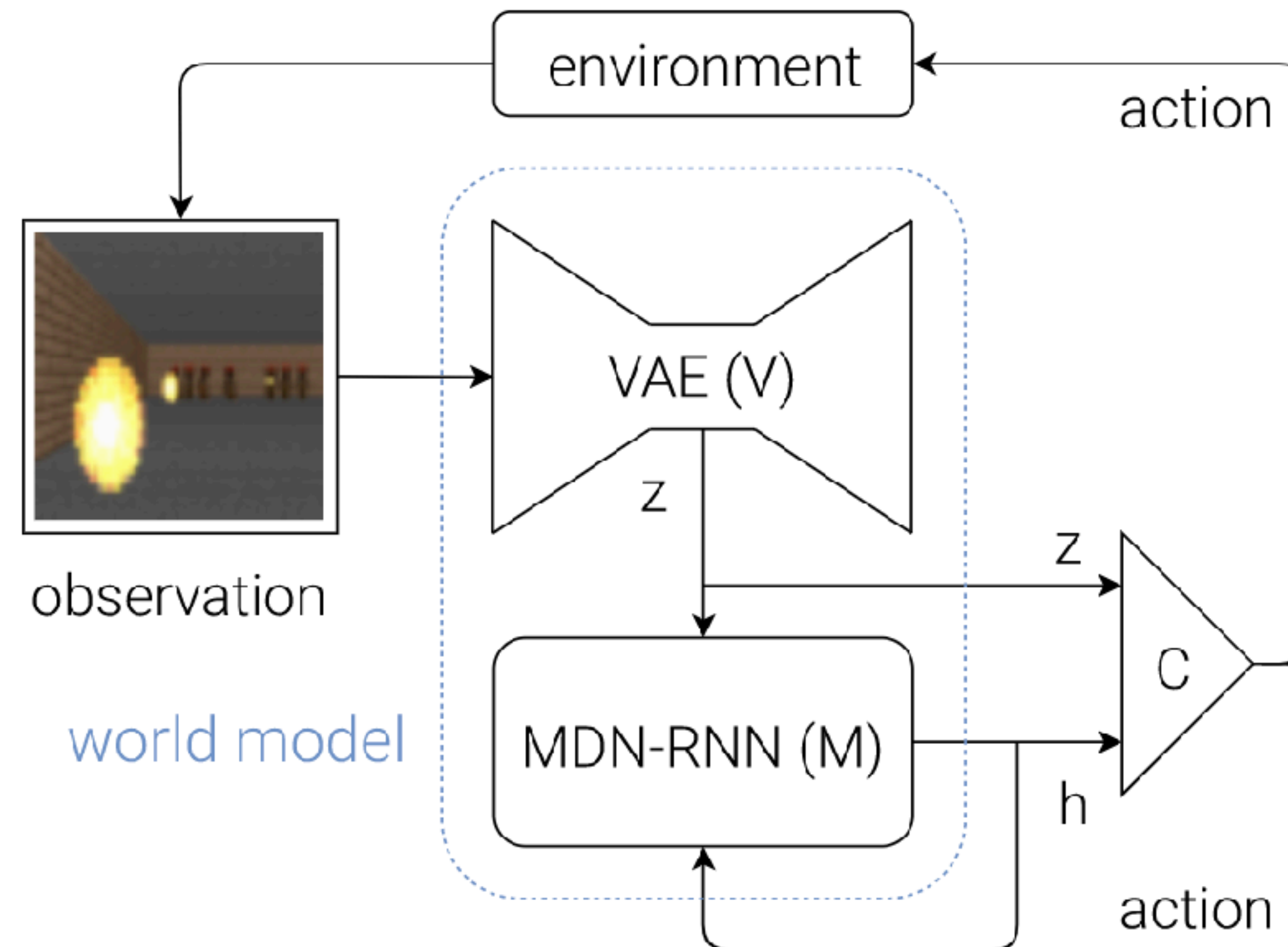


Table 1. Transfer performance

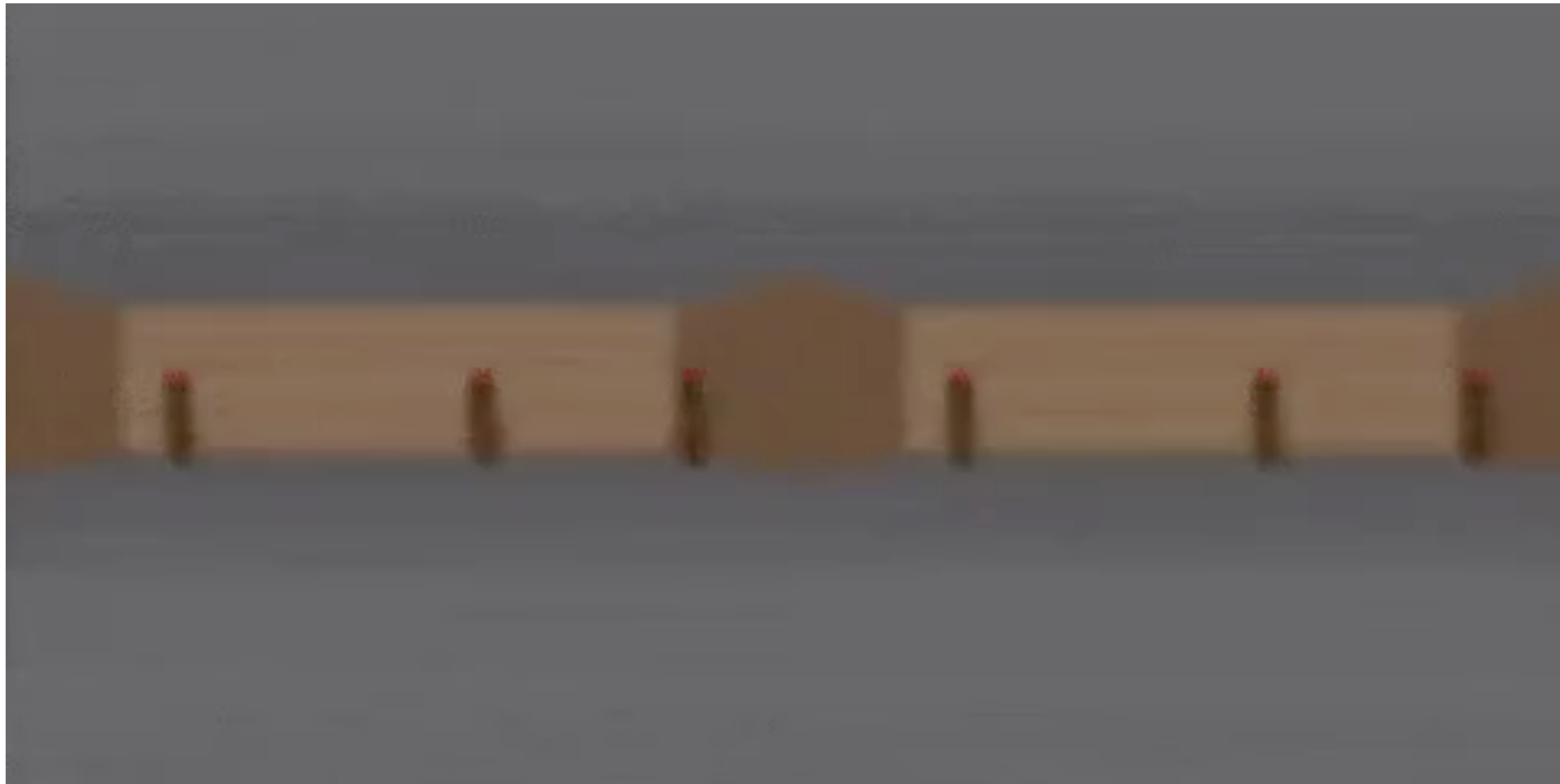
VISION TYPE	DQN	DEEPMIND LAB		JACO (A3C)	
		A3C	EC	SIM2SIM	SIM2REAL
BASELINE AGENT	$1.86 \pm 3.91$	$5.32 \pm 3.36$	$-0.41 \pm 4.21$	$97.64 \pm 9.02$	$94.56 \pm 3.55$
UNREAL	-	$4.13 \pm 3.95$	-	-	-
DARLA <sub>FT</sub>	<b><math>13.36 \pm 5.8</math></b>	$1.4 \pm 2.16$	-	$86.59 \pm 5.53$	$99.25 \pm 2.3$
DARLA <sub>ENT</sub>	$3.45 \pm 4.47$	$15.66 \pm 5.19$	$5.69 \pm 3.73$	$84.77 \pm 4.42$	$59.99 \pm 15.05$
DARLA <sub>DAE</sub>	$7.83 \pm 4.47$	$6.74 \pm 2.81$	$5.59 \pm 3.37$	$85.15 \pm 7.43$	$100.72 \pm 4.7$
<b>DARLA</b>	$10.25 \pm 5.46$	<b><math>19.7 \pm 5.43</math></b>	<b><math>11.41 \pm 3.52</math></b>	<b><math>100.85 \pm 2.92</math></b>	<b><math>108.2 \pm 5.97</math></b>

DARLA'S PERFORMANCE IS SIGNIFICANTLY DIFFERENT FROM ALL BASELINES UNDER WELCH'S UNEQUAL VARIANCES T-TEST WITH  $p < 0.01$  ( $N \in [60, 150]$ ).

# World Models (Ha and Schmidhuber 2018)



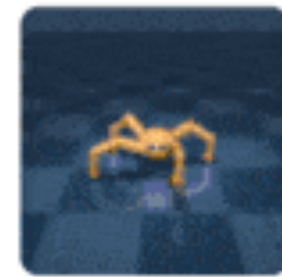
# World Models (Ha and Schmidhuber 2018)



# Dreamer (Hafner et al 2019)



# Dreamer (Hafner et al 2019)



$O_1$

# Quick background

1. Autoencoder
2. Variational Autoencoder
3. Contrastive Learning
4. Siamese Networks
5. Data-Augmentations

**Contrastive-like UL**

# Siamese Networks: Invariance to (augmented) views

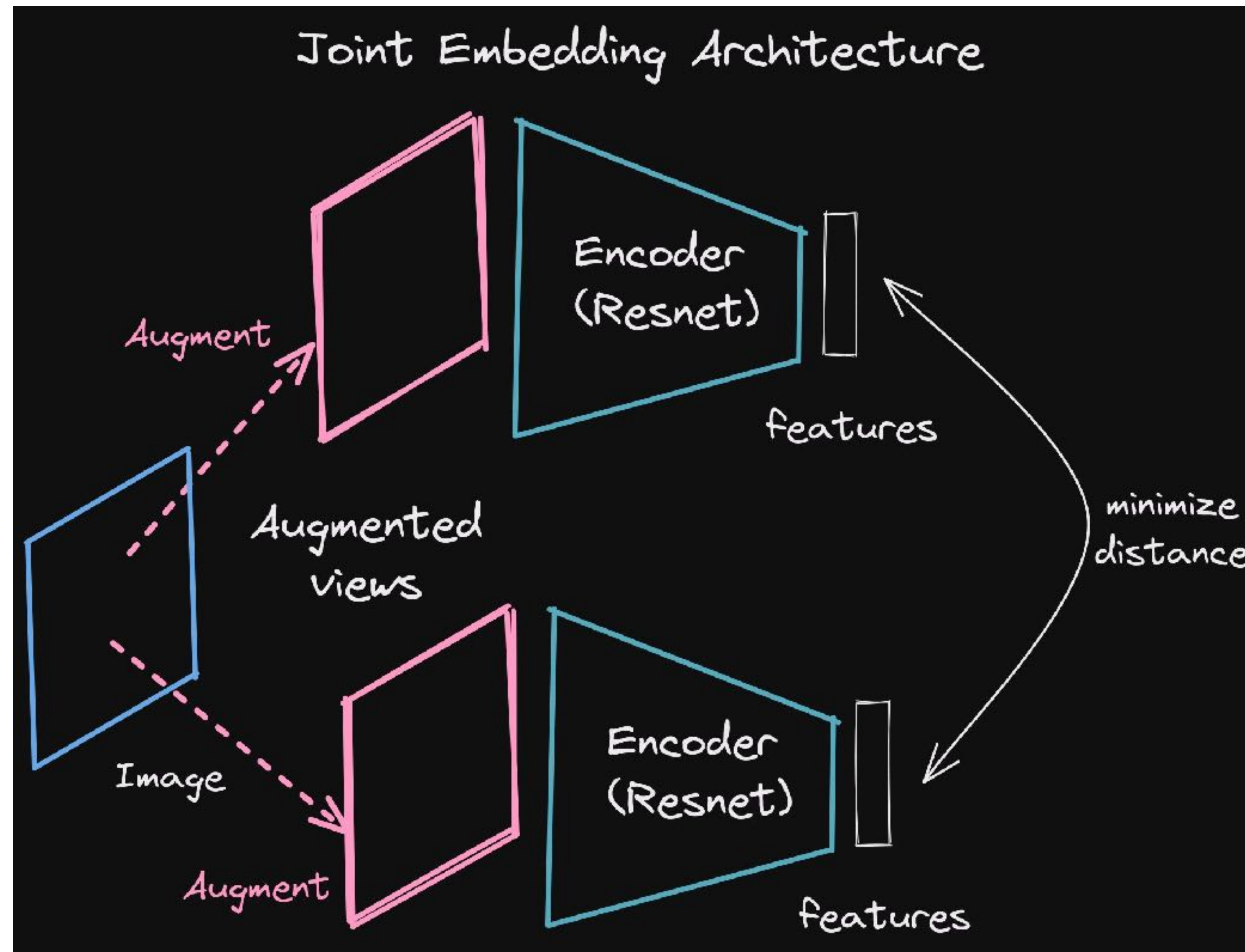
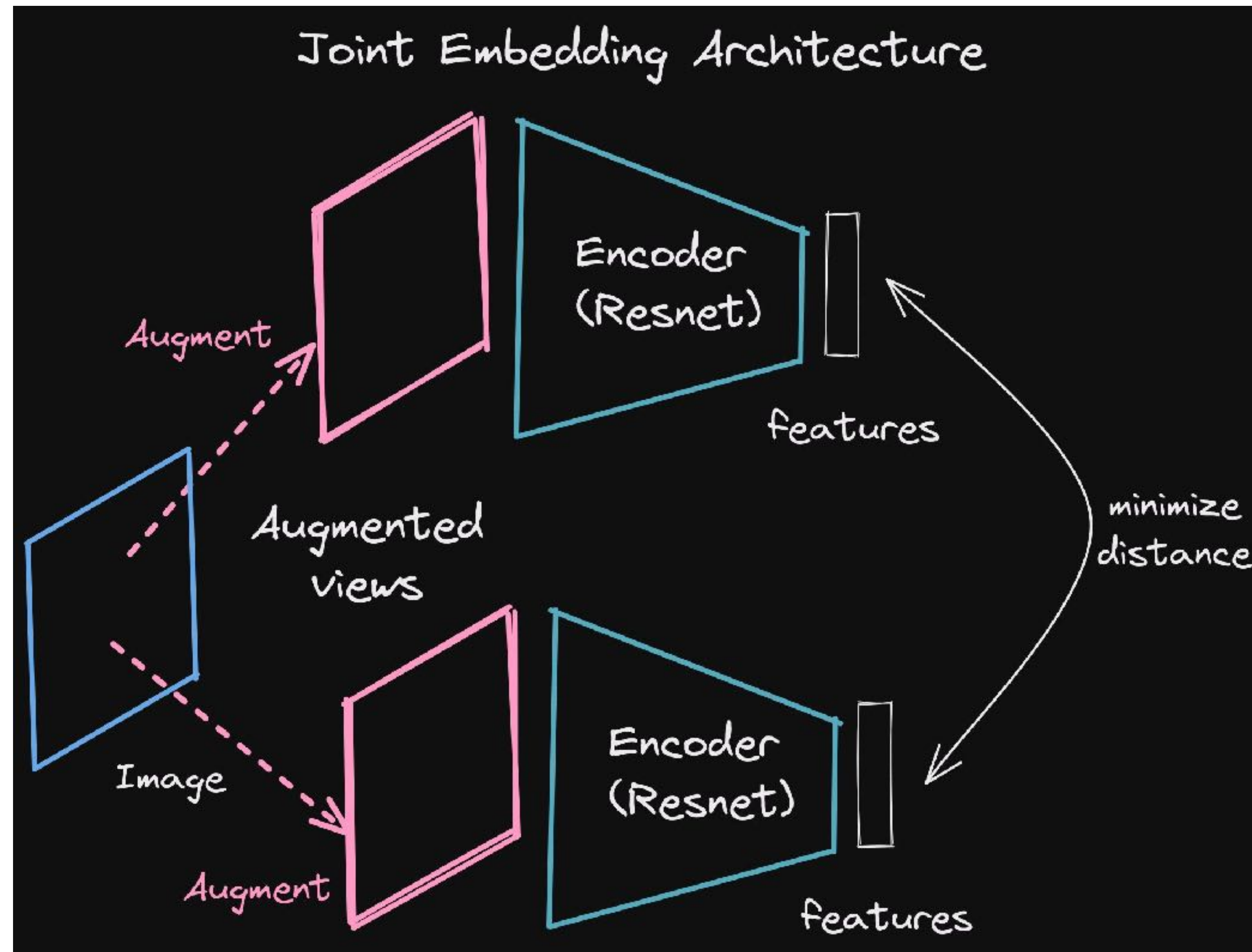


Figure from Antonin Raffin,  
@araffin2

# Siamese Networks: Invariance to (augmented) views



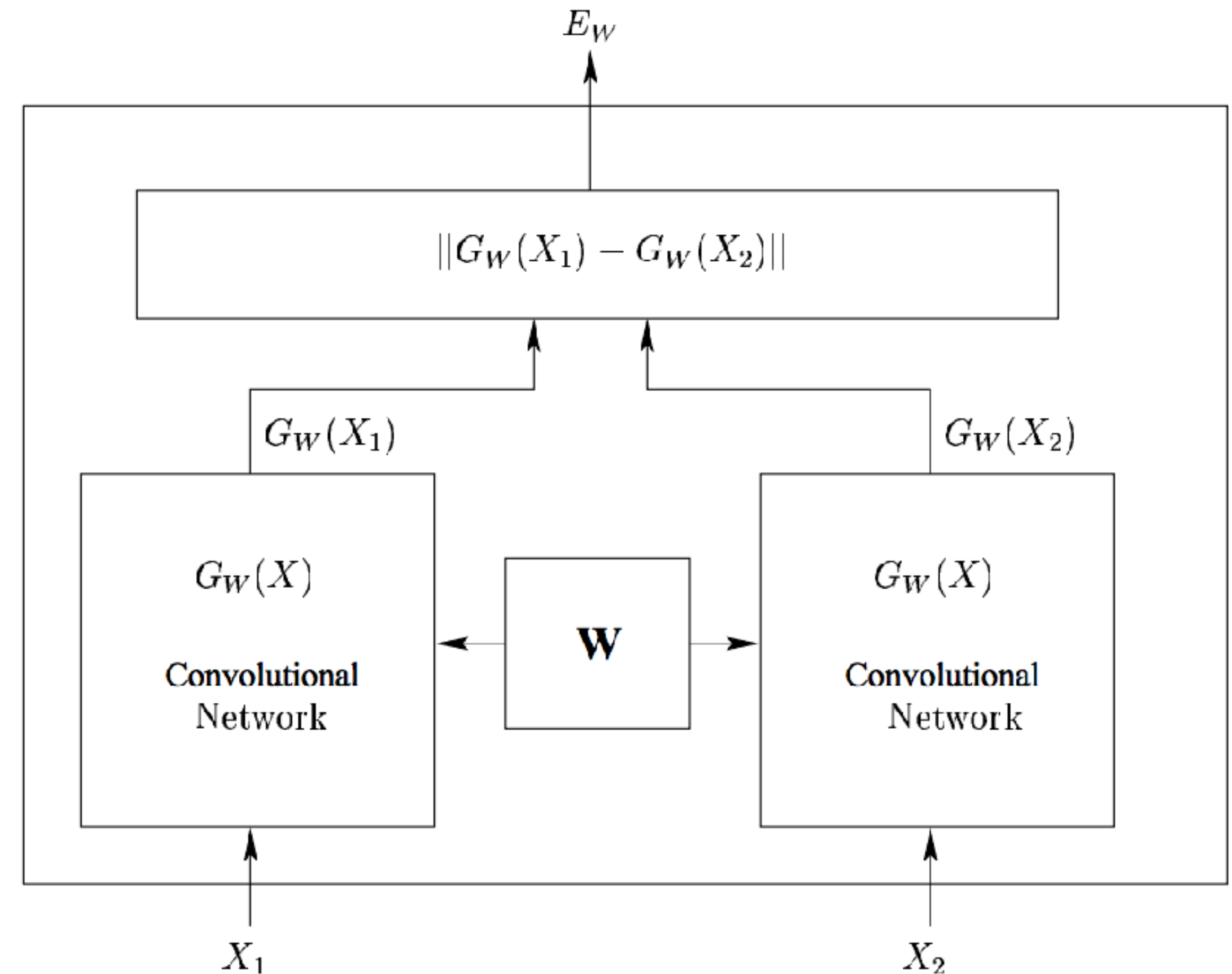
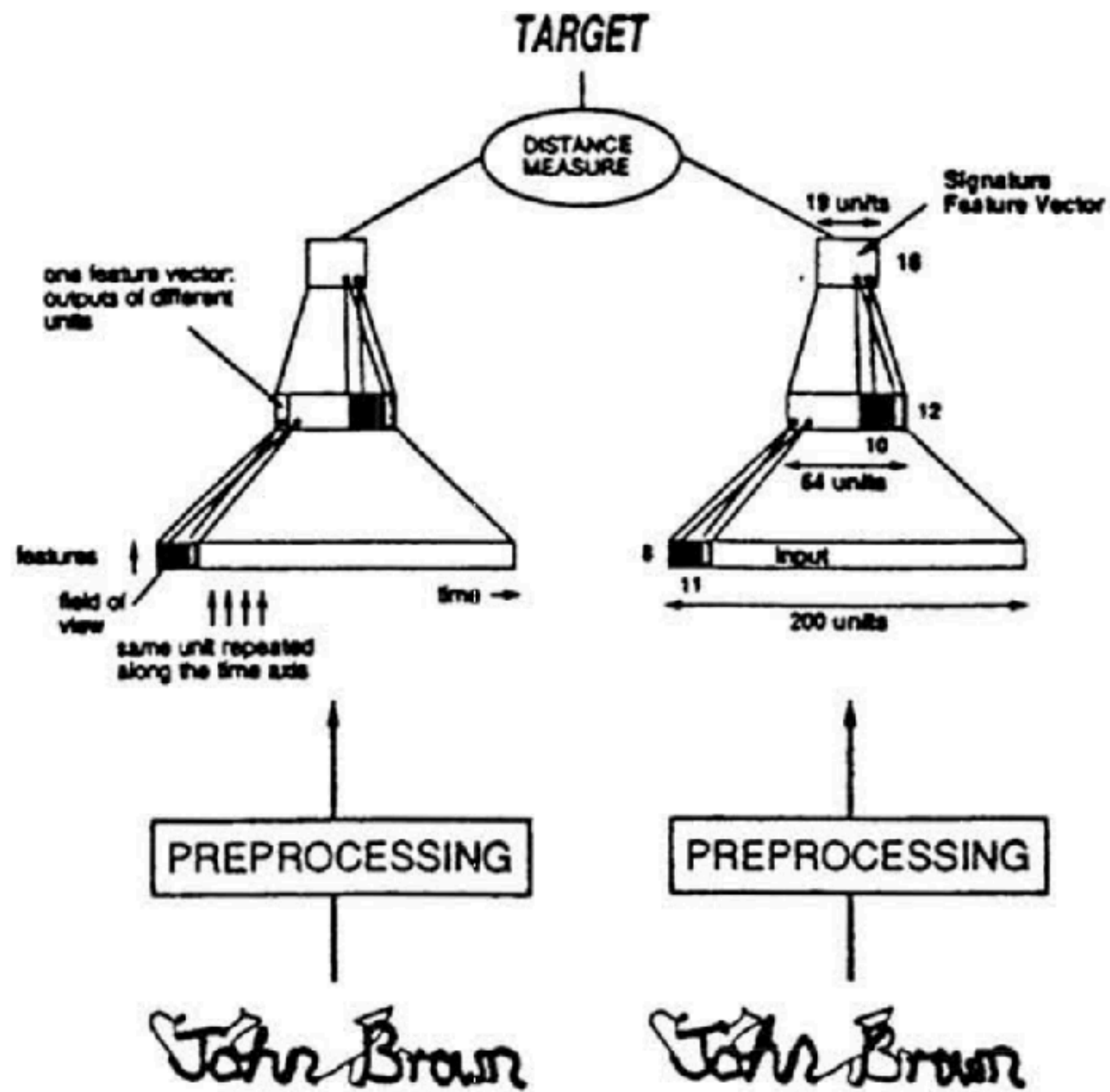
If you have some knowledge about the domain, encode it into the latent space through the form of invariances.

More often than not, this results in emergent representations that capture the most *useful aspects* of the high-dimensional inputs.

Why? Because it ignores what you ask the encoder to be invariant to - and the better you can prescribe that, the more the model knows what **not** to focus on.



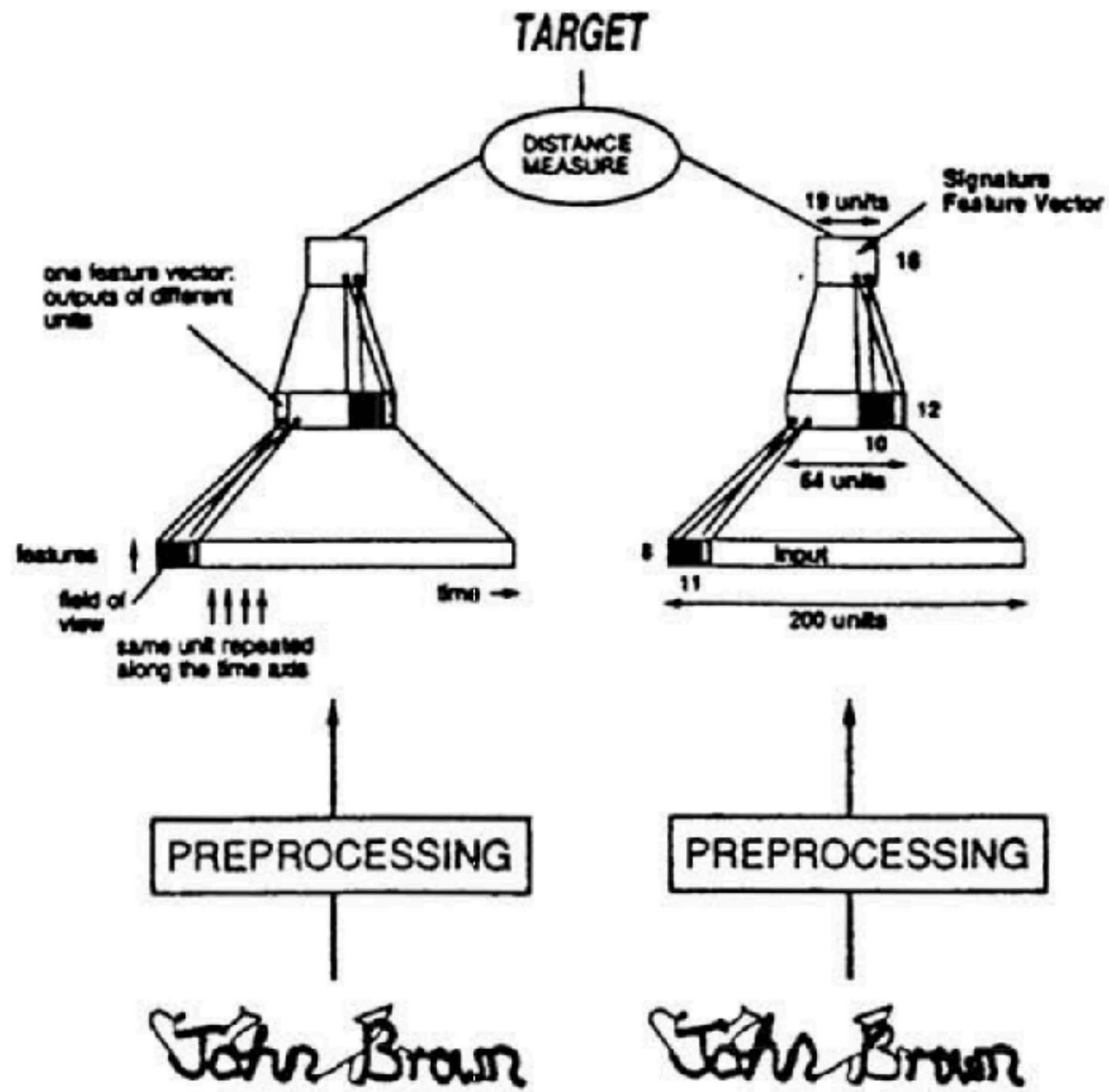
# Siamese Networks: Invariance to (augmented) views



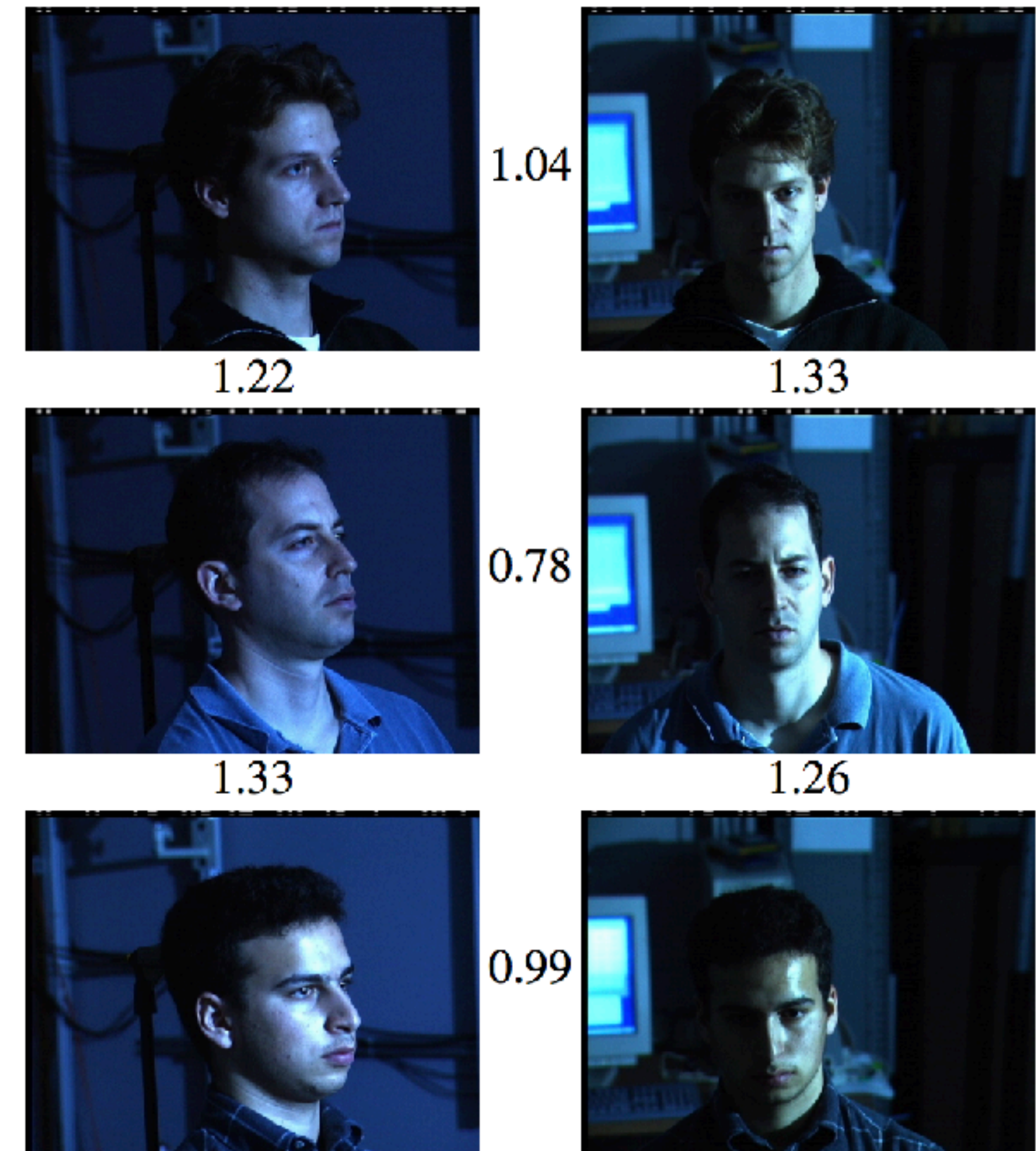
Signature Verification (Bromley, Guyon, LeCun, et al, 93)

Learning a similarity metric discriminatively,  
Chopra, Hadsell, LeCun, 2005

# Siamese Networks: Invariance to (augmented) views



Signature Verification (Bromley, Guyon, LeCun, et al, 93)



FaceNet, Schroff et al 2015

Invariance to pose, lighting, etc - model focuses on encoding facial features

One way to train Siamese Networks is Contrastive Learning

# Contrastive Learning

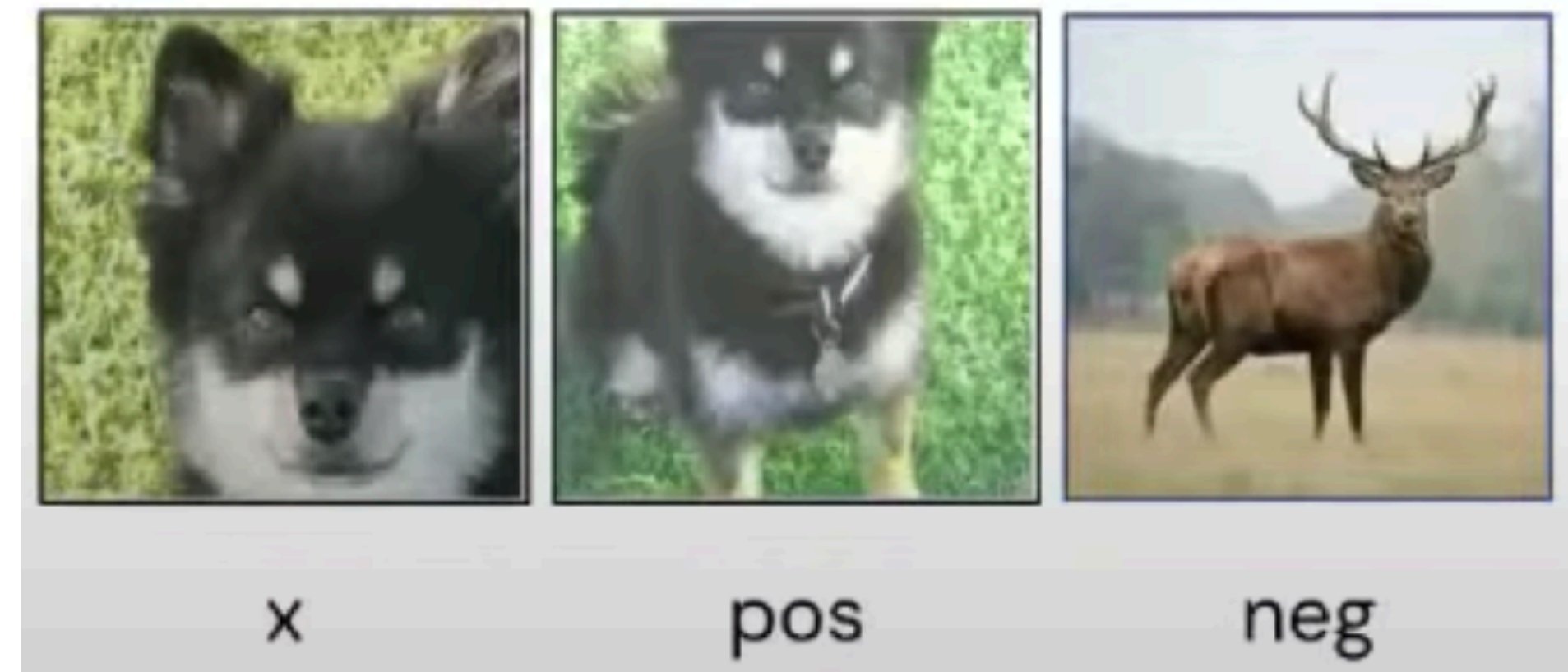
Take a datapoint (an image), and try to fit a scoring function to make sure it aligns more with a positive relative to a negative.

$$\textit{score}(x, x_{pos}) > \textit{score}(x, x_{neg})$$

# Contrastive Learning

Take a datapoint (an image), and try to fit a scoring function to make sure it aligns more with a positive relative to a negative.

$$\text{score}(x, x_{\text{pos}}) > \text{score}(x, x_{\text{neg}})$$



Learn the concept of cats and dogs

Slide adapted from Aaron van den oord

# Contrastive Learning

Take a datapoint (an image), and try to fit a scoring function to make sure it aligns more with a positive relative to a negative.

$$\textit{score}(x, x_{\text{pos}}) > \textit{score}(x, x_{\text{neg}})$$

$$L_{\text{InfoNCE}} = -\mathbb{E} \left[ \log \frac{s(x, x_{\text{pos}})}{s(x, x_{\text{pos}}) + \sum_{y_j \neq x_{\text{pos}}} s(x, y_j)} \right]$$

# Contrastive Learning

Take a datapoint (an image), and try to fit a scoring function to make sure it aligns more with a positive relative to a negative.

$$\textit{score}(x, x_{\text{pos}}) > \textit{score}(x, x_{\text{neg}})$$

$$L_{\text{InfoNCE}} = -\mathbb{E} \left[ \log \frac{s(x, x_{\text{pos}})}{s(x, x_{\text{pos}}) + \sum_{y_j \neq x_{\text{pos}}} s(x, y_j)} \right]$$

InfoNCE (or N-Pairs) Loss

1. Representation Learning with Contrastive Predictive Coding (van den Oord et al 2018)
2. Improved Deep Metric Learning with Multi-Class N-Pairs Loss - (Sohn et al 2016)
3. Deep InfoMax, AMDIM (Hjelm, Bachman, et al 2019)

# Contrastive Learning

Take a datapoint (an image), and try to fit a scoring function to make sure it aligns more with a positive relative to a negative.

$$\textit{score}(x, x_{pos}) > \textit{score}(x, x_{neg})$$

$$L_{\text{InfoNCE}} = -\mathbb{E} \left[ \log \frac{s(x, y)}{\sum_{y_j} s(x, y_j)} \right]$$

InfoNCE (or N-Pairs) Loss

1. Representation Learning with Contrastive Predictive Coding (van den Oord et al 2018)
2. Improved Deep Metric Learning with Multi-Class N-Pairs Loss - (Sohn et al 2016)
3. Deep InfoMax, AMDIM (Hjelm, Bachman, et al 2019)

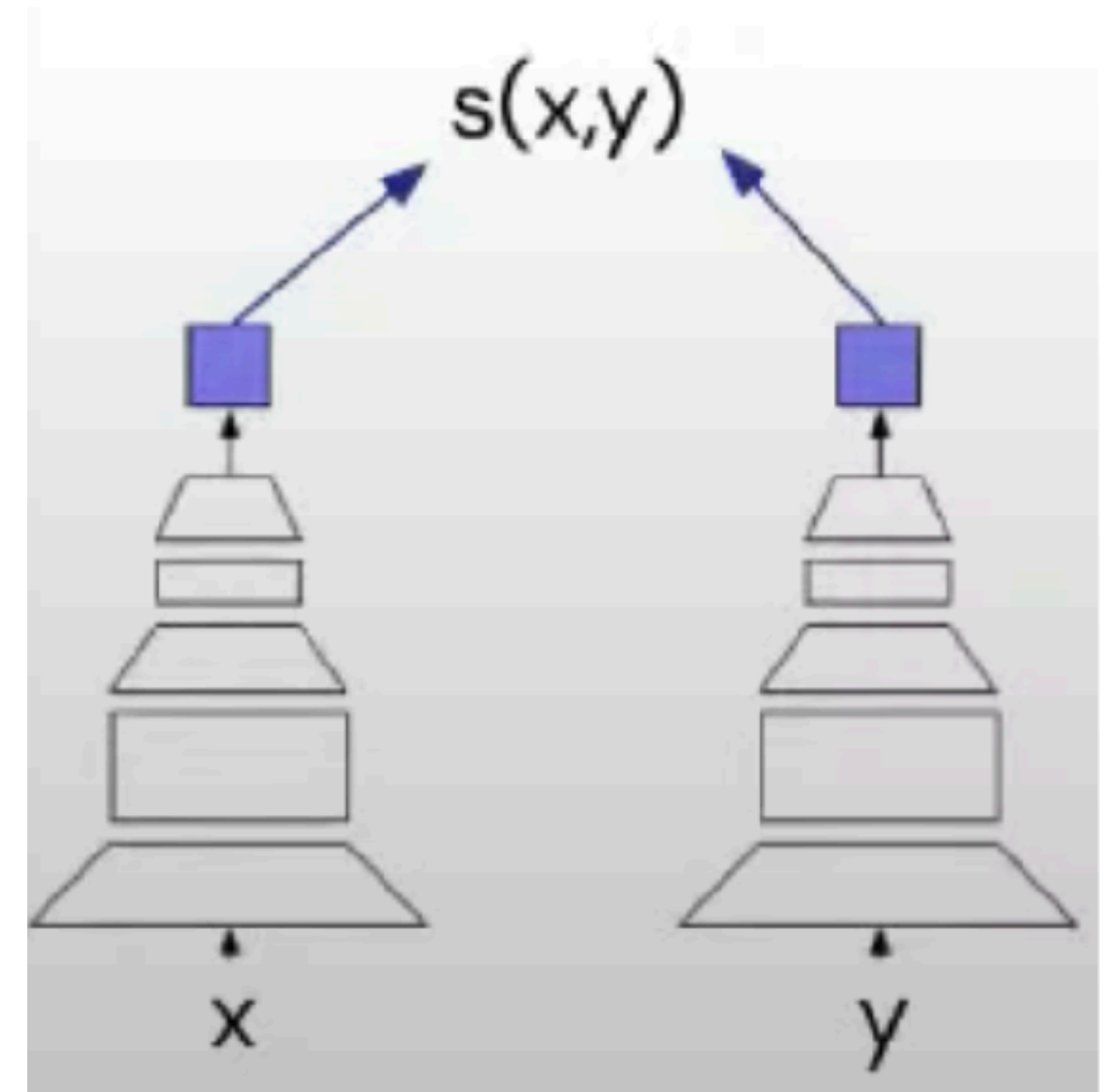


# Contrastive Learning: InfoNCE Loss

$$L_{\text{InfoNCE}} = -\mathbb{E} \left[ \log \frac{s(x, y)}{\sum_{y_j} s(x, y_j)} \right]$$

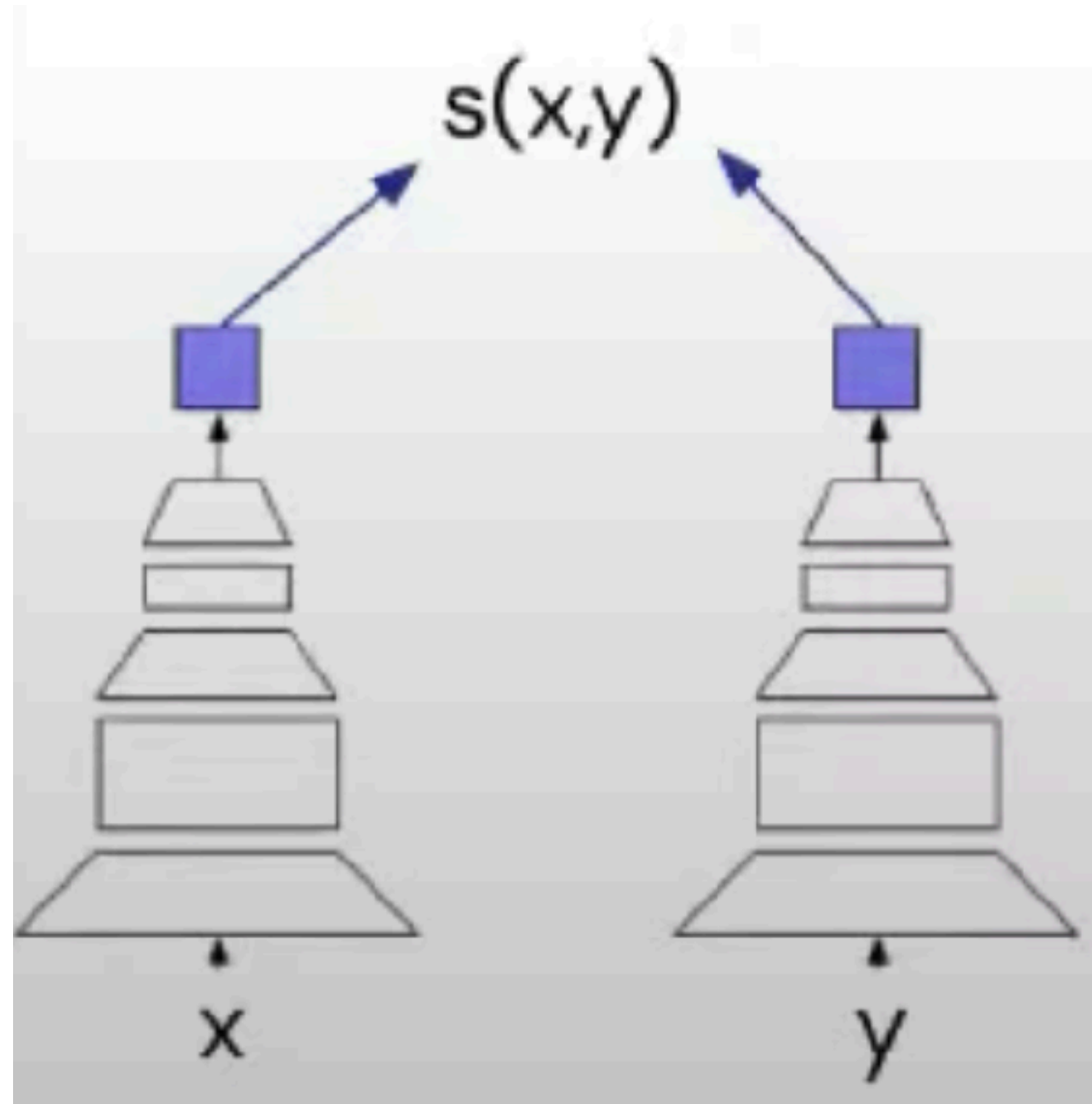
$$s(x, y) = e^{f_X(x)^T f_Y(y)}$$

$$L_{\text{InfoNCE}} = -\mathbb{E} \left[ \log \frac{e^{f_X(x)^T f_Y(y)}}{\sum_{y_j} e^{f_X(x)^T f_Y(y_j)}} \right]$$



1. Representation Learning with Contrastive Predictive Coding (van den Oord et al 2018)
2. Improved Deep Metric Learning with Multi-Class N-Pairs Loss - (Sohn et al 2016)
3. Deep InfoMax, AMDIM (Hjelm, Bachman, et al 2019)

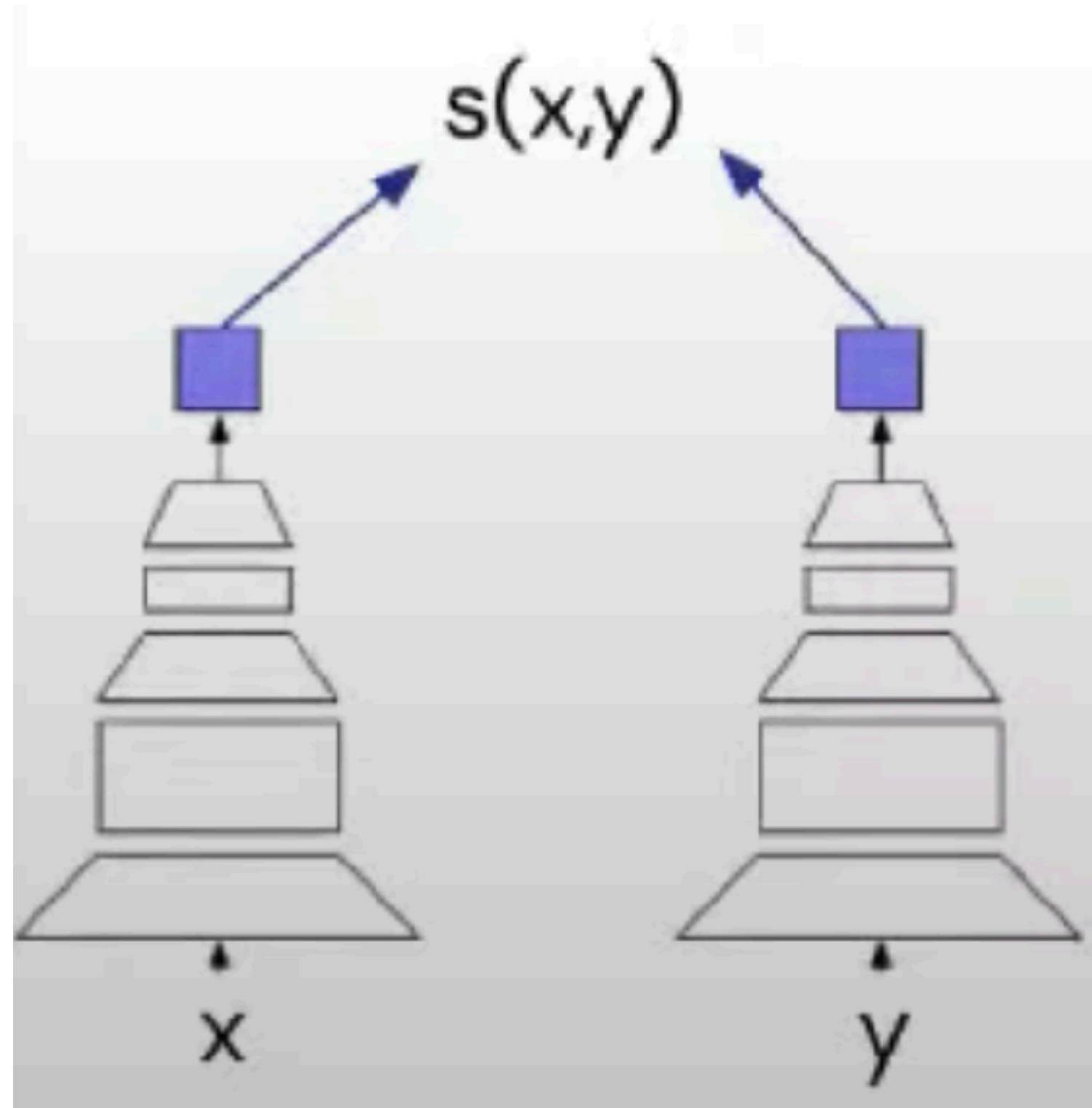
# Contrastive Learning: InfoNCE Loss



1. Maximize mutual information between the views  $x$  and  $y$ .
2. Convert it to a classification problem, optimized with a stable cross-entropy loss.
3. Representations capture things *common* to  $x$  and  $y$ .
4. By augmenting  $x$  and  $y$  in several ways (stuff you don't want to capture), you only capture the *relevant left-over* common things.

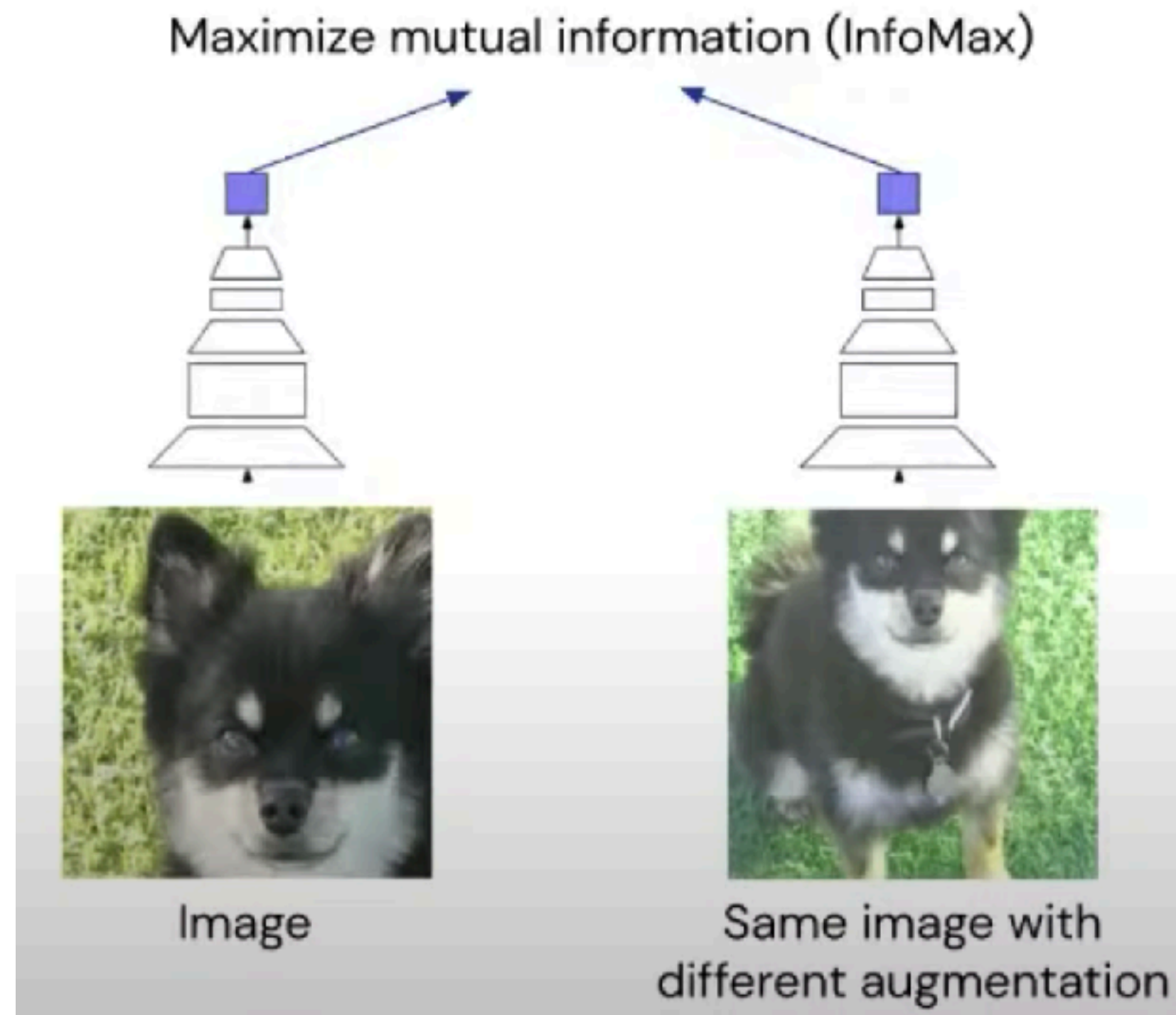
1. Representation Learning with Contrastive Predictive Coding (van den Oord et al 2018)
2. Improved Deep Metric Learning with Multi-Class N-Pairs Loss - (Sohn et al 2016)
3. Deep InfoMax, AMDIM (Hjelm, Bachman, et al 2019)

# Contrastive Learning: InfoNCE Loss



1. Plenty of choices for what  $x$  and  $y$  can be.
2.  $x$  and  $y$  could be two augmented views of the same image.
3.  $x$  could be the past frame (aug),  $y$  could be the future frame (aug)
4.  $x$  could be past frame (aug) + action,  $y$  could be future (aug)...
5. Whatever it is, you are optimizing for MI between  $x$  and  $y$  (lower bound)

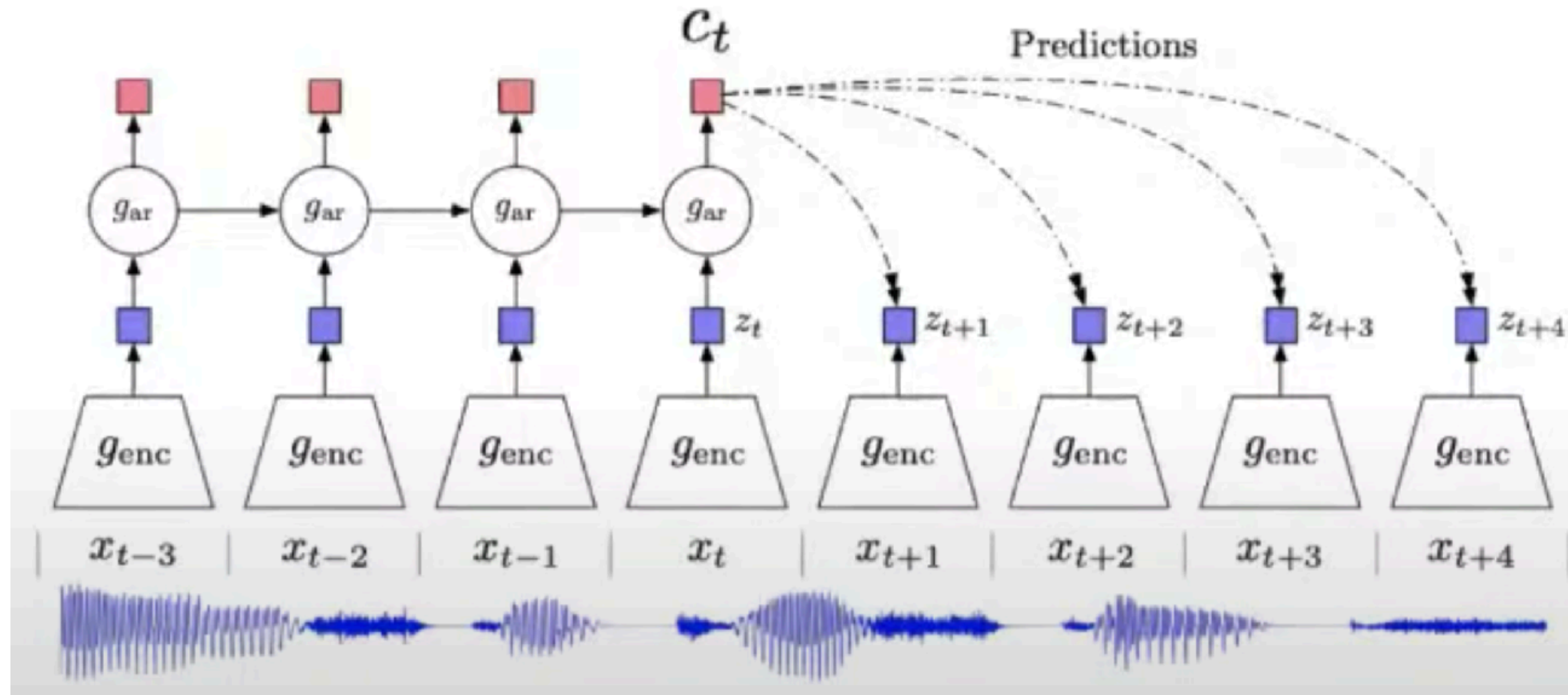
# Contrastive Learning: InfoNCE Loss



1. Plenty of choices for what  $x$  and  $y$  can be.
2.  $x$  and  $y$  could be two augmented views of the same image.
3.  $x$  could be the past frame (aug),  $y$  could be the future frame (aug)
4.  $x$  could be past frame (aug) + action,  $y$  could be future (aug)...
5. Whatever it is, you are optimizing for MI between  $x$  and  $y$  (lower bound)

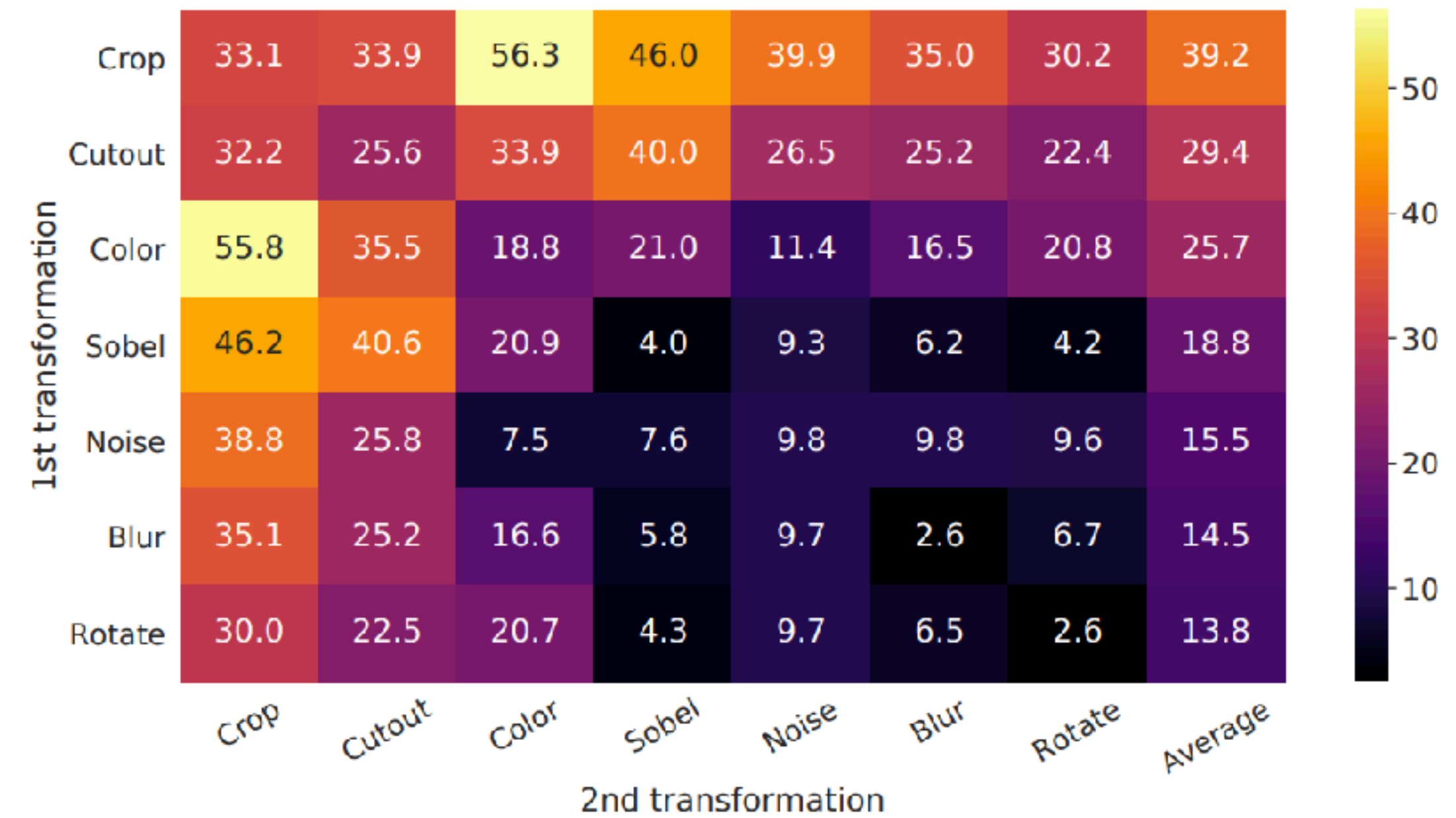
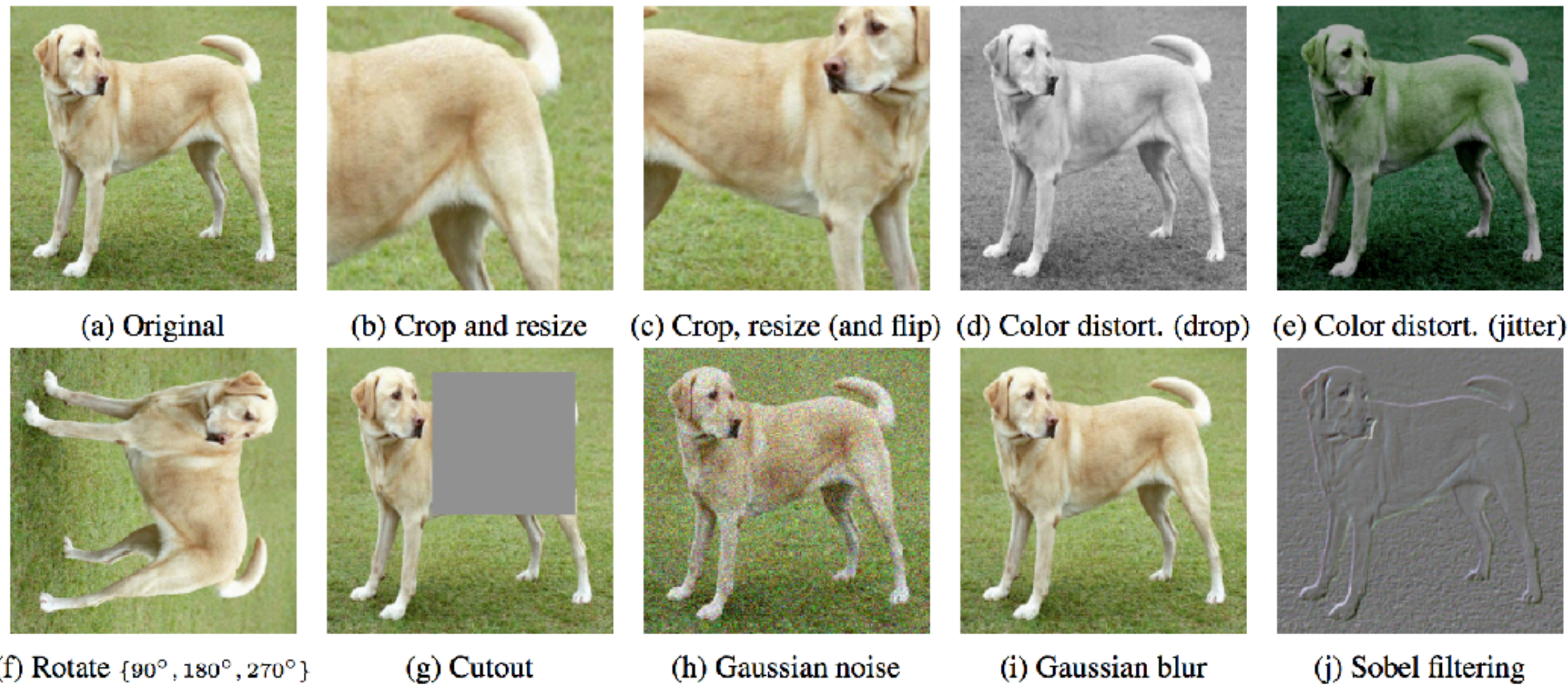
# Contrastive Learning: SimCLR

# Contrastive Learning: Contrastive Predictive Coding (CPC)

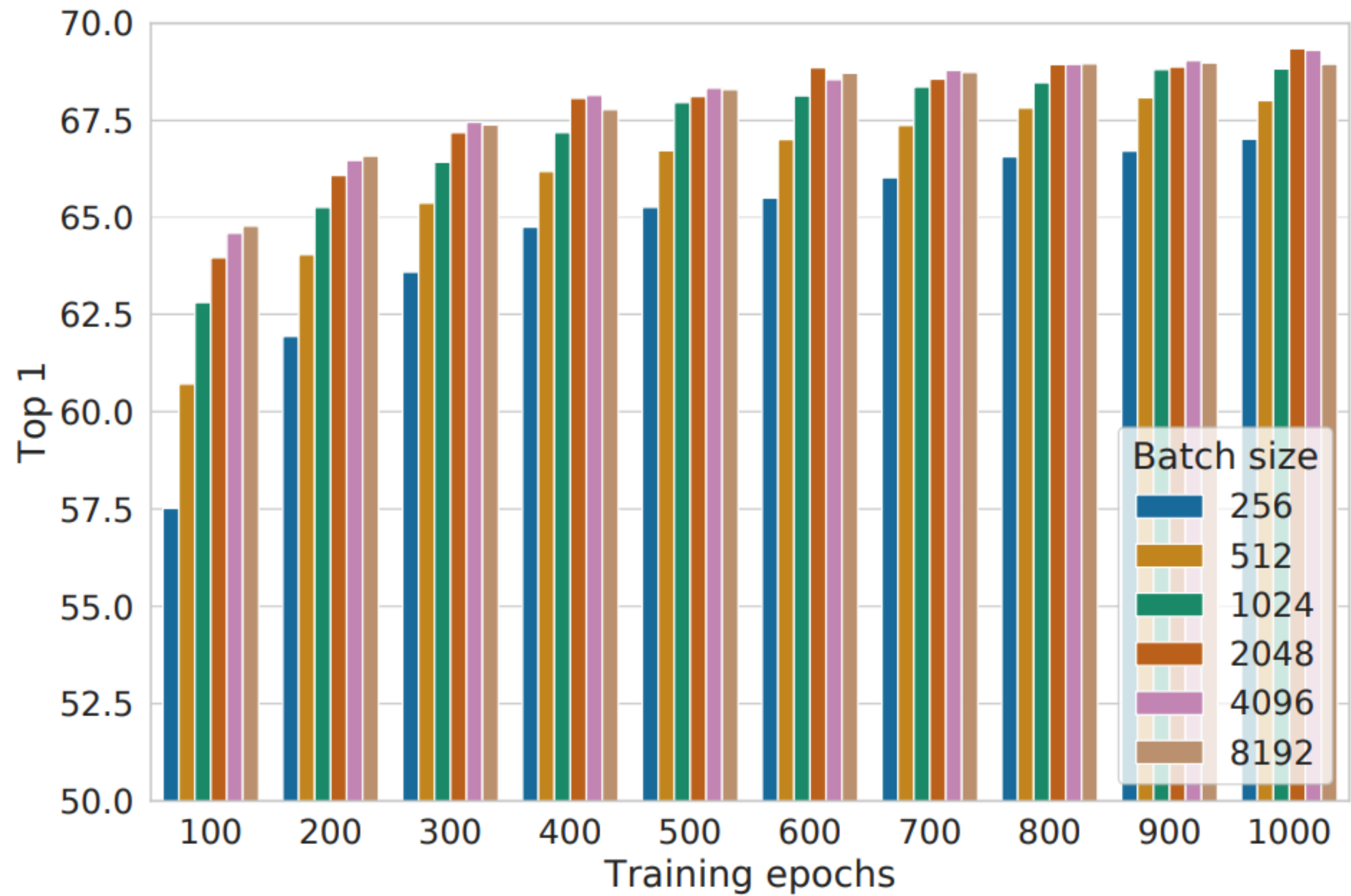


1. Representation Learning using Contrastive Predictive Coding, van den Oord et al 2018
2. Data-Efficient Image Recognition using Contrastive Predictive Coding, Henaff et al 2020

# SimCLR: Importance of Augs



# SimCLR: Importance of Negatives





# Challenges in Contrastive Learning

# Contrastive Learning: Negatives are tricky



x



pos



neg

# Contrastive Learning: Negatives are tricky



x



pos



neg

# Contrastive Learning: Negatives are tricky



x



pos



neg

# Contrastive Learning: Negatives are tricky



x

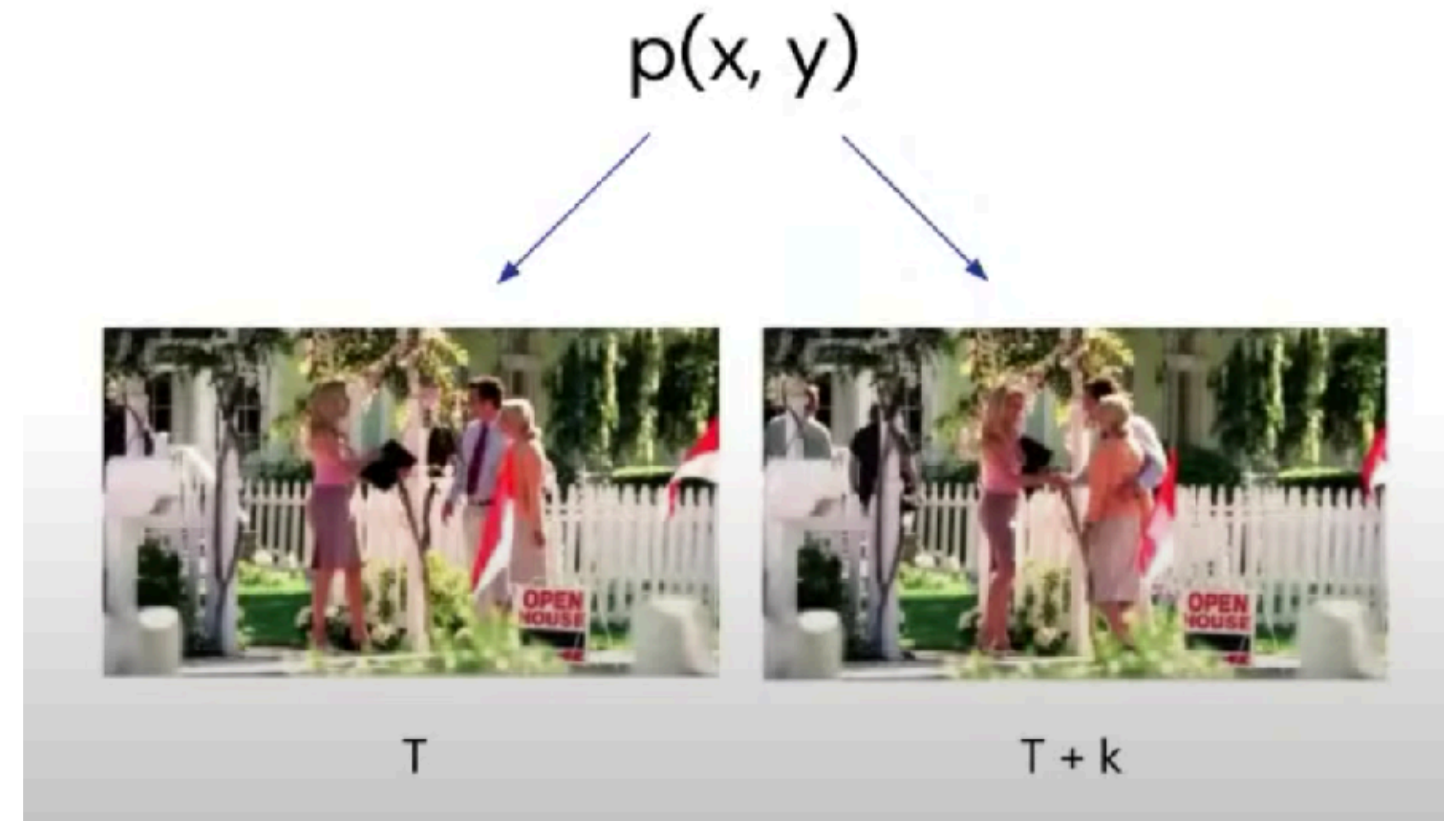
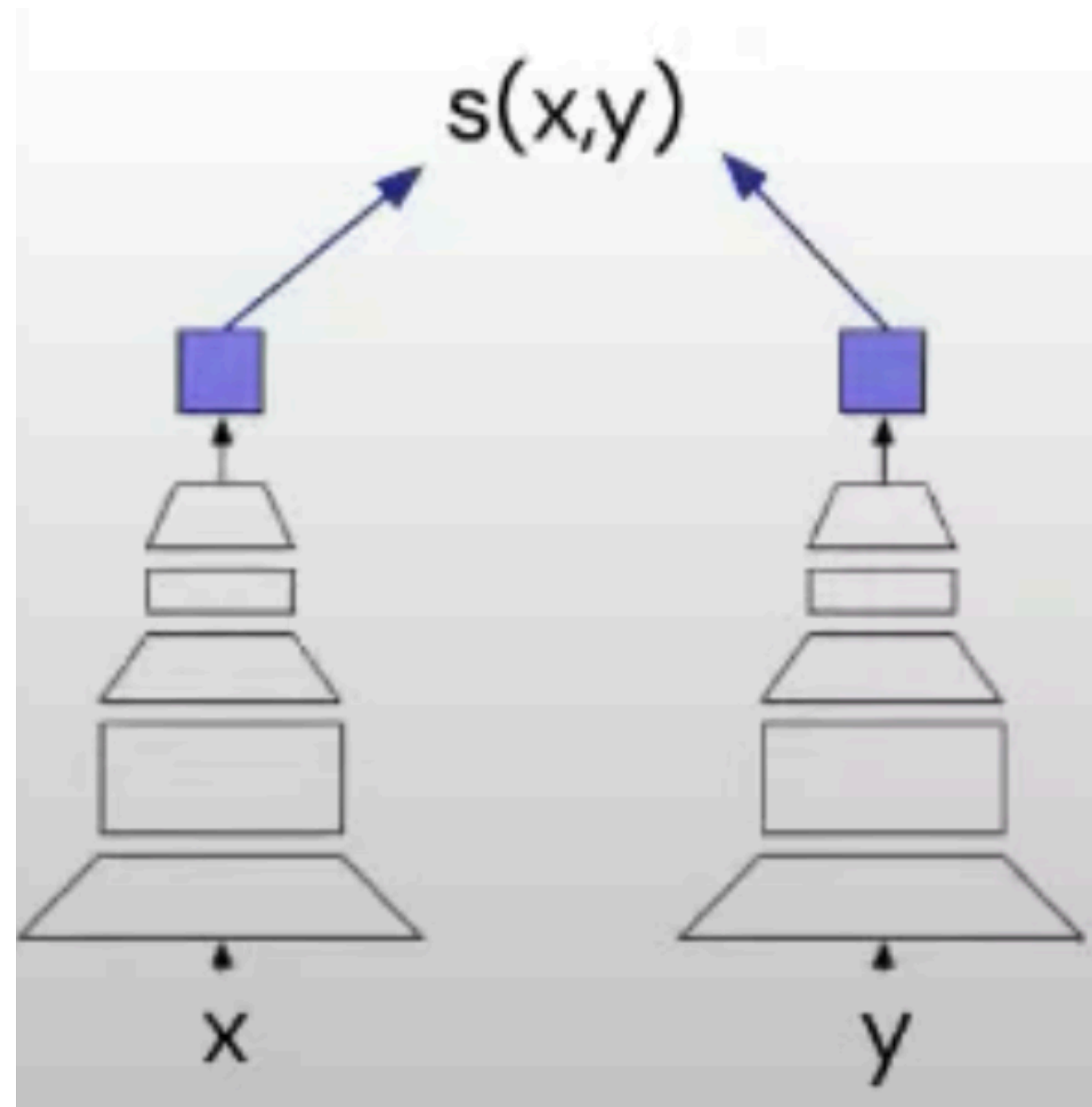


pos

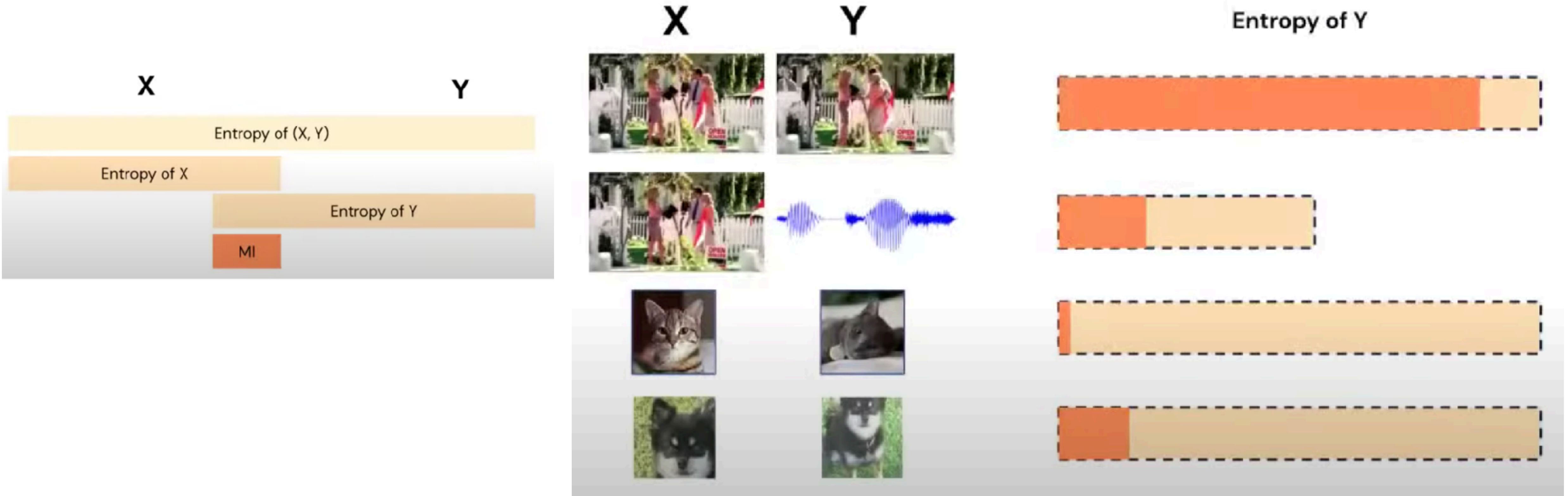


neg

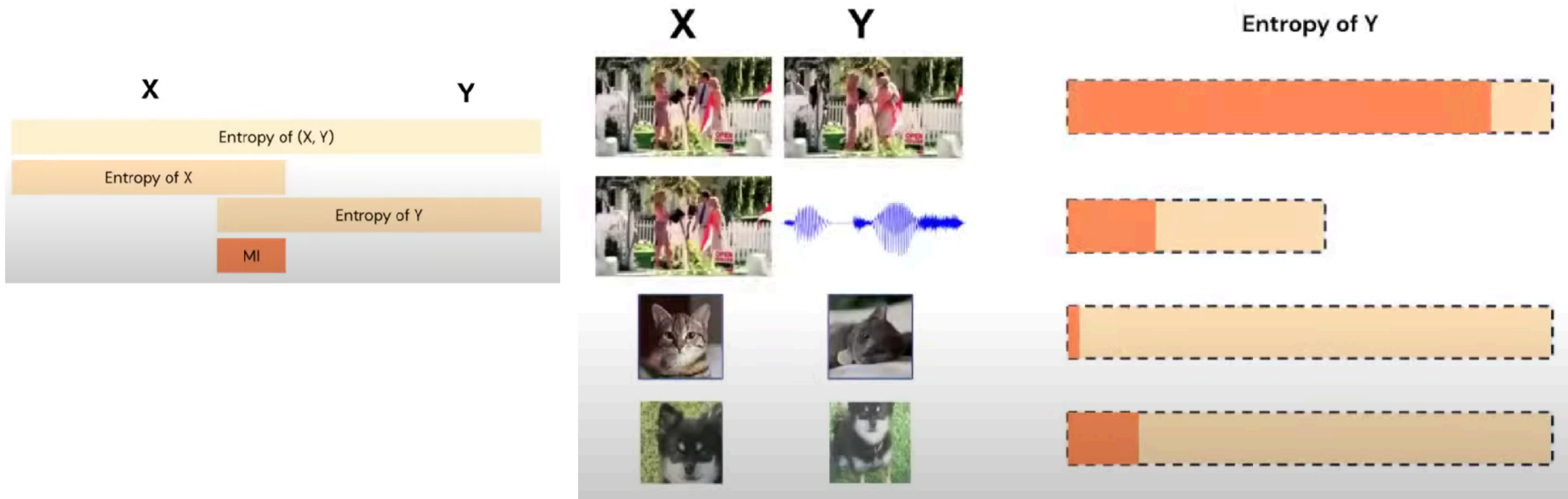
# High MI between anchor and positive



# High MI between anchor and positive



# High MI between anchor and positive



Please do check out Aaron Van den Oord's talks on this topic

Slide adapted from Aaron van den Oord



**Contrastive Learning is just one way to learn Siamese-style Networks**

**There are approaches that work without negatives (BYOL, SimSiam, DINO, Barlow Twins)**

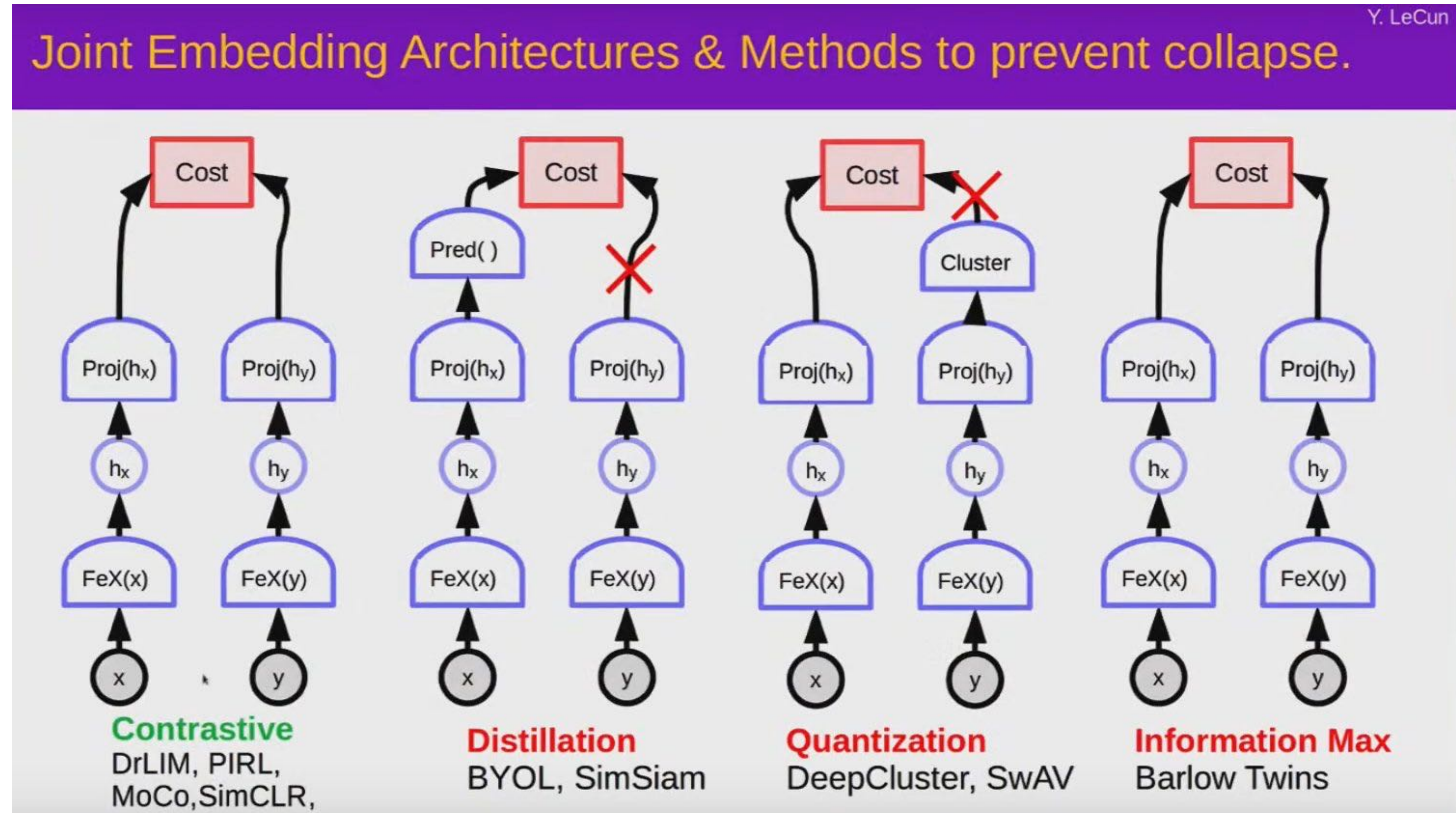
# Cluster Swapping using Augmentations (SwaV)

1. Unsupervised Learning of Visual Features by Contrasting Cluster Assignments (Caron et al 2020)

# BYOL, SimSiam, DINO

1. Bootstrap Your Own Latent (Grill et al 2020)
2. Exploring Simple Siamese Representation Learning - (Chen & He 2021)
3. Emerging properties in self-supervised vision transformers (Caron et al 2021)

# Yann LeCun's summary of all Siamese approaches



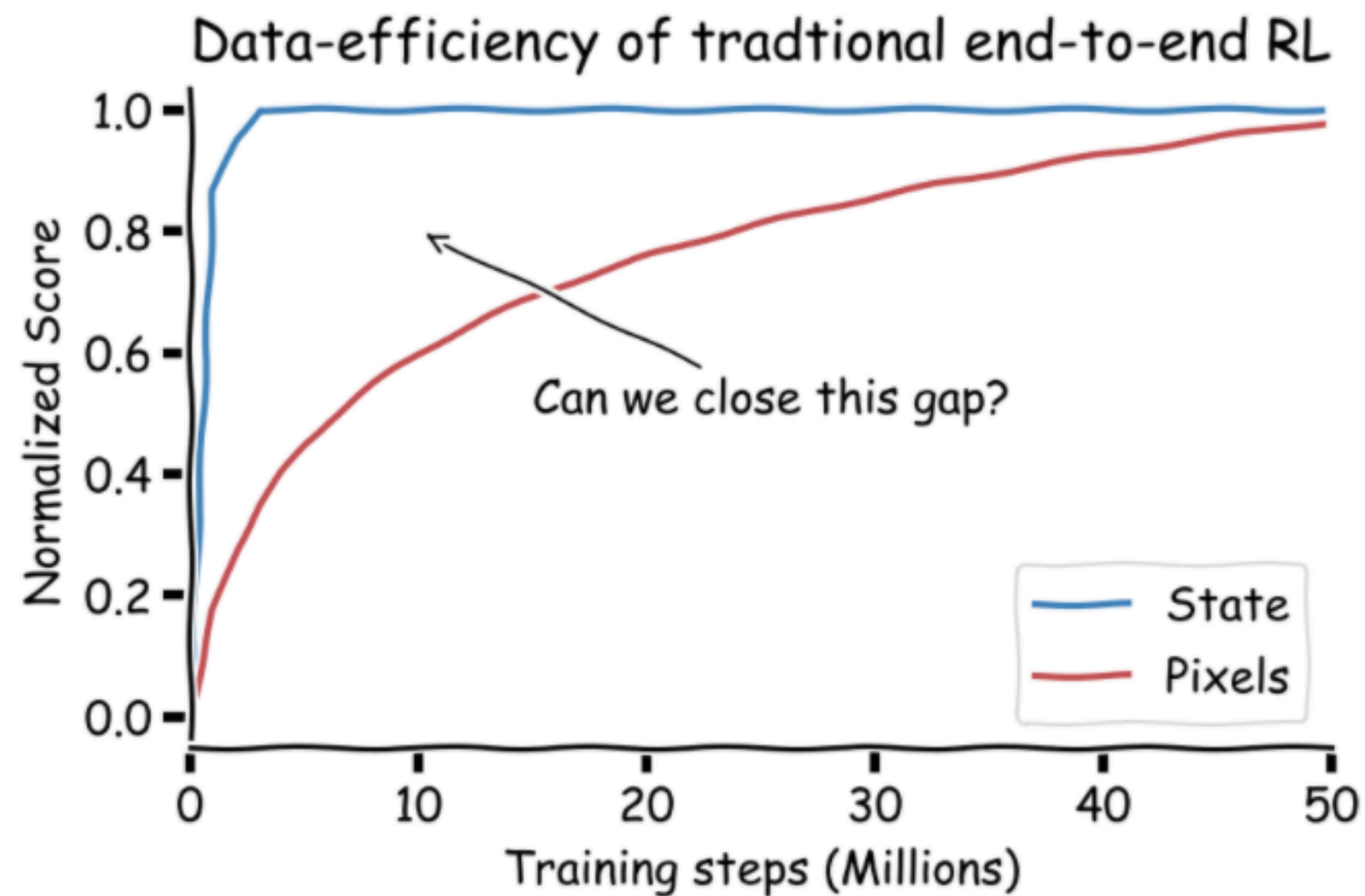
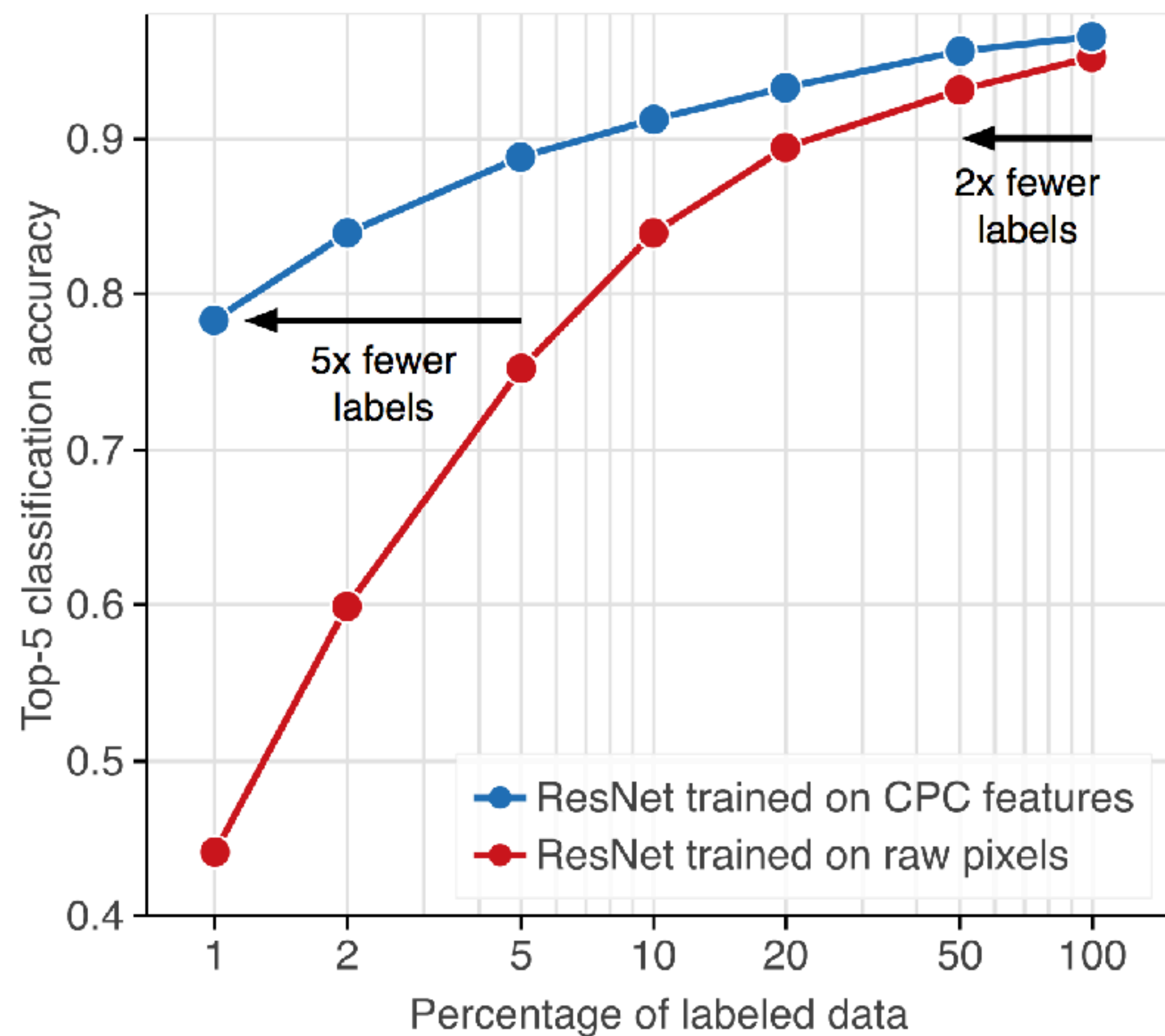
# Quick background

1. Autoencoder
2. Variational Autoencoder
3. Contrastive Learning
4. Siamese Networks
5. Data-Augmentations

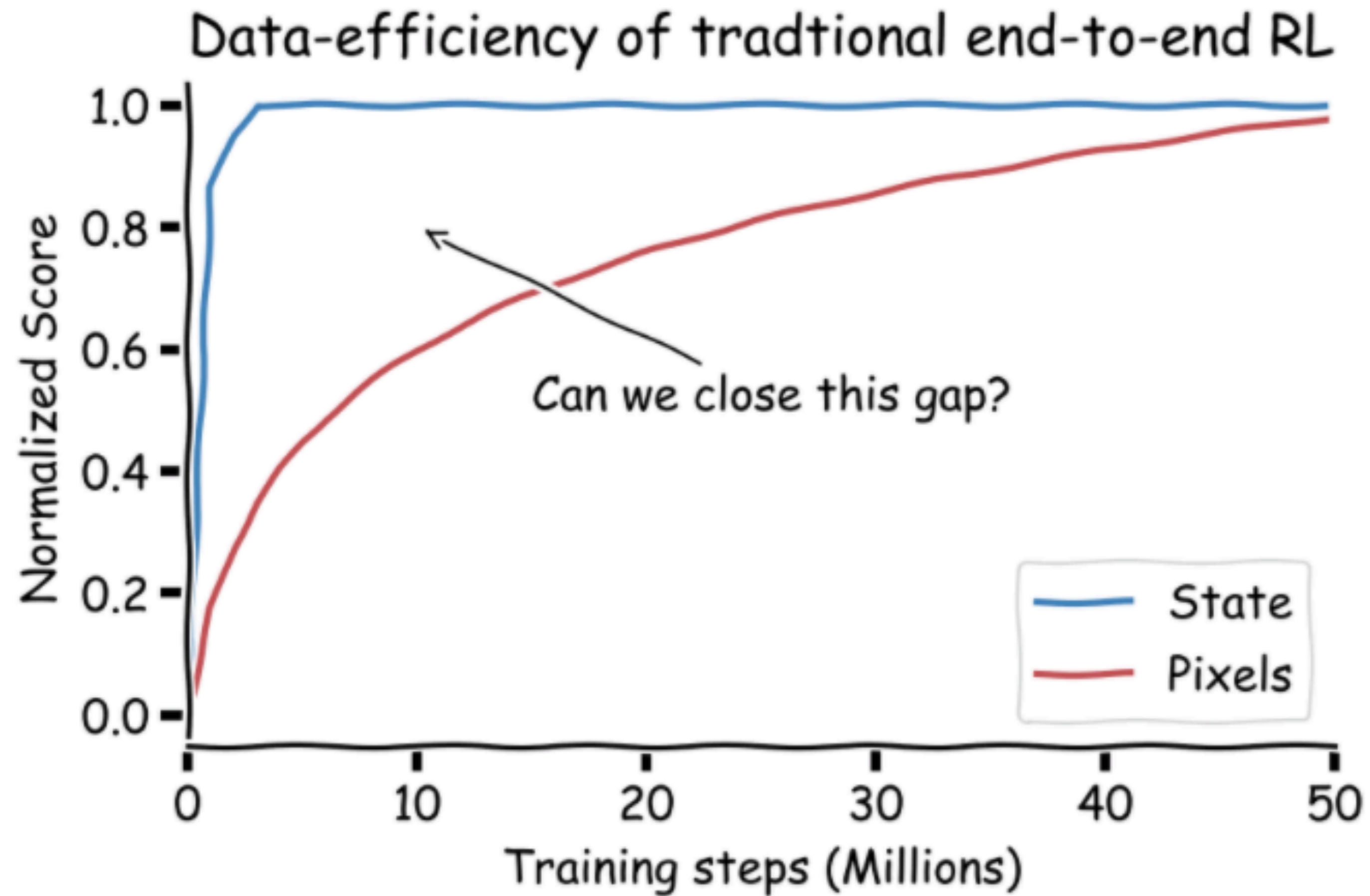
**Contrastive-like UL**

How have people applied Contrastive Learning, Siamese Nets, Data Augmentations, etc in Reinforcement Learning?

# Capture useful aspects of high dimensional sensory stream

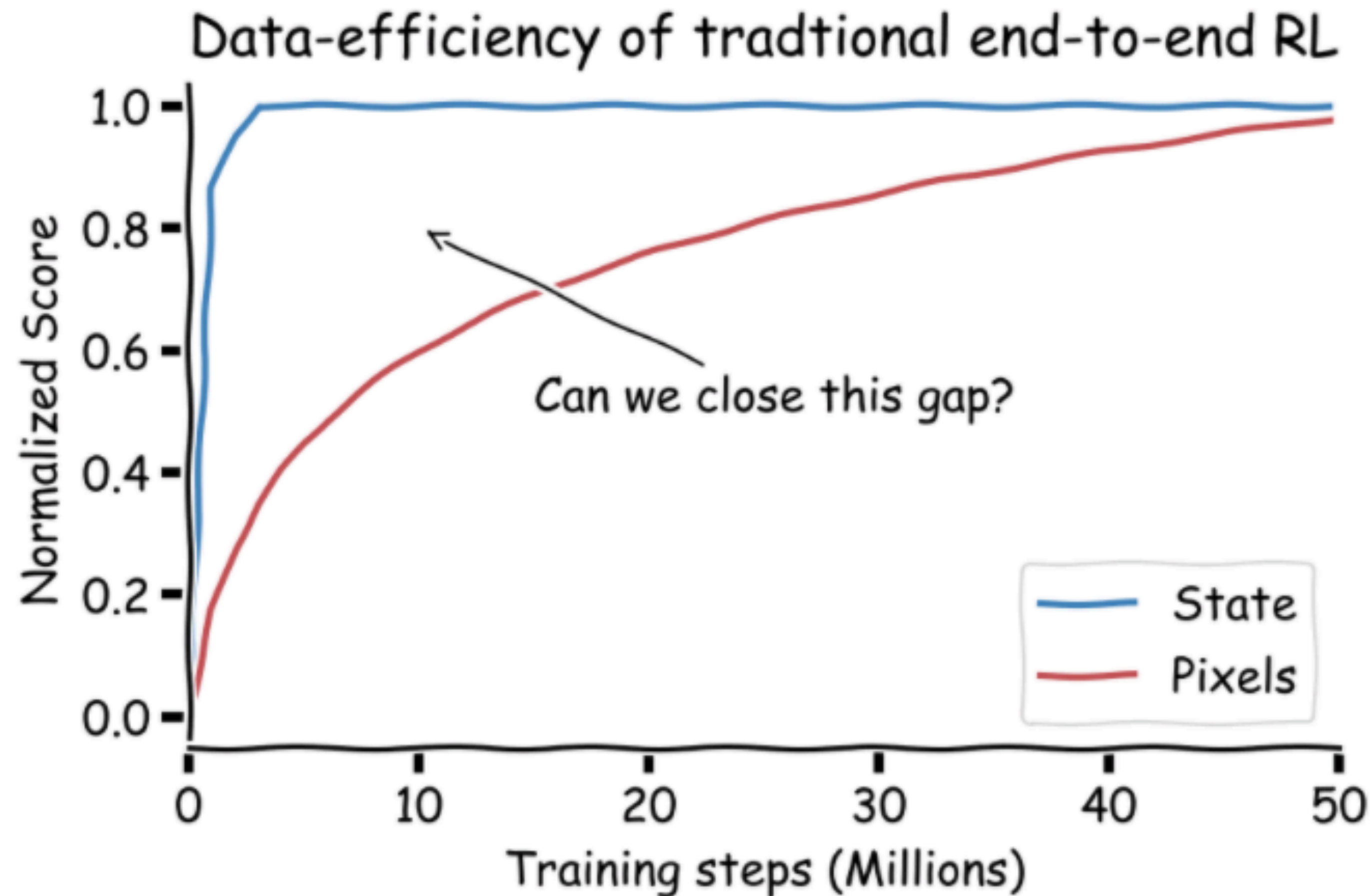


# Data-Efficiency in Reinforcement Learning





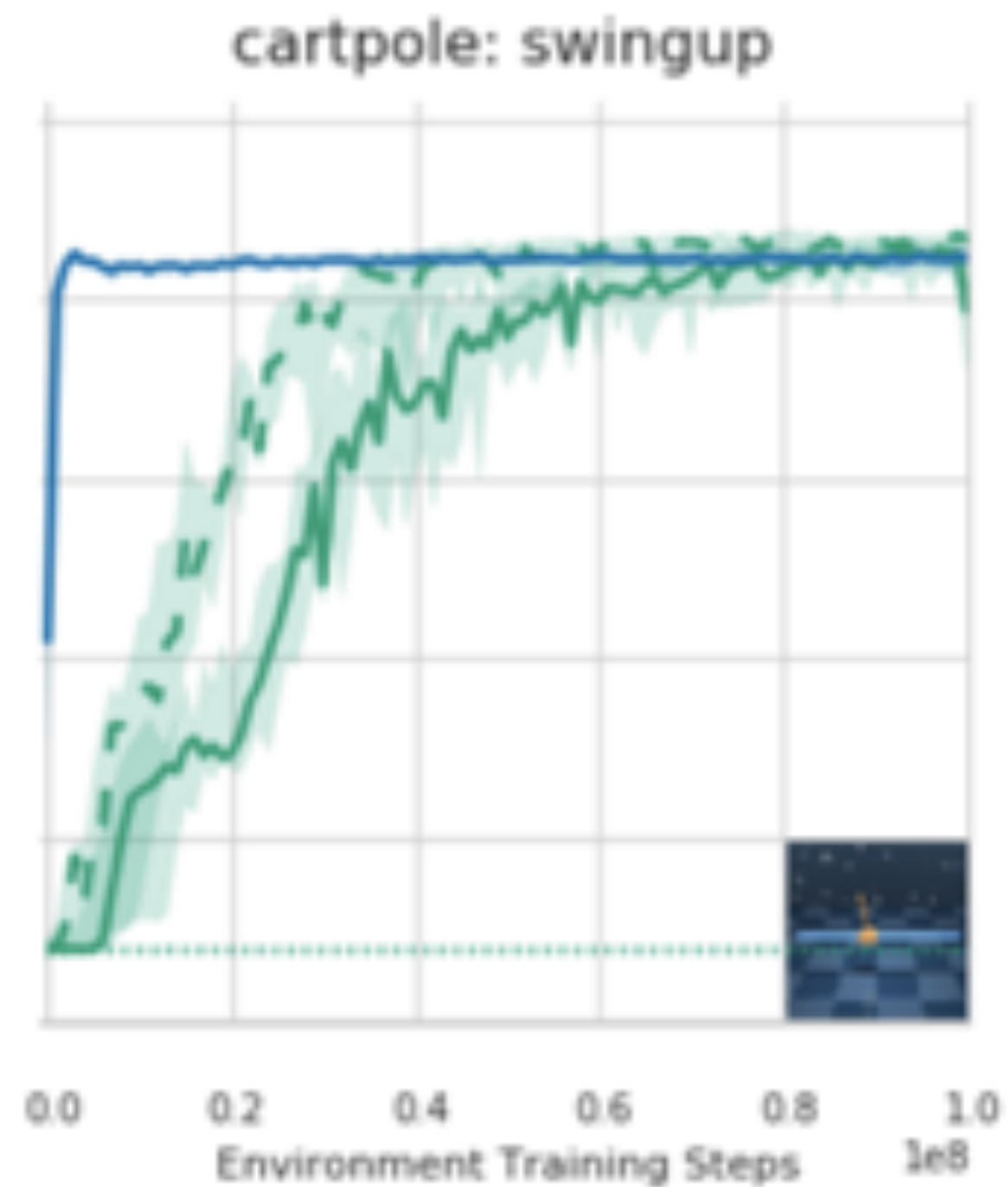
# Data-Efficiency in Reinforcement Learning



1. Code becomes a lot simpler. Less feature engineering.
2. Calibration, sim2real etc become a lot simpler.
3. Take advantage of all the DL scale infra that's been built for vision and NLP, for robotics.
4. Sensor costs.
5. MLPs are not as parameter and FLOP efficient as transformers and ConvNets.

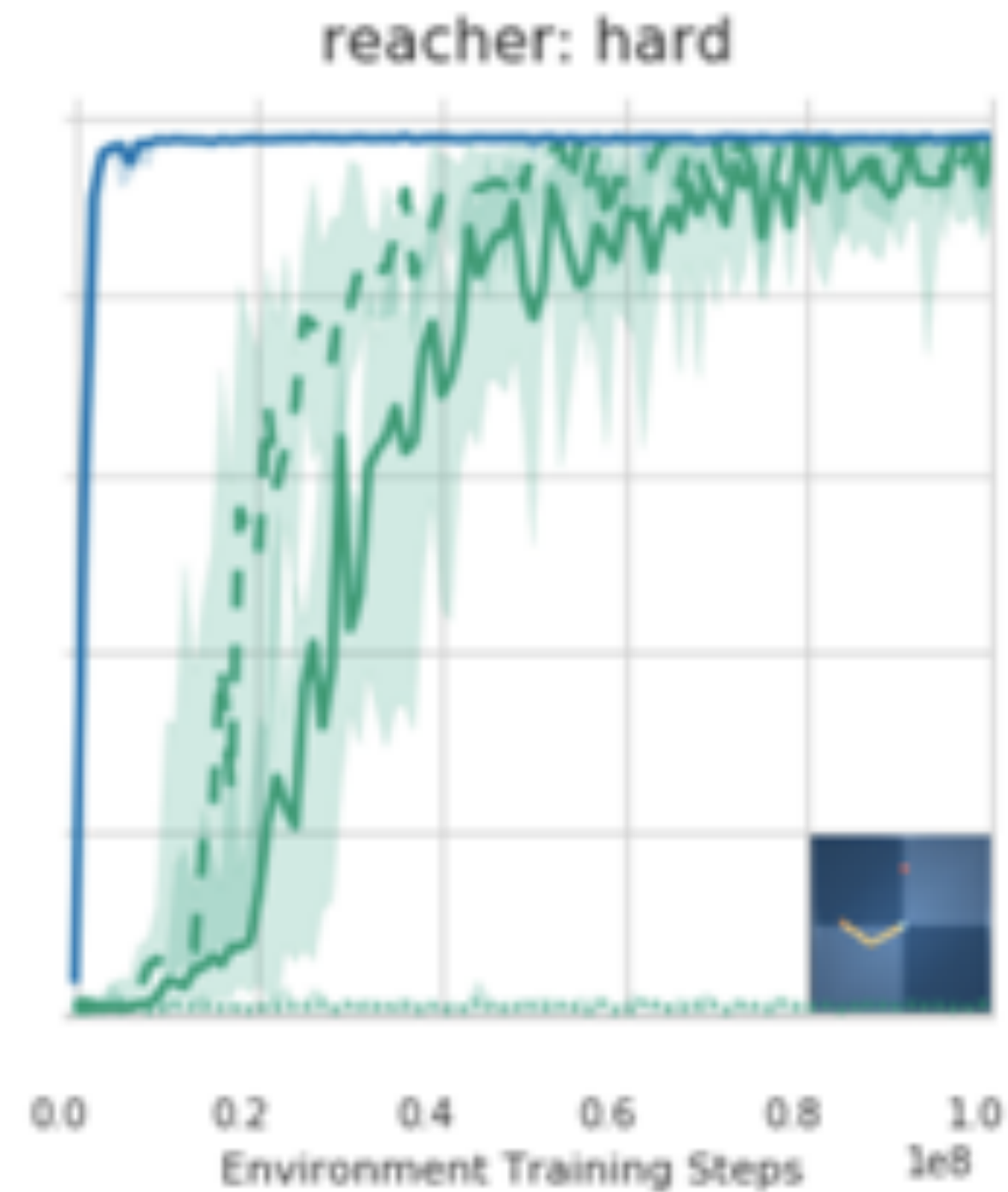
# Data-Efficiency in Reinforcement Learning

60M training steps



— D4PG (Pixels)  
— D4PG (State)

60M training steps

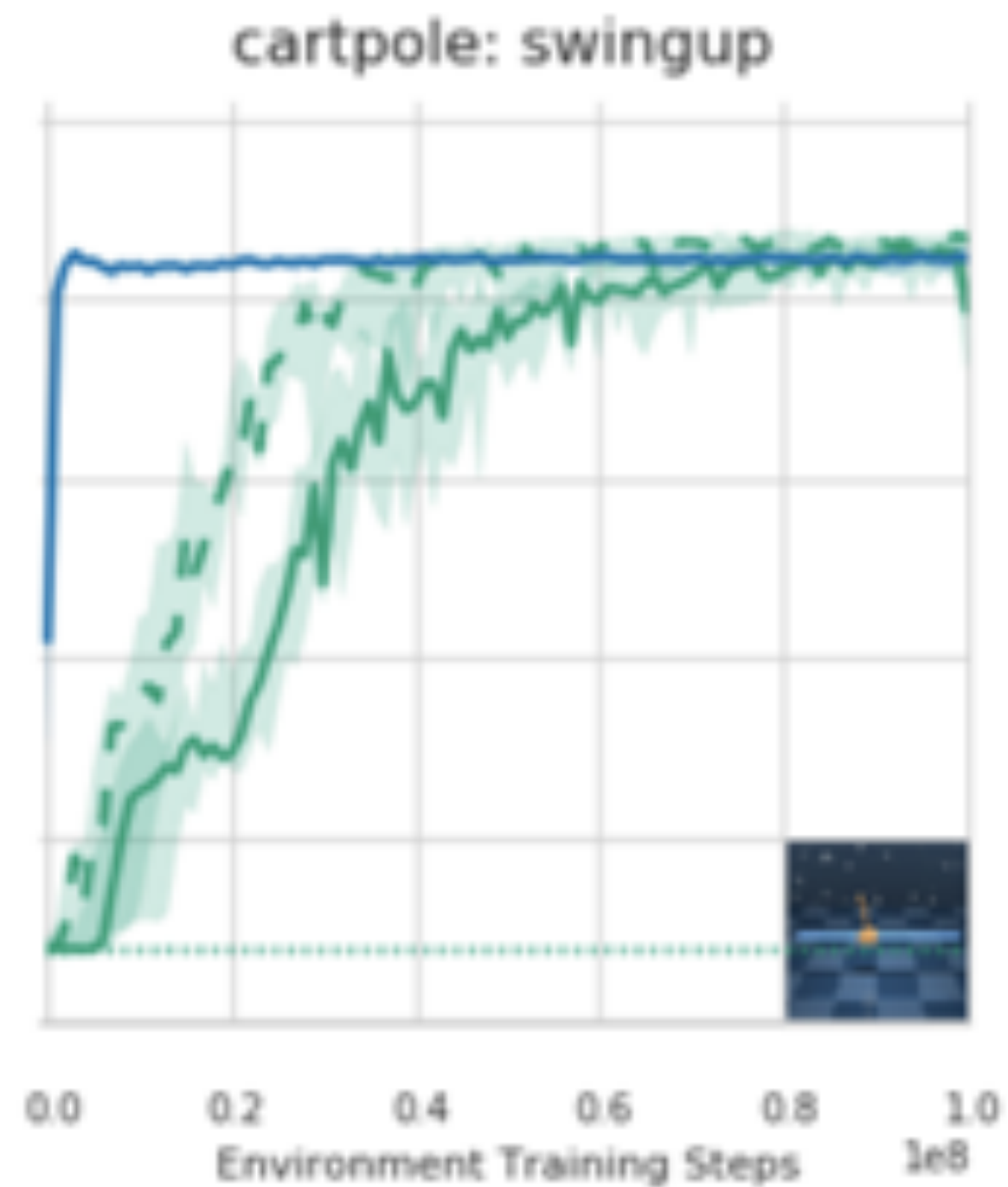


..... D4PG (Pixels, actor)  
- - D4PG (Pixels, critic)

[Tassa et al., 2018] Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D.D.L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A. and Lillicrap, T [DeepMind Control Suite](#), arxiv:1801.00690, 2018.

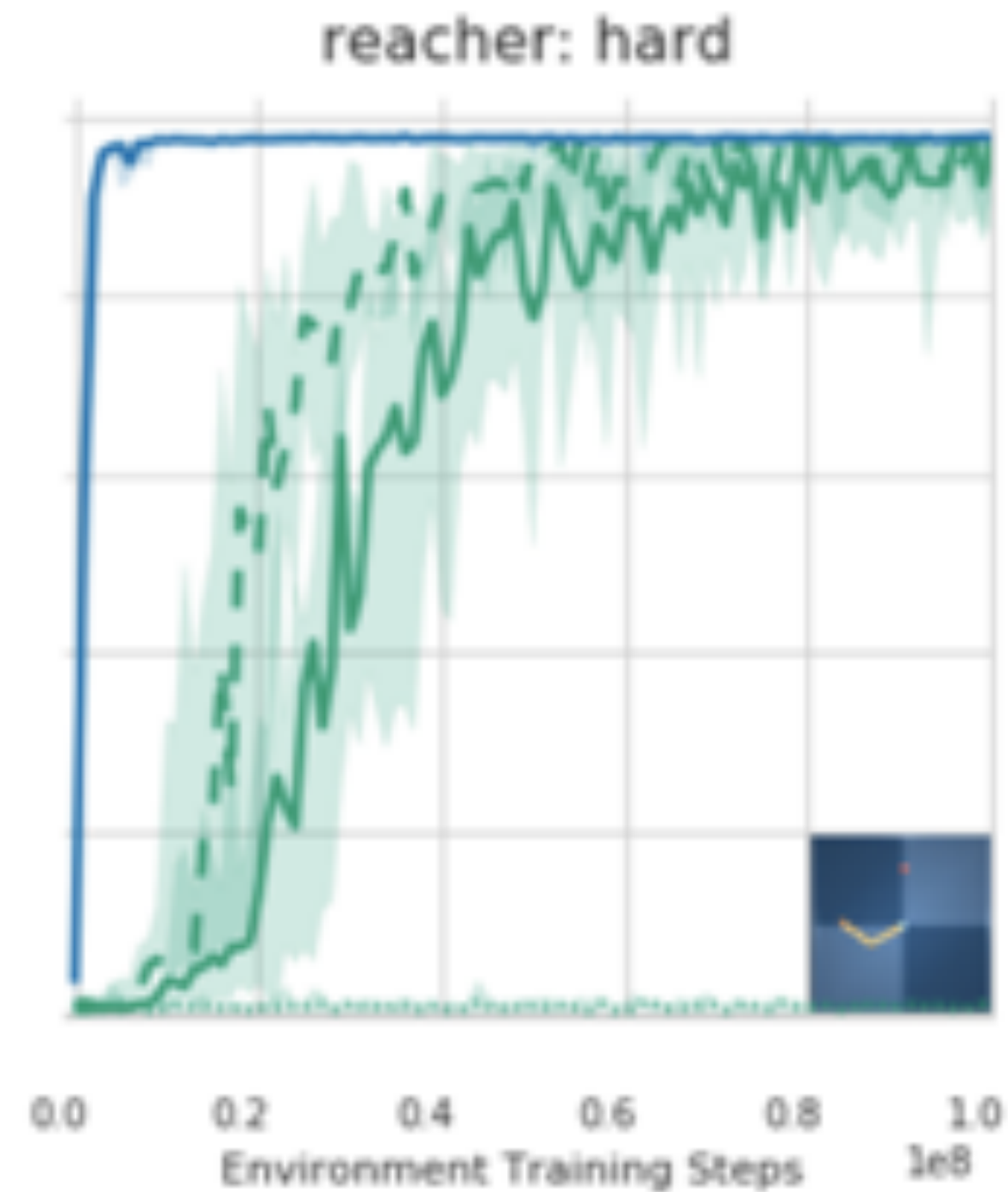
# Data-Efficiency in Reinforcement Learning

60M training steps



— D4PG (Pixels)  
— D4PG (State)

60M training steps



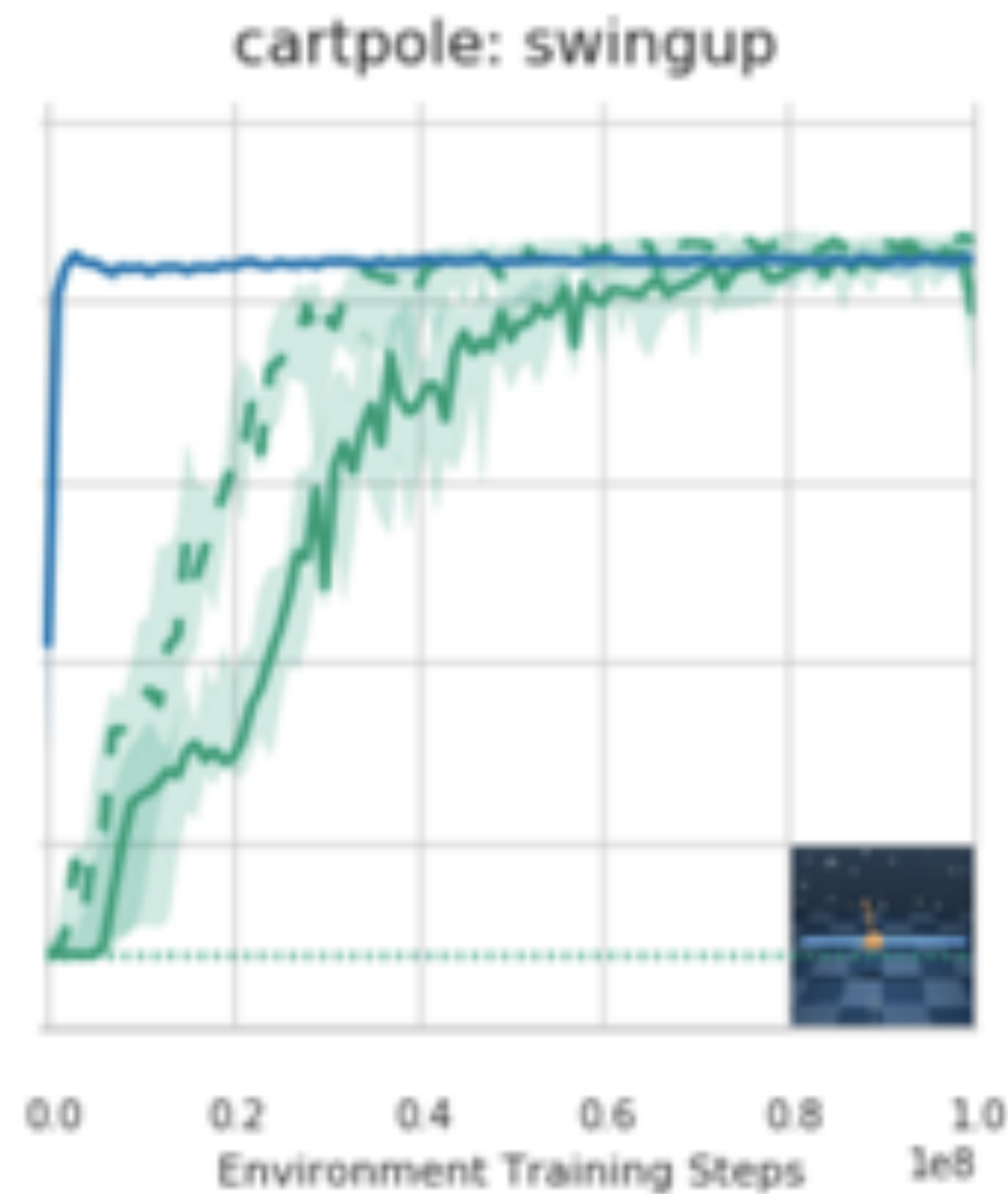
..... D4PG (Pixels, actor)  
- - D4PG (Pixels, critic)

Blue: RL from state,  
Green: RL from pixels

[Tassa et al., 2018] Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D.D.L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A. and Lillicrap, T [DeepMind Control Suite](#), arxiv:1801.00690, 2018.

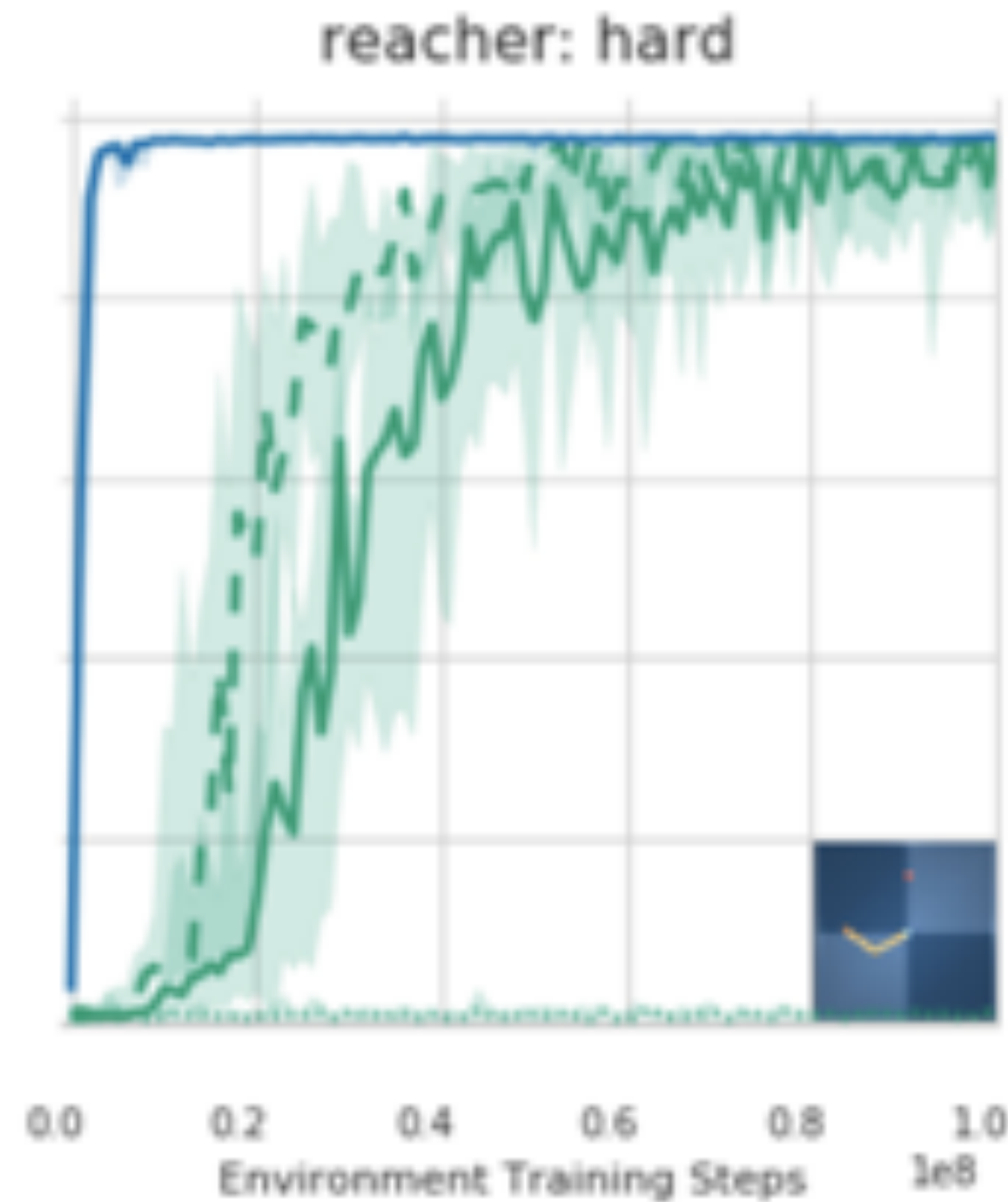
# Data-Efficiency in Reinforcement Learning

60M training steps



— D4PG (Pixels)  
— D4PG (State)

60M training steps



..... D4PG (Pixels, actor)  
- - D4PG (Pixels, critic)

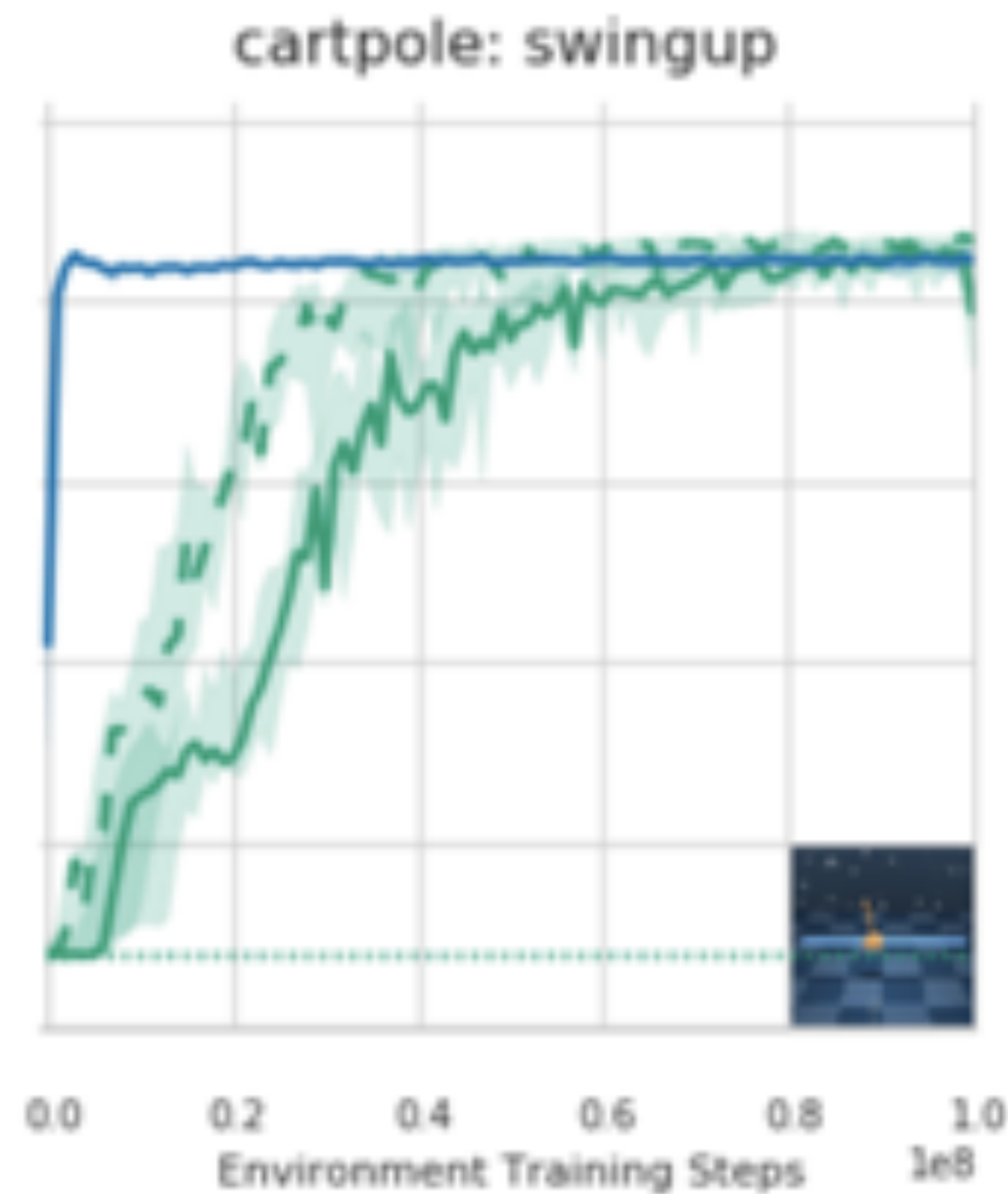
Blue: RL from state,  
Green: RL from pixels

Difference in sample  
complexity is at least >  
50M training steps

[Tassa et al., 2018] Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D.D.L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A. and Lillicrap, T [DeepMind Control Suite](#), arxiv:1801.00690, 2018.

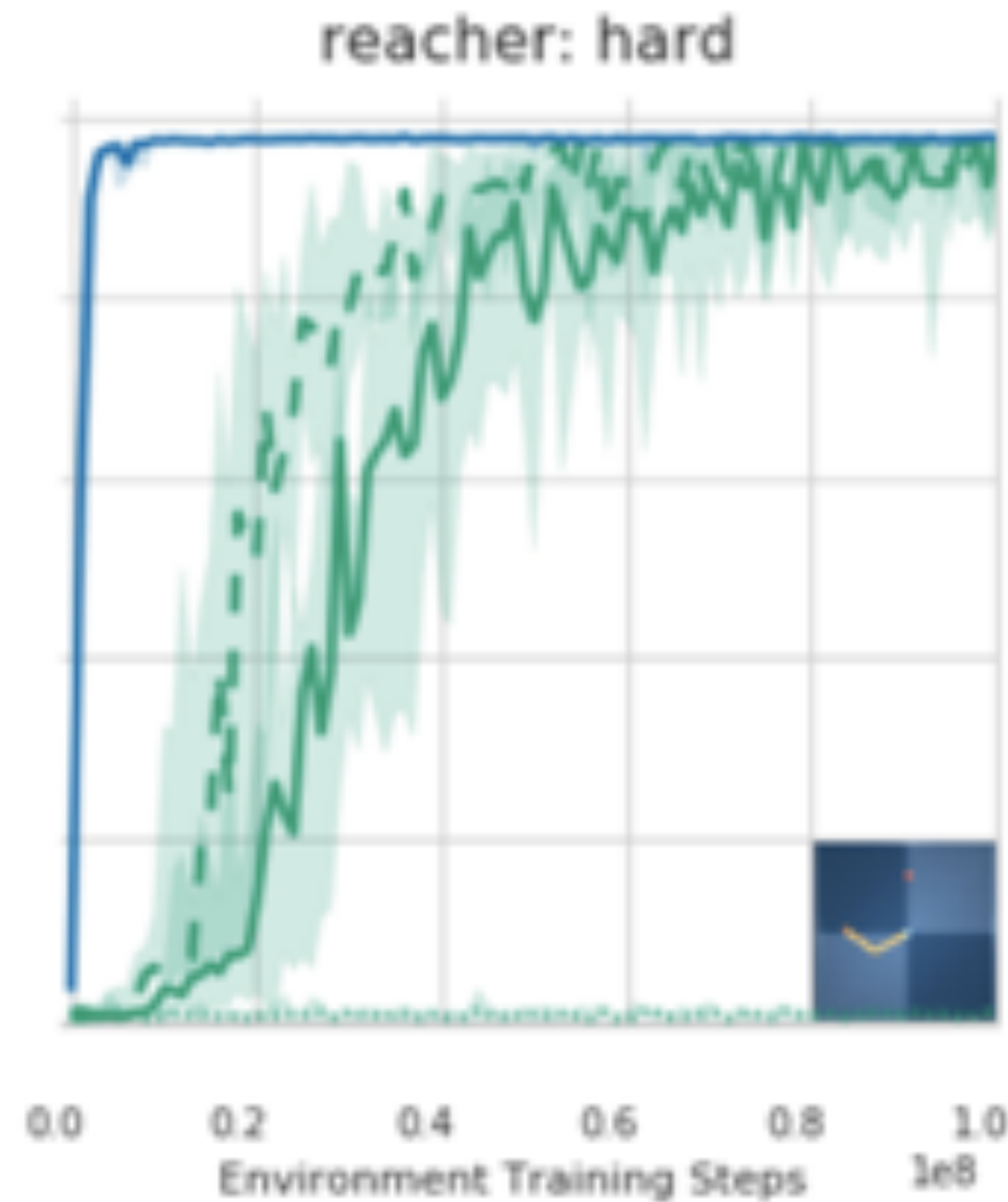
# Data-Efficiency in Reinforcement Learning

60M training steps



— D4PG (Pixels)  
— D4PG (State)

60M training steps



..... D4PG (Pixels, actor)  
- - D4PG (Pixels, critic)

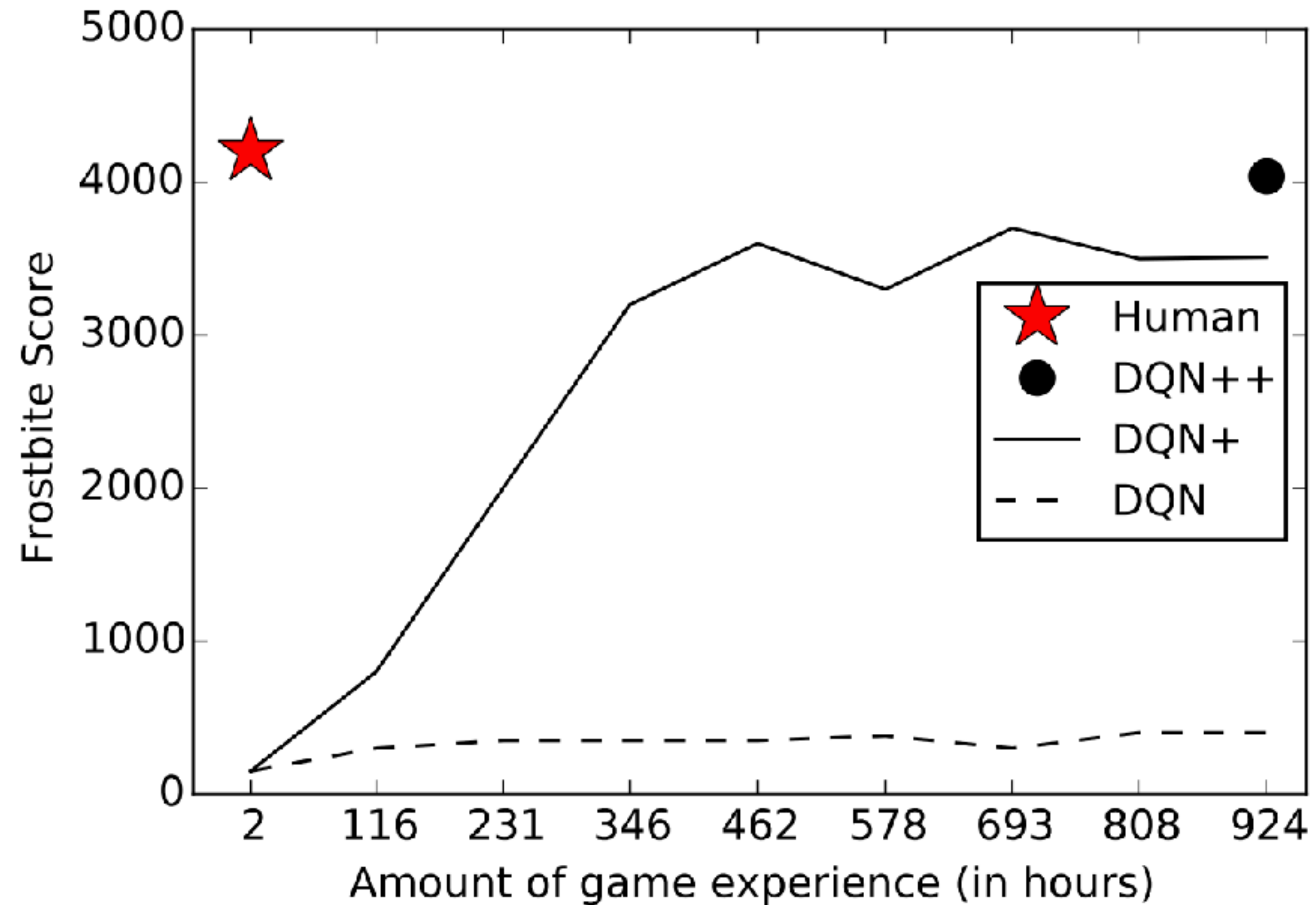
Blue: RL from state,  
Green: RL from pixels

Difference in sample  
complexity is at least >  
50M training steps

Can contrastive learning fix this?

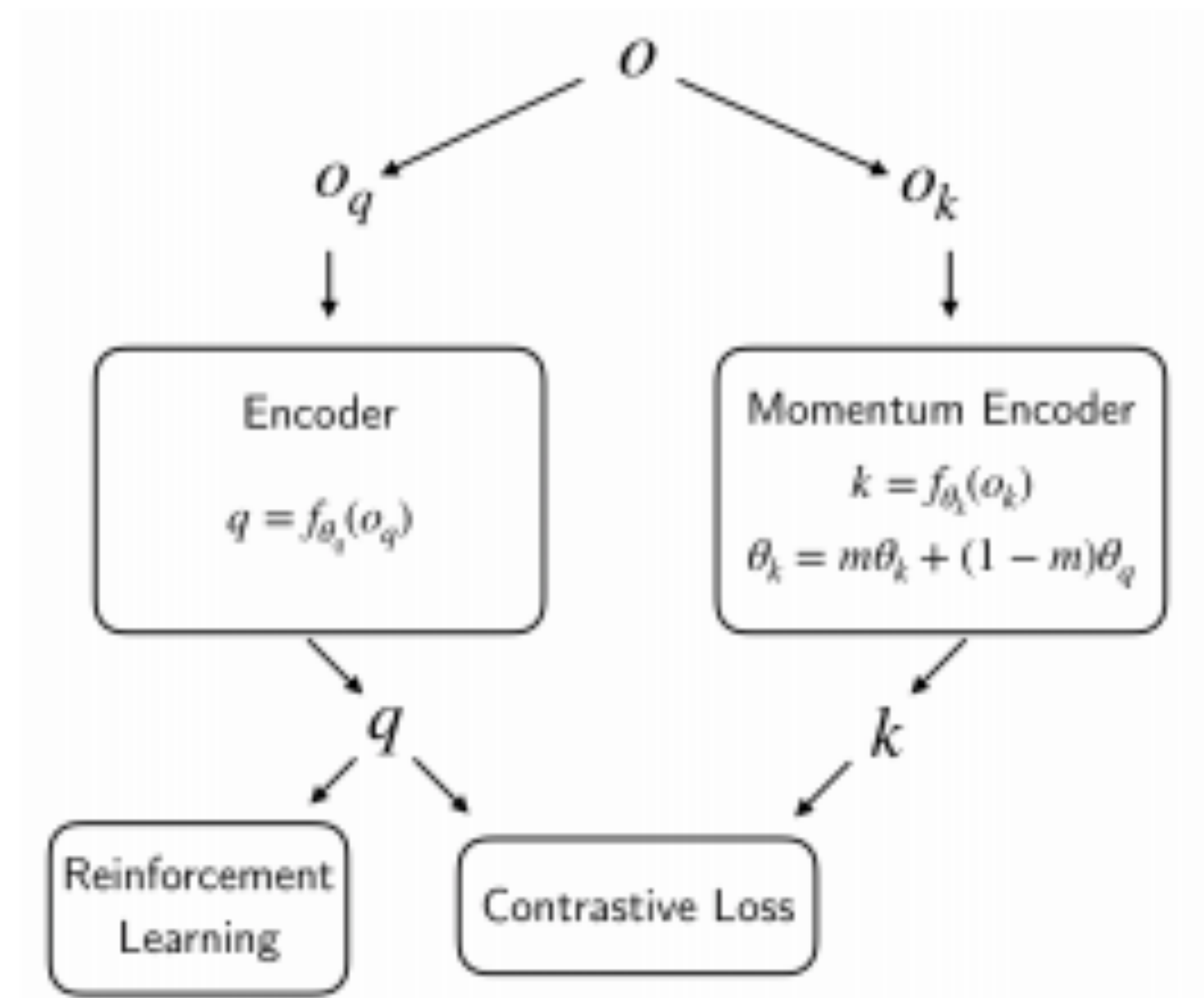
[Tassa et al., 2018] Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D.D.L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A. and Lillicrap, T [DeepMind Control Suite](#), arxiv:1801.00690, 2018.

# Data-Efficiency in Reinforcement Learning

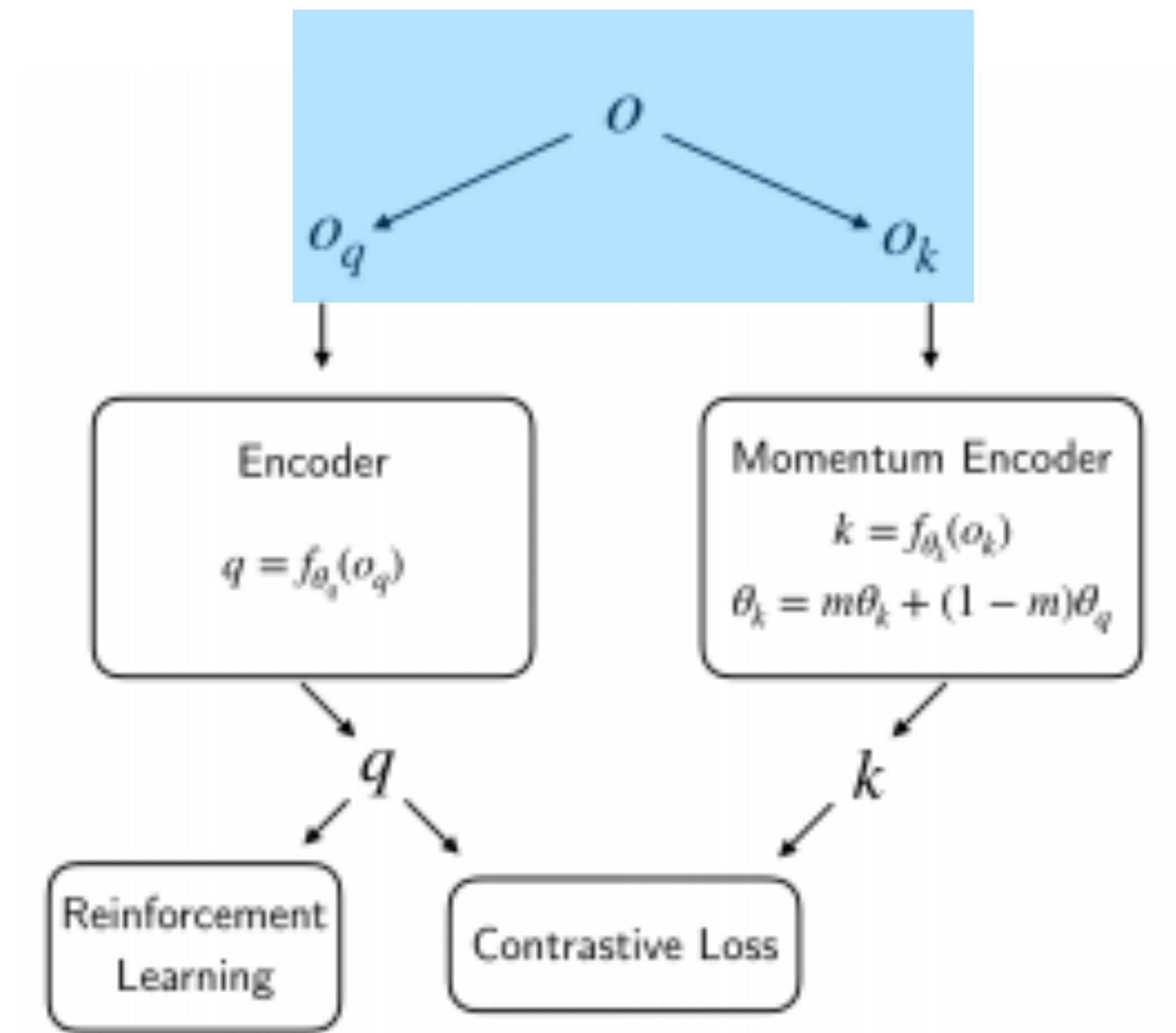


Building machines that can learn and think like people, Lake et al 2016

# Contrastive Learning + Reinforcement Learning

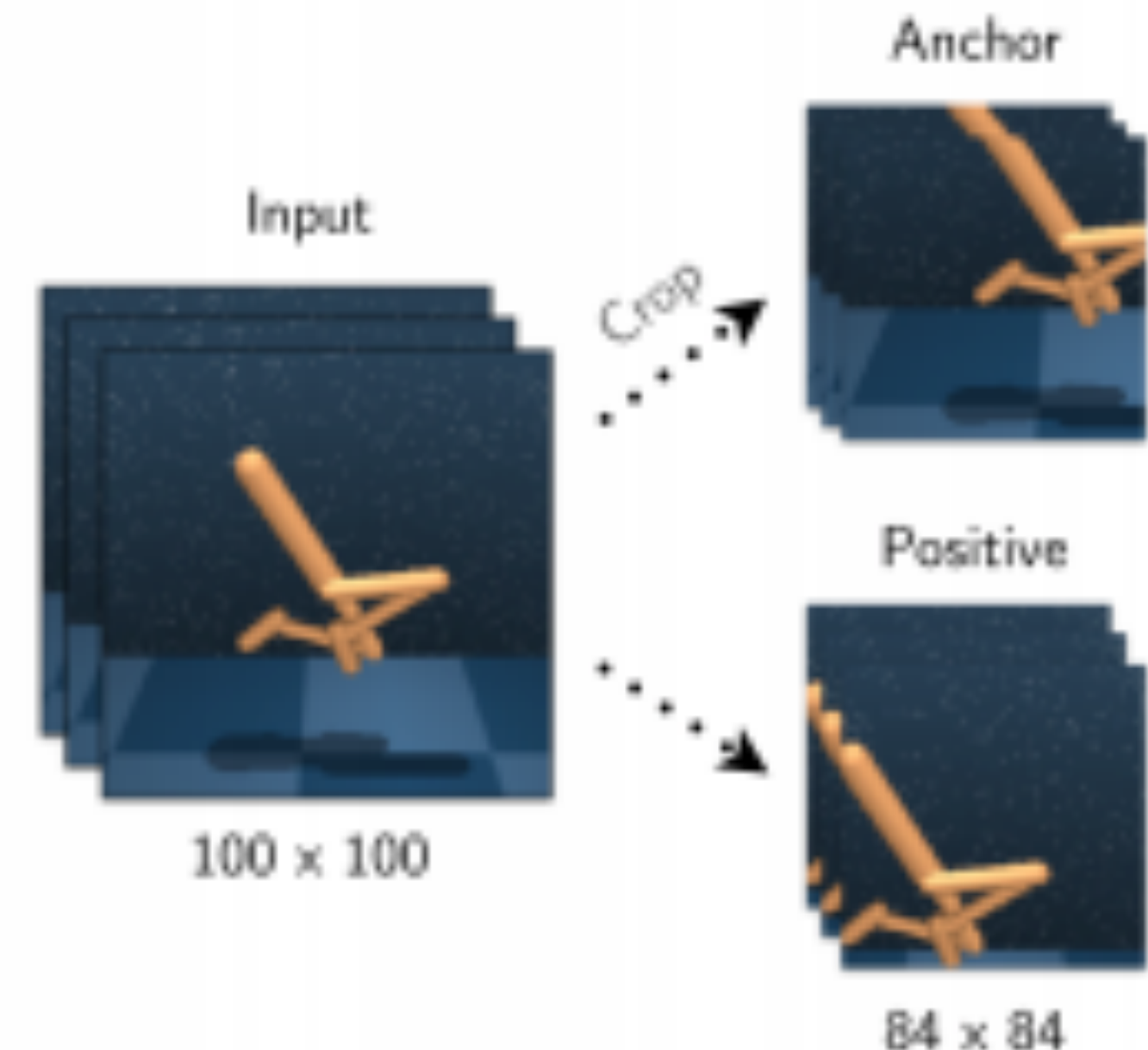
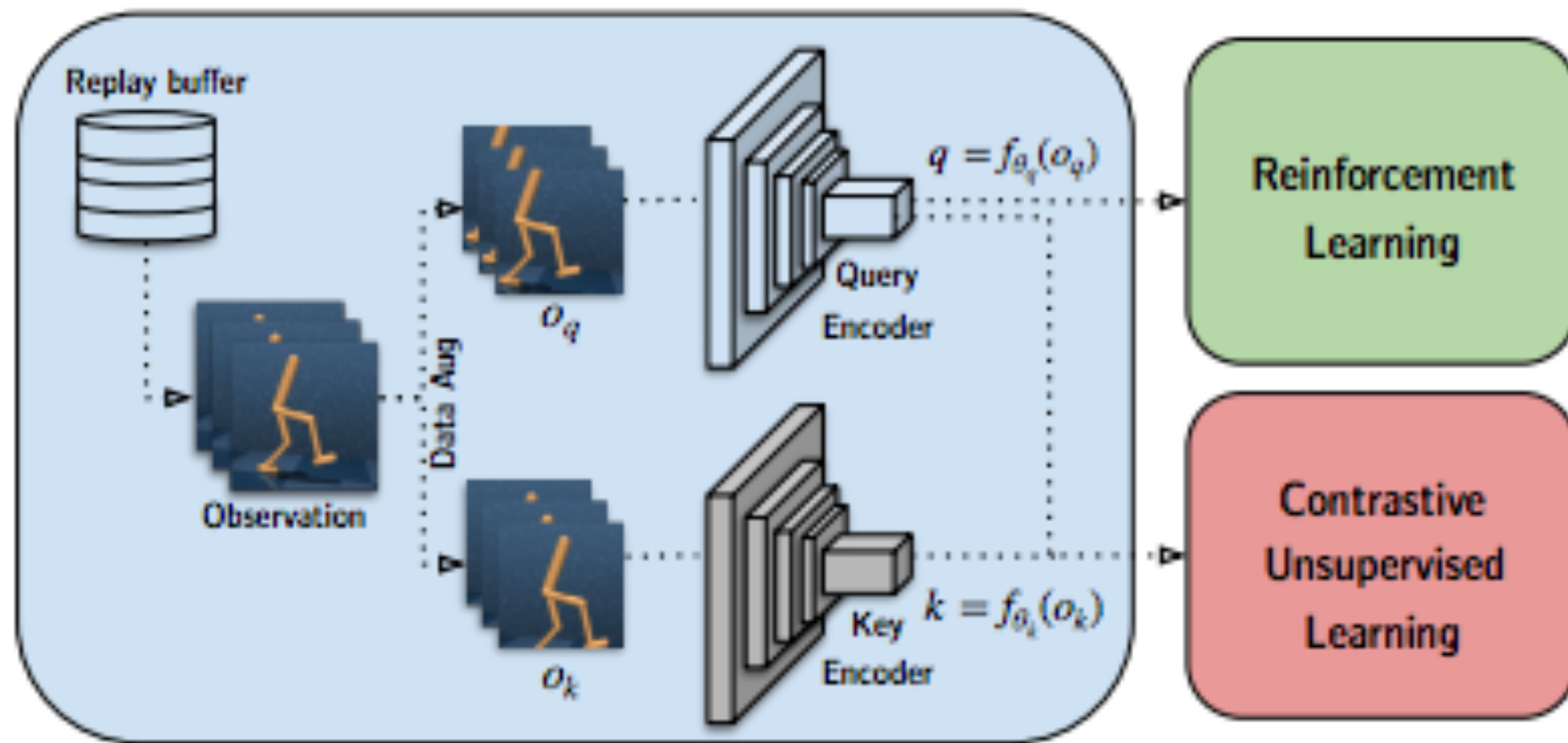


# Contrastive Learning + Reinforcement Learning



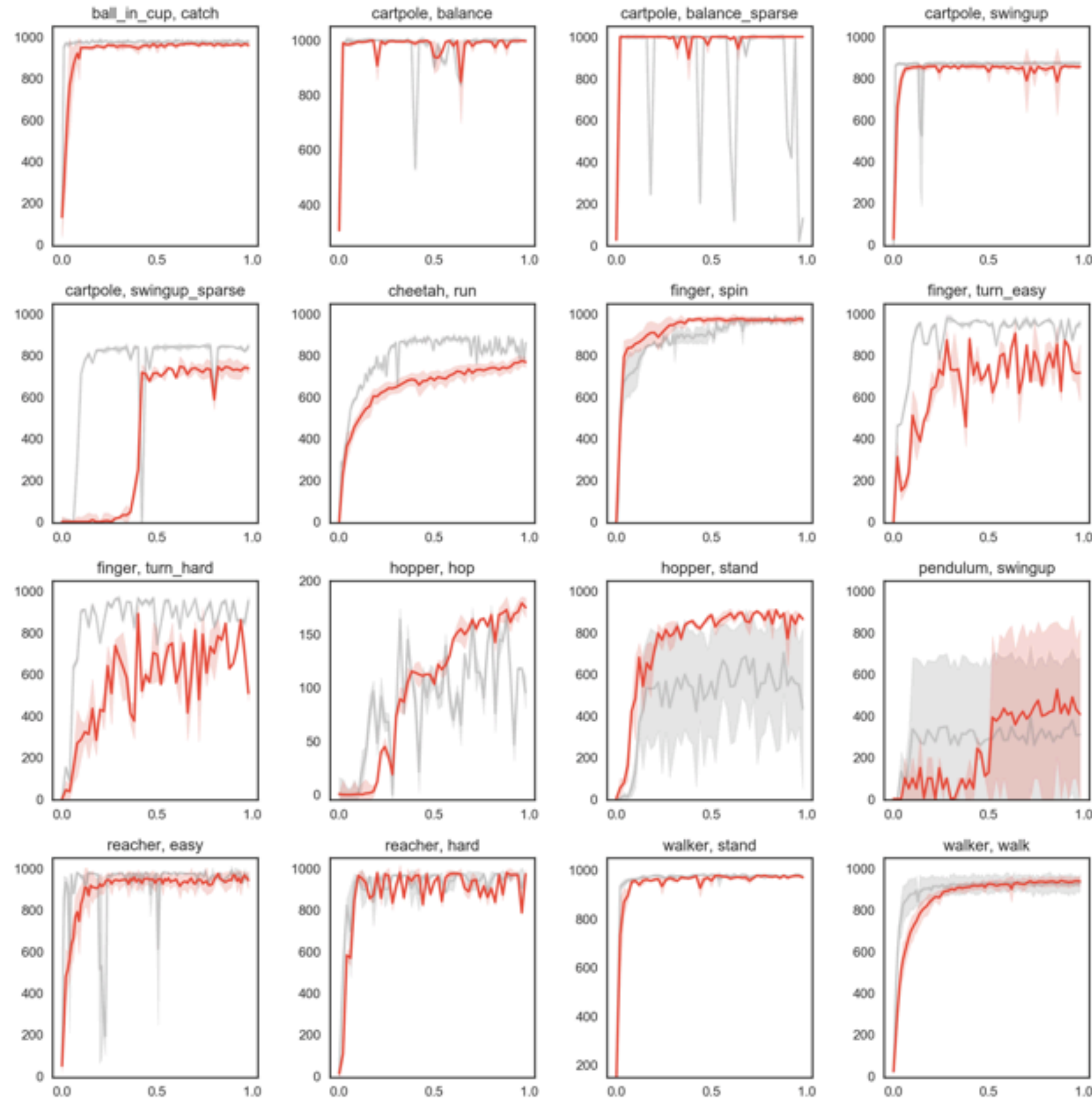


# Contrastive Learning + Reinforcement Learning



1. Stacked Frames instead of single image
2. Temporally Consistent Spatial Random Crop
3. Common practice in Video Recognition

# Contrastive Learning + Reinforcement Learning



GRAY: SAC State

RED: CURL

# Contrastive Learning + Reinforcement Learning

500K STEP SCORES	CURL	PLANET	DREAMER	SAC+AE	SLAC	PIXEL SAC	STATE SAC
FINGER, SPIN	<b>971 ± 18</b>	693 ± 27	343 ± 43	884 ± 128	892 ± 130	509 ± 148	932 ± 32
CARTPOLE, SWINGUP	<b>853 ± 10</b>	794 ± 14	688 ± 207	735 ± 63	-	382 ± 79	870 ± 11
REACHER, EASY	<b>945 ± 27</b>	833 ± 101	480 ± 128	627 ± 58	-	201 ± 94	944 ± 30
CHEETAH, RUN	<b>694 ± 42</b>	608 ± 20	628 ± 32	550 ± 34	617 ± 14	292 ± 31	826 ± 22
WALKER, WALK	<b>925 ± 21</b>	912 ± 35	865 ± 44	847 ± 48	877 ± 54	226 ± 15	935 ± 31
BALL IN CUP, CATCH	<b>956 ± 18</b>	725 ± 309	955 ± 19	794 ± 58	900 ± 181	118 ± 92	984 ± 16
100K STEP SCORES							
FINGER, SPIN	<b>845 ± 42</b>	632 ± 112	141 ± 73	740 ± 64	693 ± 141	403 ± 67	741 ± 54
CARTPOLE, SWINGUP	<b>855 ± 8</b>	498 ± 95	288 ± 167	311 ± 11	-	119 ± 22	869 ± 10
REACHER, EASY	<b>819 ± 72</b>	336 ± 122	134 ± 64	274 ± 14	-	120 ± 23	966 ± 29
CHEETAH, RUN	<b>495 ± 31</b>	294 ± 98	113 ± 25	267 ± 24	296 ± 34	49 ± 14	617 ± 9
WALKER, WALK	<b>715 ± 45</b>	290 ± 111	102 ± 12	394 ± 22	529 ± 47	76 ± 19	884 ± 64
BALL IN CUP, CATCH	<b>942 ± 27</b>	405 ± 375	298 ± 40	391 ± 82	834 ± 128	131 ± 52	979 ± 14

# Contrastive Learning + Reinforcement Learning

500K STEP SCORES	CURL	PLANET	DREAMER	SAC+AE	SLAC	PIXEL SAC	STATE SAC
FINGER, SPIN	<b>971 ± 18</b>	693 ± 27	343 ± 43	884 ± 128	892 ± 130	509 ± 148	932 ± 32
CARTPOLE, SWINGUP	<b>853 ± 10</b>	794 ± 14	688 ± 207	735 ± 63	-	382 ± 79	870 ± 11
REACHER, EASY	<b>945 ± 27</b>	833 ± 101	480 ± 128	627 ± 58	-	201 ± 94	944 ± 30
CHEETAH, RUN	<b>694 ± 42</b>	608 ± 20	628 ± 32	550 ± 34	617 ± 14	292 ± 31	826 ± 22
WALKER, WALK	<b>925 ± 21</b>	912 ± 35	865 ± 44	847 ± 48	877 ± 54	226 ± 15	935 ± 31
BALL IN CUP, CATCH	<b>956 ± 18</b>	725 ± 309	955 ± 19	794 ± 58	900 ± 181	118 ± 92	984 ± 16
100K STEP SCORES							
FINGER, SPIN	<b>845 ± 42</b>	632 ± 112	141 ± 73	740 ± 64	693 ± 141	403 ± 67	741 ± 54
CARTPOLE, SWINGUP	<b>855 ± 8</b>	498 ± 95	288 ± 167	311 ± 11	-	119 ± 22	869 ± 10
REACHER, EASY	<b>819 ± 72</b>	336 ± 122	134 ± 64	274 ± 14	-	120 ± 23	966 ± 29
CHEETAH, RUN	<b>495 ± 31</b>	294 ± 98	113 ± 25	267 ± 24	296 ± 34	49 ± 14	617 ± 9
WALKER, WALK	<b>715 ± 45</b>	290 ± 111	102 ± 12	394 ± 22	529 ± 47	76 ± 19	884 ± 64
BALL IN CUP, CATCH	<b>942 ± 27</b>	405 ± 375	298 ± 40	391 ± 82	834 ± 128	131 ± 52	979 ± 14

Planet, Dreamer, SLAC, SAC+AE are strong but relatively complex baselines.

All of them reconstruct pixels (either autoencoder or future frames). Contrastive just learns what's necessary and high-level.

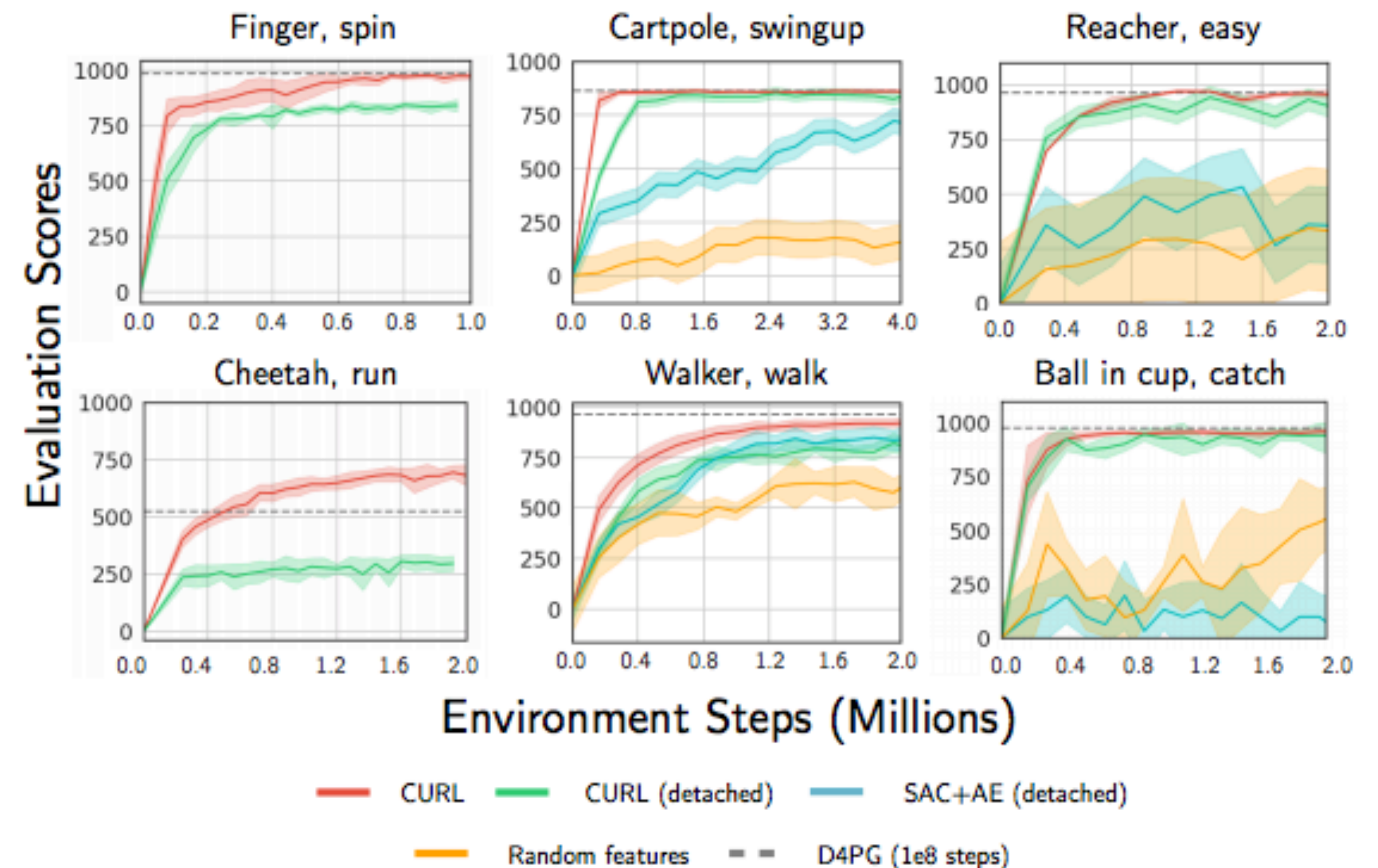
CURL beats Planet, Dreamer, SLAC (model-based methods) despite being model-free and minimal.

# Decoupled Representation Learning

Just learn the encoder with contrastive, and do not backpropogate the RL loss to the encoder.  
Towards a Feature Layer Moment (promise, not yet happened) for RL from pixels.



The Le Cake (Figure from David Ha)



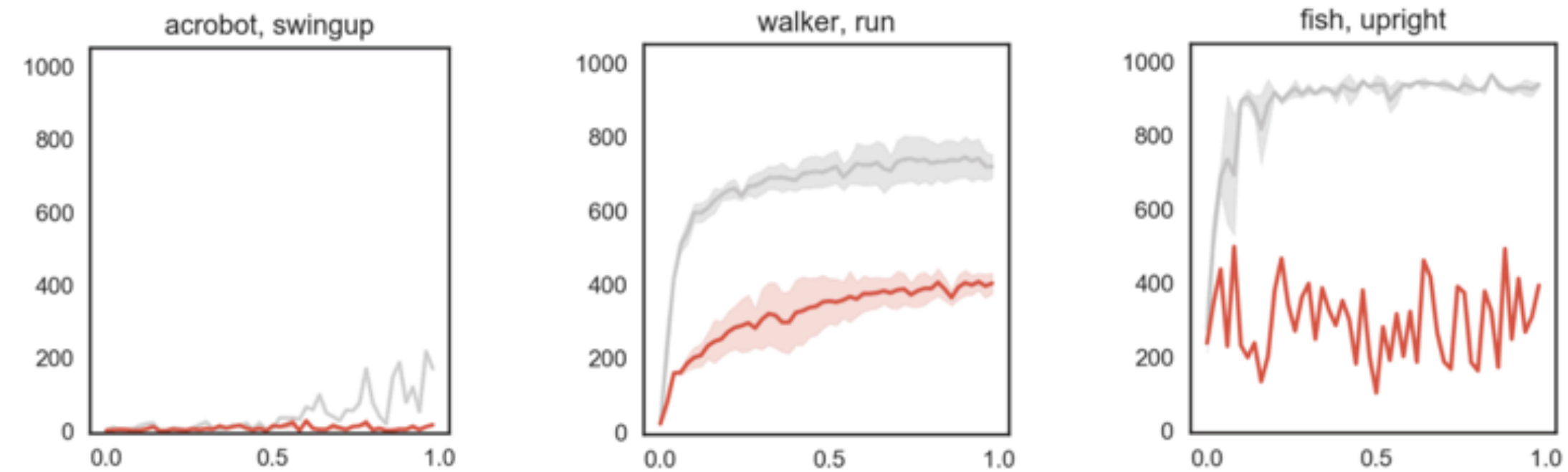
On 4/6 tasks, it is almost as good. On remaining 2, it is not (ex Cheetah Run, Walk Walk).

# Failure Cases

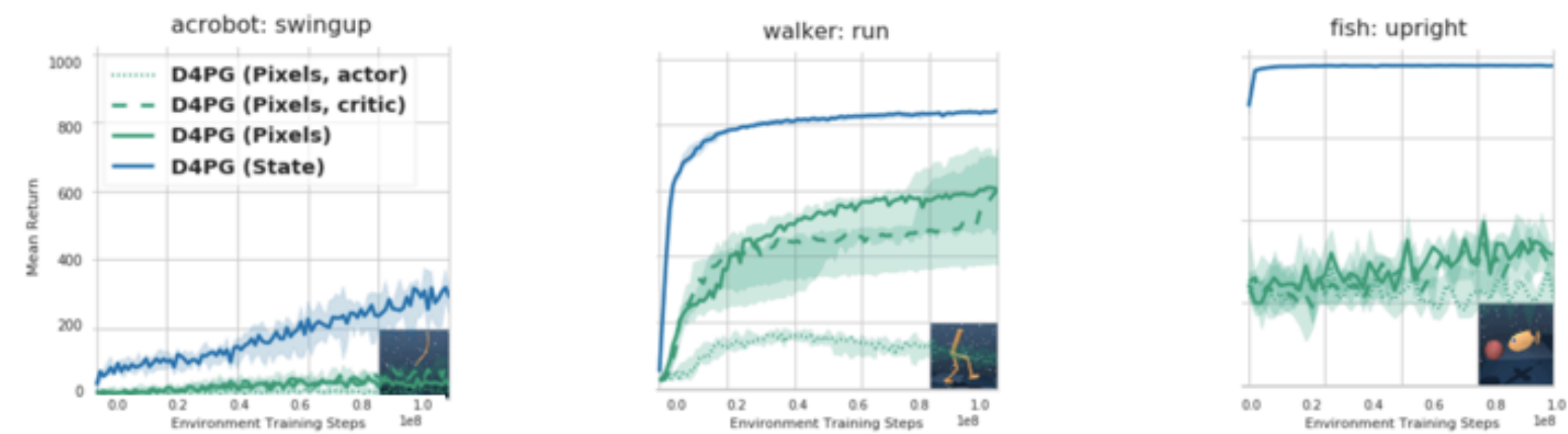
GRAY: SAC State

RED: CURL

Environment training steps 1 = 1M

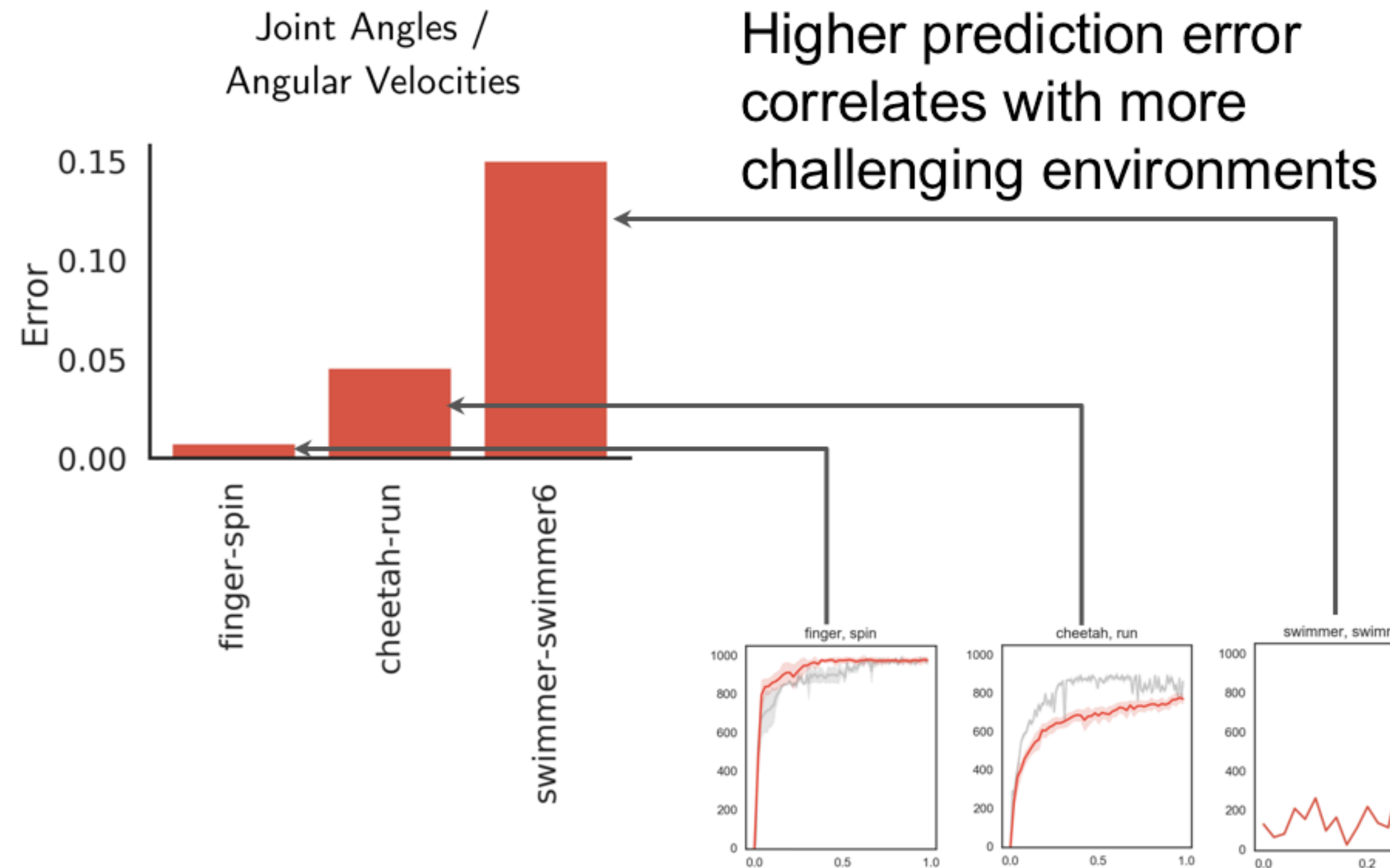


Environment training steps 1 = 100M



Not good on tasks with high frequency dynamics, for example, running, swimming, acrobot, etc.

# Failure Cases



The representation learning struggles on these harder environments.

# Self-Predictive Representations for RL

---

## Data-Efficient Reinforcement Learning with Self-Predictive Representations

---

Max Schwarzer\*  
Mila, Université de Montréal

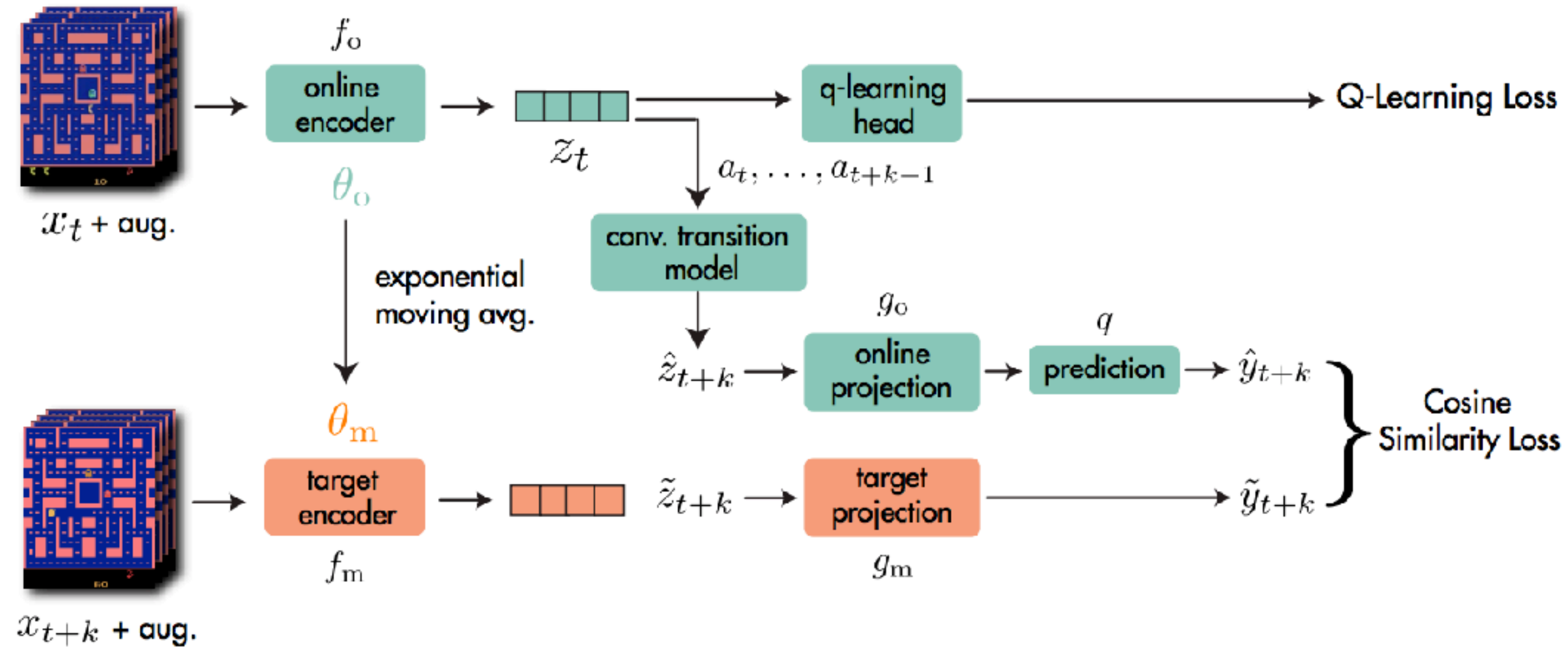
Ankesh Anand\*  
Mila, Université de Montréal  
Microsoft Research

Rishab Goel  
Mila

R Devon Hjelm  
Microsoft Research  
Mila, Université de Montréal

Aaron Courville  
Mila, Université de Montréal  
CIFAR Fellow

Philip Bachman  
Microsoft Research





# Self-Predictive Representations for RL

---

## Data-Efficient Reinforcement Learning with Self-Predictive Representations

---

Max Schwarzer\*  
Mila, Université de Montréal

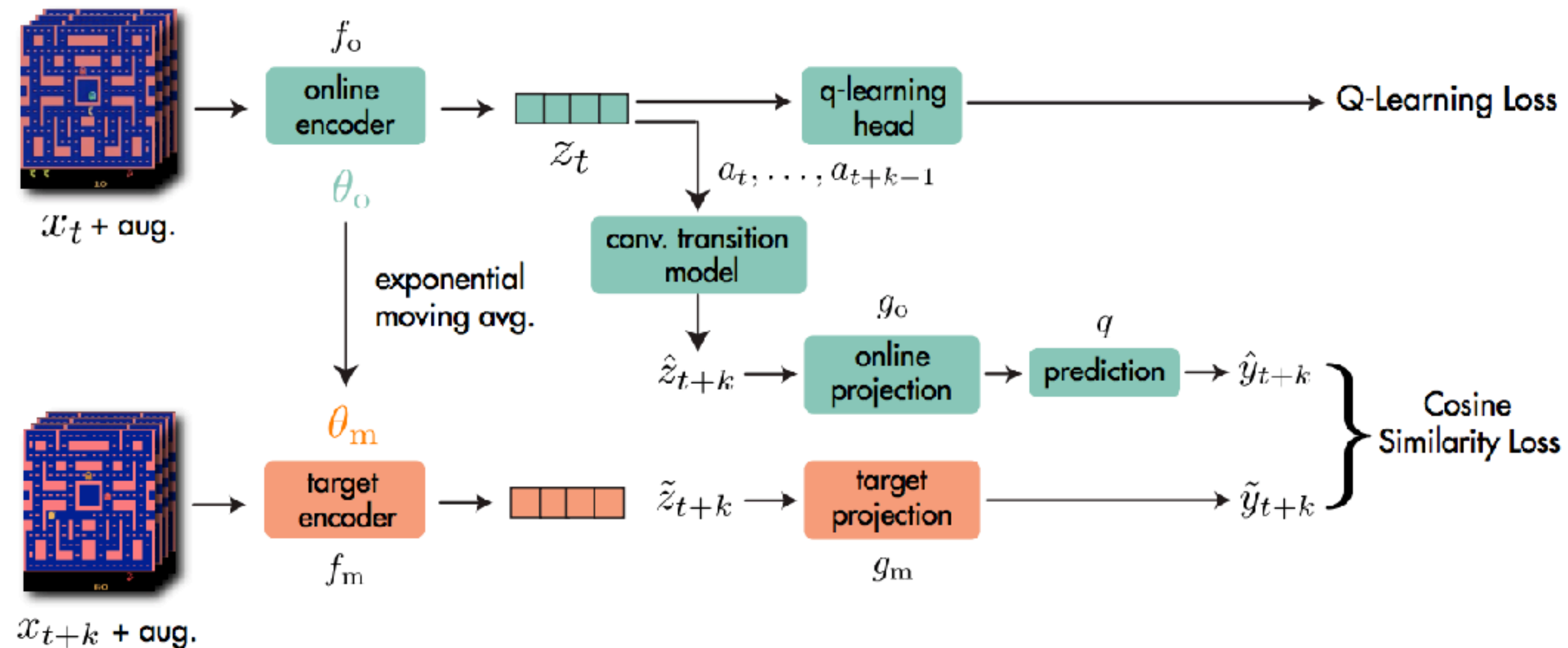
Ankesh Anand\*  
Mila, Université de Montréal  
Microsoft Research

Rishab Goel  
Mila

R Devon Hjelm  
Microsoft Research  
Mila, Université de Montréal

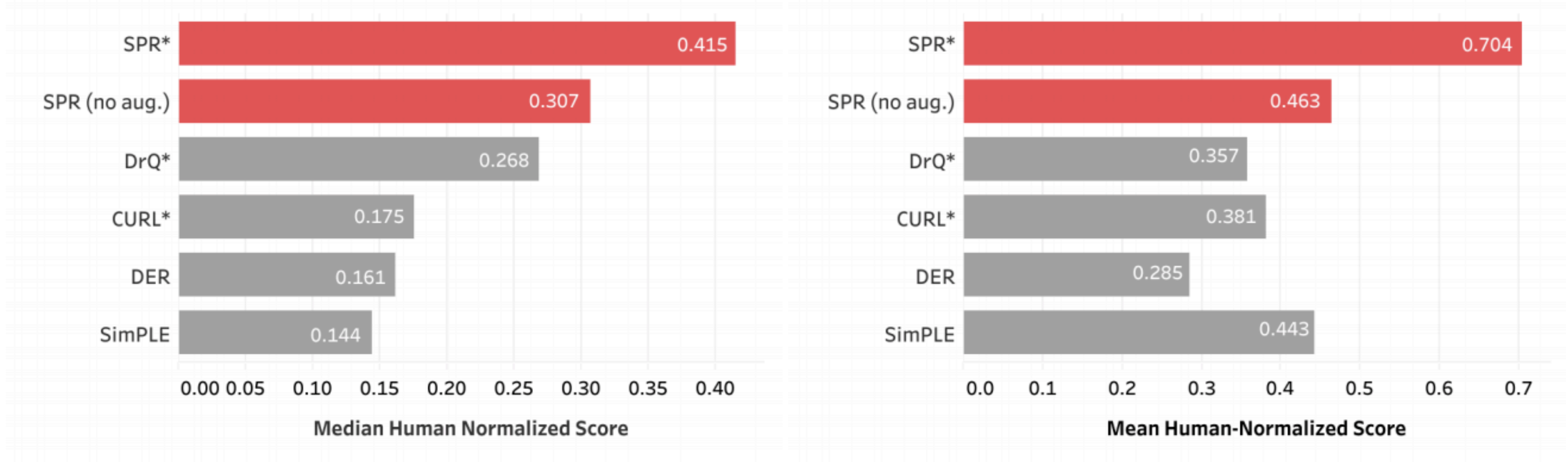
Aaron Courville  
Mila, Université de Montréal  
CIFAR Fellow

Philip Bachman  
Microsoft Research



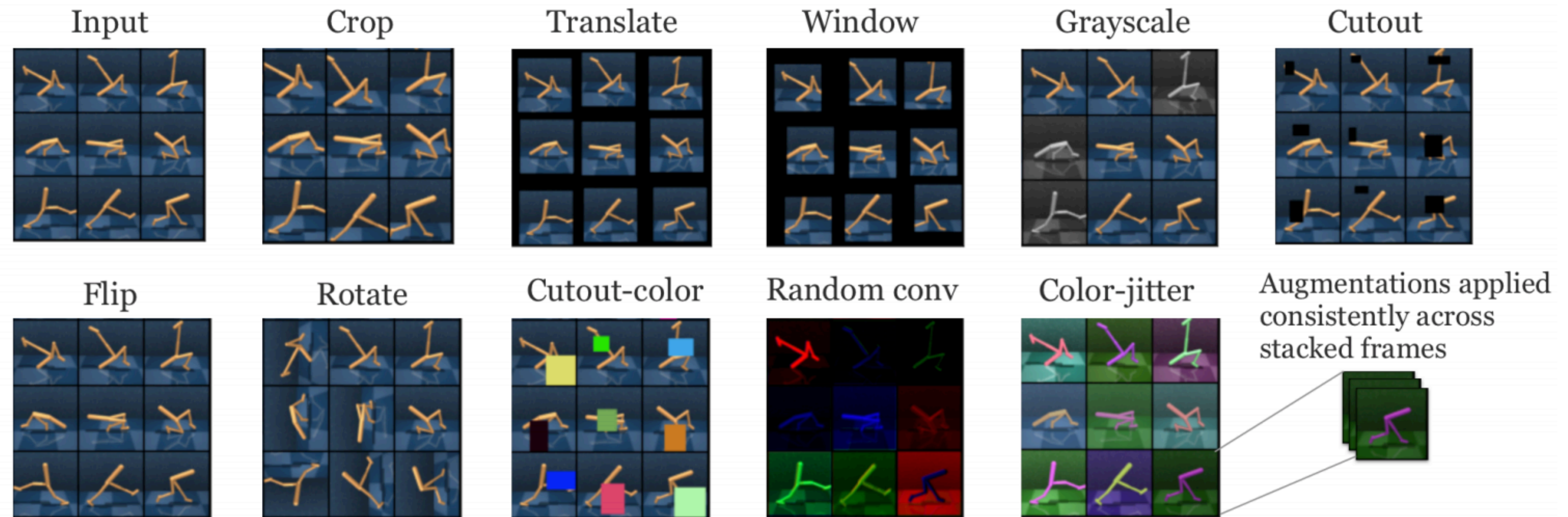
Temporal information helps. Predict the future from past. (CPC-like, but executed in BYOL fashion)

# Self-Predictive Representations for RL



On an average, contrastive learning (or like methods) help get to 70% data-efficiency as a human.

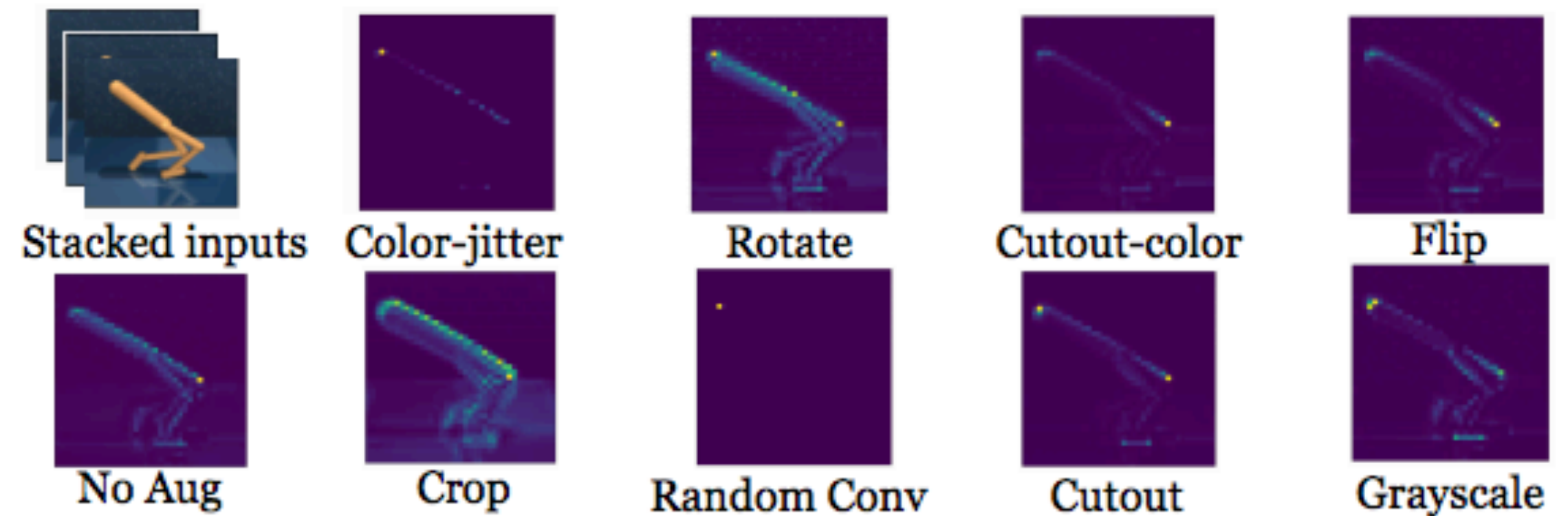
# Reinforcement Learning with Augmented Data



# Reinforcement Learning with Augmented Data

Crop	920	849	635	855	797	650
Grayscale	856	175	349	187	214	231
Rotate	604	403	268	293	392	394
Cutout	722	206	376	215	31	284
Color-jitter	831	227	420	407	194	265
Flip	828	210	391	244	264	223
	Crop	Gray scale	Rotate	Cutout	Color- jitter	Flip

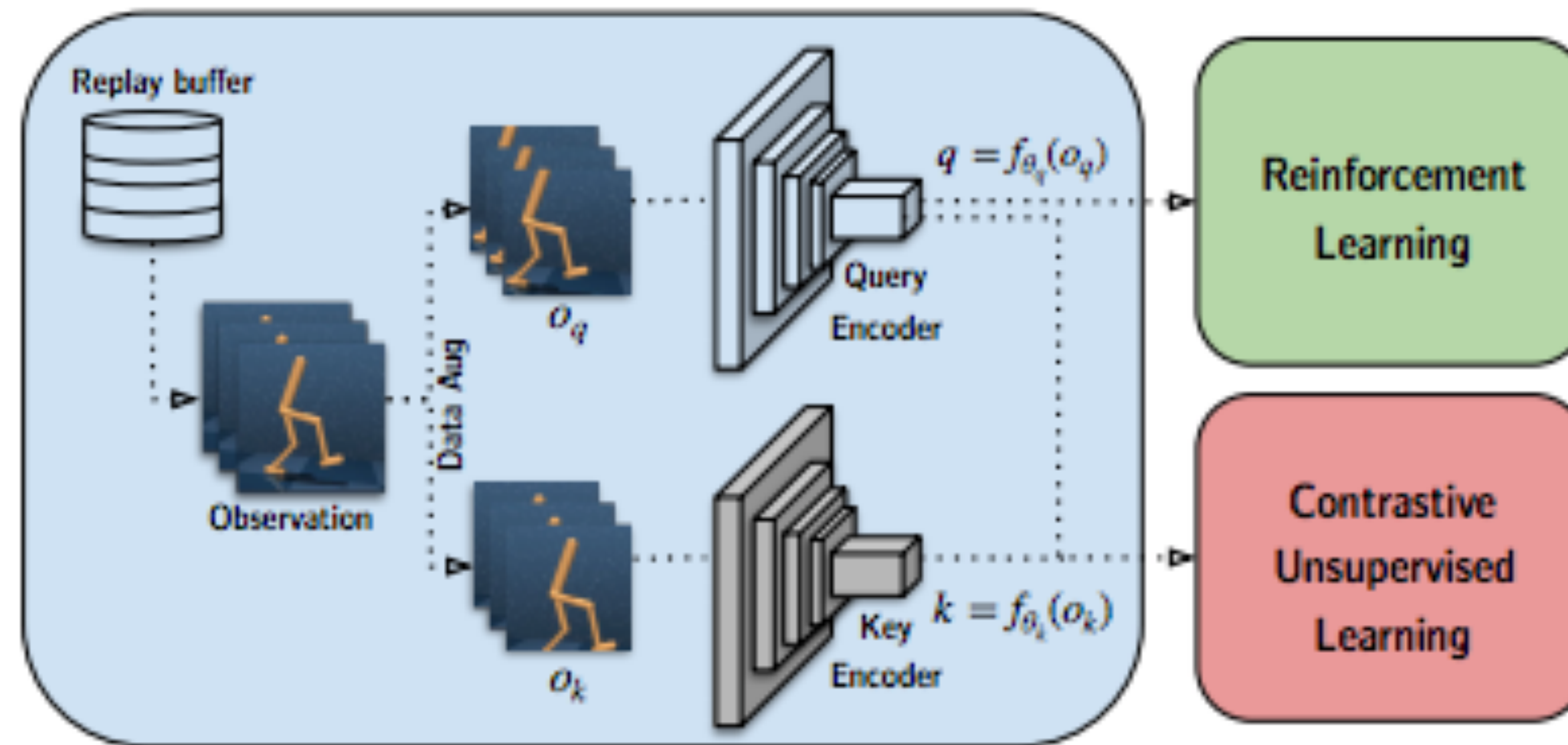
(a) Scores on DMControl500k for Walker, walk.



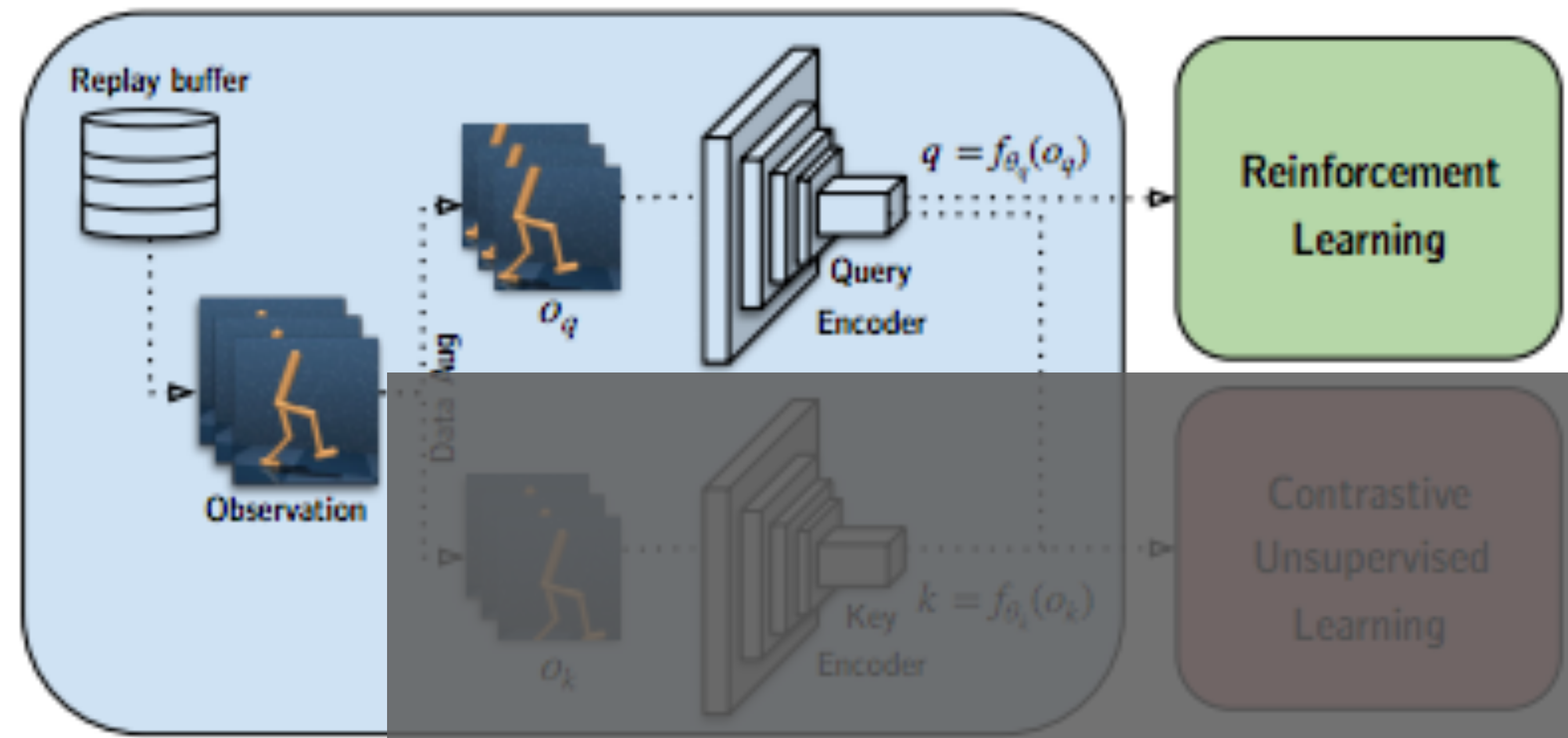
(b) Spatial attention map for Walker, walk.

No surprises, random crops helps. And the network focuses on the right saliency regions.

# Reinforcement Learning with Augmented Data



# Reinforcement Learning with Augmented Data



Reinforcement Learning with Augmented Data, Laskin, Lee, Stooke, Pinto, Abbeel, Srinivas, NeurIPS 2020  
Image Augmentation is all you need; Kostrikov, Yarats, Fergus; ICLR 2021

# Reinforcement Learning with Augmented Data

500K STEP SCORES	RAD	CURL	PLANET	DREAMER	SAC+AE	SLACv1	PIXEL SAC	STATE SAC
FINGER, SPIN	<b>947</b> ± 101	926 ± 45	561 ± 284	796 ± 183	884 ± 128	673 ± 92	192 ± 166	923 ± 211
CARTPOLE, SWING	<b>863</b> ± 9	845 ± 45	475 ± 71	762 ± 27	735 ± 63	-	419 ± 40	848 ± 15
REACHER, EASY	<b>955</b> ± 71	929 ± 44	210 ± 44	793 ± 164	627 ± 58	-	145 ± 30	923 ± 24
CHEETAH, RUN	<b>728</b> ± 71	518 ± 28	305 ± 131	570 ± 253	550 ± 34	640 ± 19	197 ± 15	795 ± 30
WALKER, WALK	<b>918</b> ± 16	902 ± 43	351 ± 58	897 ± 49	847 ± 48	842 ± 51	42 ± 12	948 ± 54
CUP, CATCH	<b>974</b> ± 12	959 ± 27	460 ± 380	879 ± 87	794 ± 58	852 ± 71	312 ± 63	974 ± 33
100K STEP SCORES								
FINGER, SPIN	<b>856</b> ± 73	767 ± 56	136 ± 216	341 ± 70	740 ± 64	693 ± 141	224 ± 101	811 ± 46
CARTPOLE, SWING	<b>828</b> ± 27	582 ± 146	297 ± 39	326 ± 27	311 ± 11	-	200 ± 72	835 ± 22
REACHER, EASY	<b>826</b> ± 219	538 ± 233	20 ± 50	314 ± 155	274 ± 14	-	136 ± 15	746 ± 25
CHEETAH, RUN	<b>447</b> ± 88	299 ± 48	138 ± 88	235 ± 137	267 ± 24	319 ± 56	130 ± 12	616 ± 18
WALKER, WALK	<b>504</b> ± 191	403 ± 24	224 ± 48	277 ± 12	394 ± 22	361 ± 73	127 ± 24	891 ± 82
CUP, CATCH	<b>840</b> ± 179	769 ± 43	0 ± 0	246 ± 174	391 ± 82	512 ± 110	97 ± 27	746 ± 91

Reinforcement Learning with Augmented Data, Laskin, Lee, Stooke, Pinto, Abbeel, Srinivas, NeurIPS 2020  
Image Augmentation is all you need; Kostrikov, Yarats, Fergus; ICLR 2021

# Reinforcement Learning with Augmented Data

500K STEP SCORES	RAD	CURL	PLANET	DREAMER	SAC+AE	SLACv1	PIXEL SAC	STATE SAC
FINGER, SPIN	<b>947</b> ± 101	926 ± 45	561 ± 284	796 ± 183	884 ± 128	673 ± 92	192 ± 166	923 ± 211
CARTPOLE, SWING	<b>863</b> ± 9	845 ± 45	475 ± 71	762 ± 27	735 ± 63	-	419 ± 40	848 ± 15
REACHER, EASY	<b>955</b> ± 71	929 ± 44	210 ± 44	793 ± 164	627 ± 58	-	145 ± 30	923 ± 24
CHEETAH, RUN	<b>728</b> ± 71	518 ± 28	305 ± 131	570 ± 253	550 ± 34	640 ± 19	197 ± 15	795 ± 30
WALKER, WALK	<b>918</b> ± 16	902 ± 43	351 ± 58	897 ± 49	847 ± 48	842 ± 51	42 ± 12	948 ± 54
CUP, CATCH	<b>974</b> ± 12	959 ± 27	460 ± 380	879 ± 87	794 ± 58	852 ± 71	312 ± 63	974 ± 33
100K STEP SCORES								
FINGER, SPIN	<b>856</b> ± 73	767 ± 56	136 ± 216	341 ± 70	740 ± 64	693 ± 141	224 ± 101	811 ± 46
CARTPOLE, SWING	<b>828</b> ± 27	582 ± 146	297 ± 39	326 ± 27	311 ± 11	-	200 ± 72	835 ± 22
REACHER, EASY	<b>826</b> ± 219	538 ± 233	20 ± 50	314 ± 155	274 ± 14	-	136 ± 15	746 ± 25
CHEETAH, RUN	<b>447</b> ± 88	299 ± 48	138 ± 88	235 ± 137	267 ± 24	319 ± 56	130 ± 12	616 ± 18
WALKER, WALK	<b>504</b> ± 191	403 ± 24	224 ± 48	277 ± 12	394 ± 22	361 ± 73	127 ± 24	891 ± 82
CUP, CATCH	<b>840</b> ± 179	769 ± 43	0 ± 0	246 ± 174	391 ± 82	512 ± 110	97 ± 27	746 ± 91

Reinforcement Learning with Augmented Data, Laskin, Lee, Stooke, Pinto, Abbeel, **Srinivas**, NeurIPS 2020  
Image Augmentation is all you need; Kostrikov, Yarats, Fergus; ICLR 2021



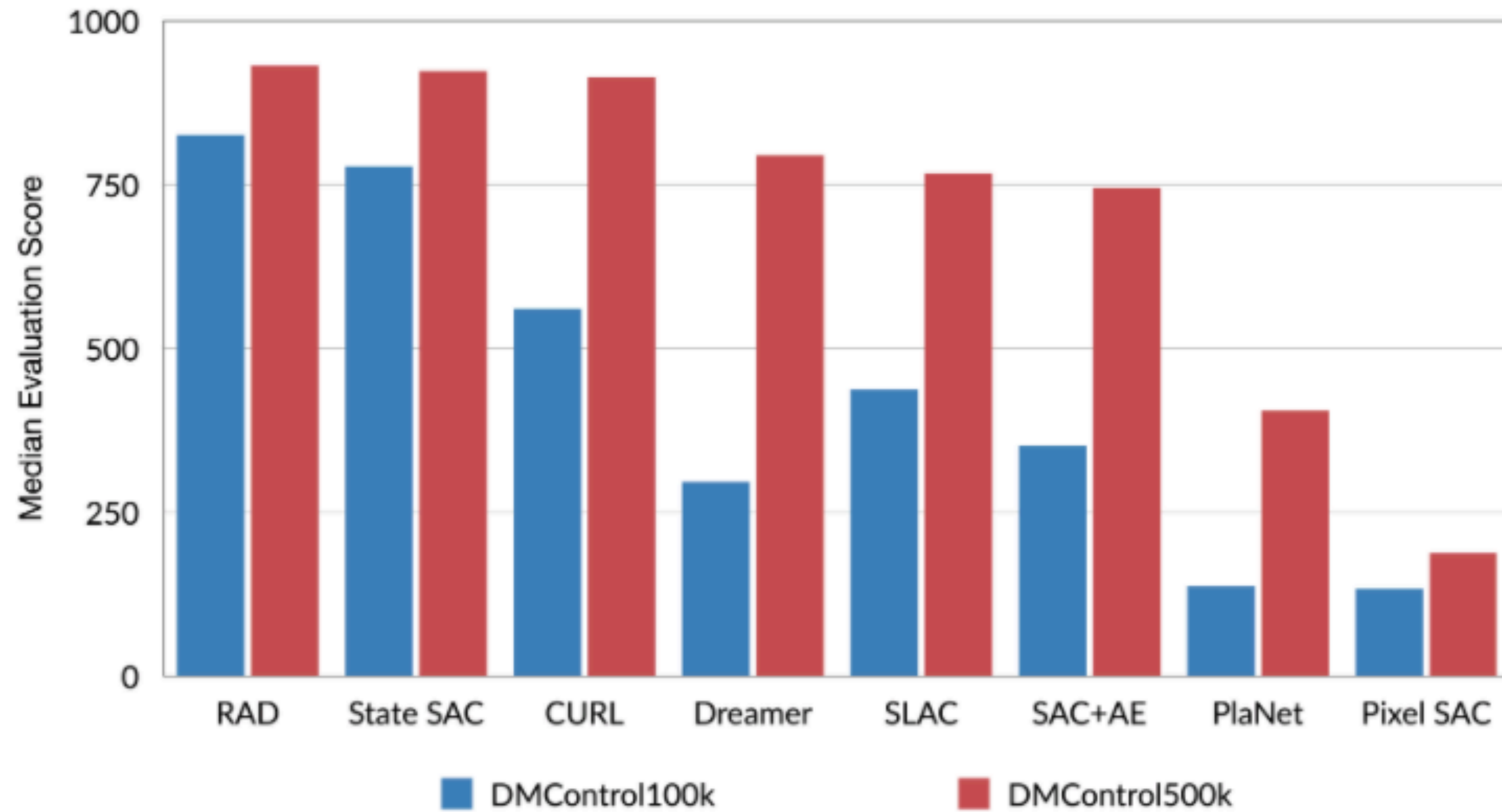
# Reinforcement Learning with Augmented Data

500K STEP SCORES	RAD	CURL	PLANET	DREAMER	SAC+AE	SLACv1	PIXEL SAC	STATE SAC
FINGER, SPIN	<b>947</b> ± 101	926 ± 45	561 ± 284	796 ± 183	884 ± 128	673 ± 92	192 ± 166	923 ± 211
CARTPOLE, SWING	<b>863</b> ± 9	845 ± 45	475 ± 71	762 ± 27	735 ± 63	-	419 ± 40	848 ± 15
REACHER, EASY	<b>955</b> ± 71	929 ± 44	210 ± 44	793 ± 164	627 ± 58	-	145 ± 30	923 ± 24
CHEETAH, RUN	<b>728</b> ± 71	518 ± 28	305 ± 131	570 ± 253	550 ± 34	640 ± 19	197 ± 15	795 ± 30
WALKER, WALK	<b>918</b> ± 16	902 ± 43	351 ± 58	897 ± 49	847 ± 48	842 ± 51	42 ± 12	948 ± 54
CUP, CATCH	<b>974</b> ± 12	959 ± 27	460 ± 380	879 ± 87	794 ± 58	852 ± 71	312 ± 63	974 ± 33
100K STEP SCORES								
FINGER, SPIN	<b>856</b> ± 73	767 ± 56	136 ± 216	341 ± 70	740 ± 64	693 ± 141	224 ± 101	811 ± 46
CARTPOLE, SWING	<b>828</b> ± 27	582 ± 146	297 ± 39	326 ± 27	311 ± 11	-	200 ± 72	835 ± 22
REACHER, EASY	<b>826</b> ± 219	538 ± 233	20 ± 50	314 ± 155	274 ± 14	-	136 ± 15	746 ± 25
CHEETAH, RUN	<b>447</b> ± 88	299 ± 48	138 ± 88	235 ± 137	267 ± 24	319 ± 56	130 ± 12	616 ± 18
WALKER, WALK	<b>504</b> ± 191	403 ± 24	224 ± 48	277 ± 12	394 ± 22	361 ± 73	127 ± 24	891 ± 82
CUP, CATCH	<b>840</b> ± 179	769 ± 43	0 ± 0	246 ± 174	391 ± 82	512 ± 110	97 ± 27	746 ± 91

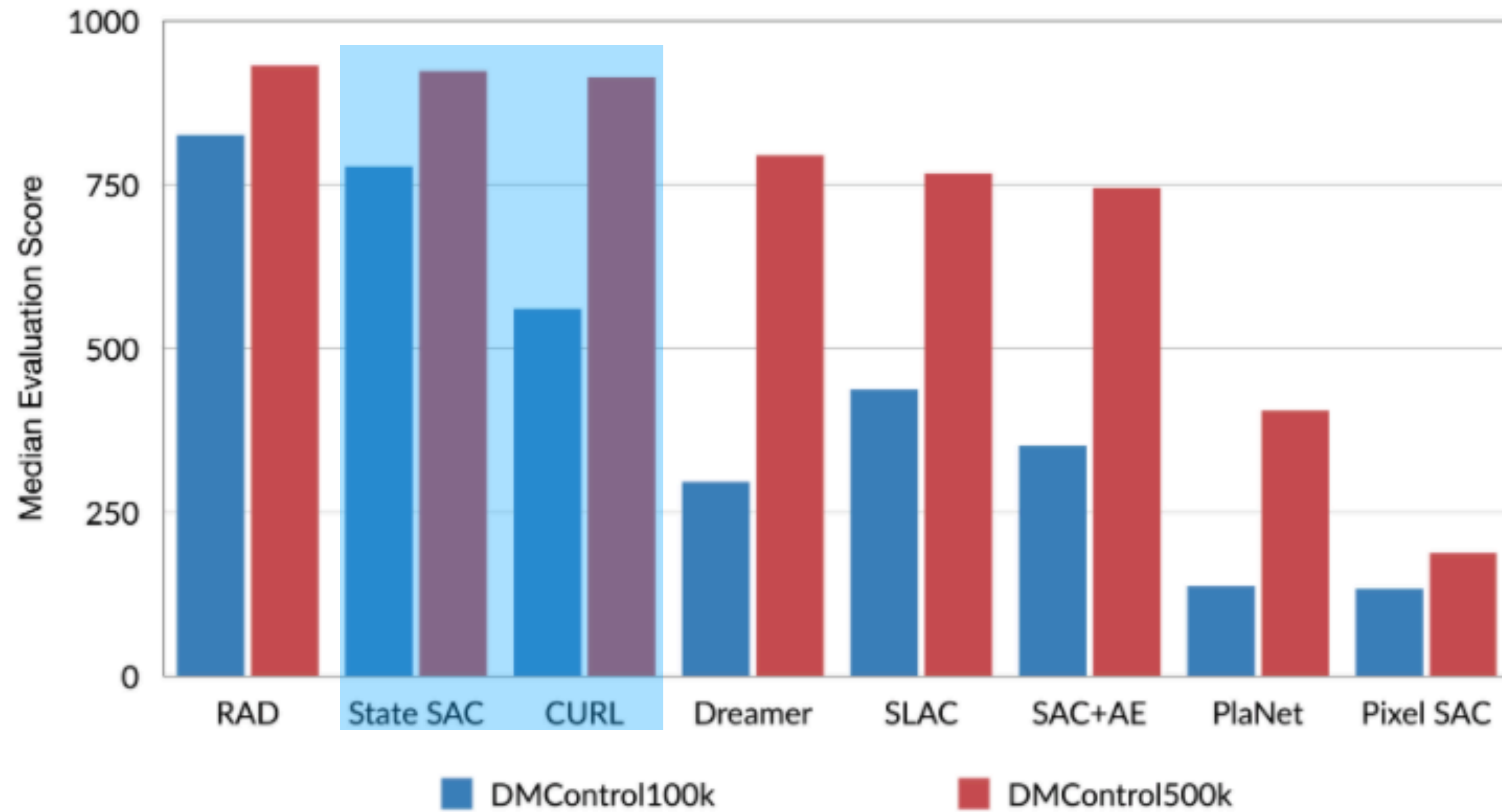
Reinforcement Learning with Augmented Data, Laskin, Lee, Stooke, Pinto, Abbeel, **Srinivas**, NeurIPS 2020

Image Augmentation is all you need; Kostrikov, Yarats, Fergus; ICLR 2021

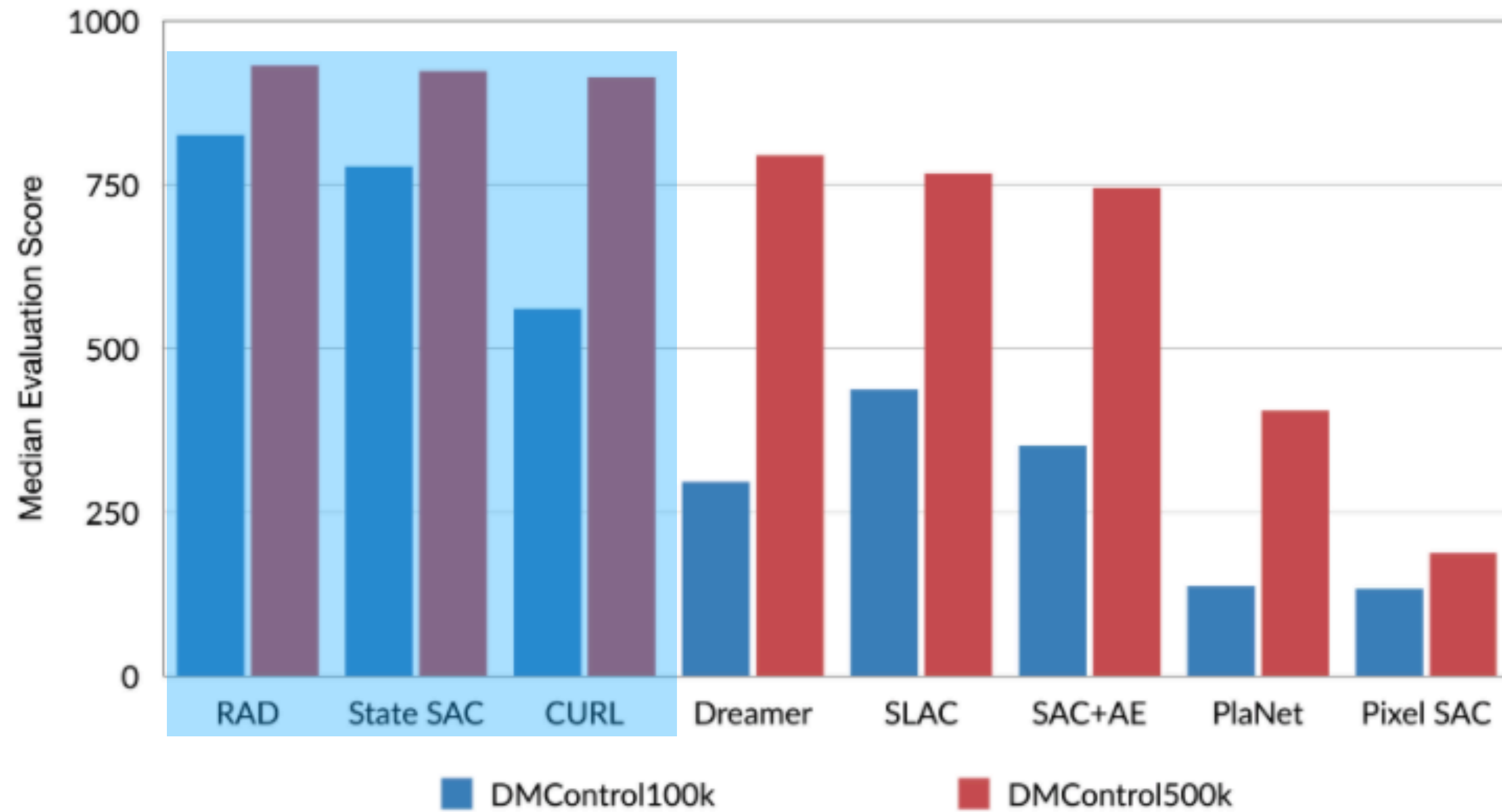
# Reinforcement Learning with Augmented Data



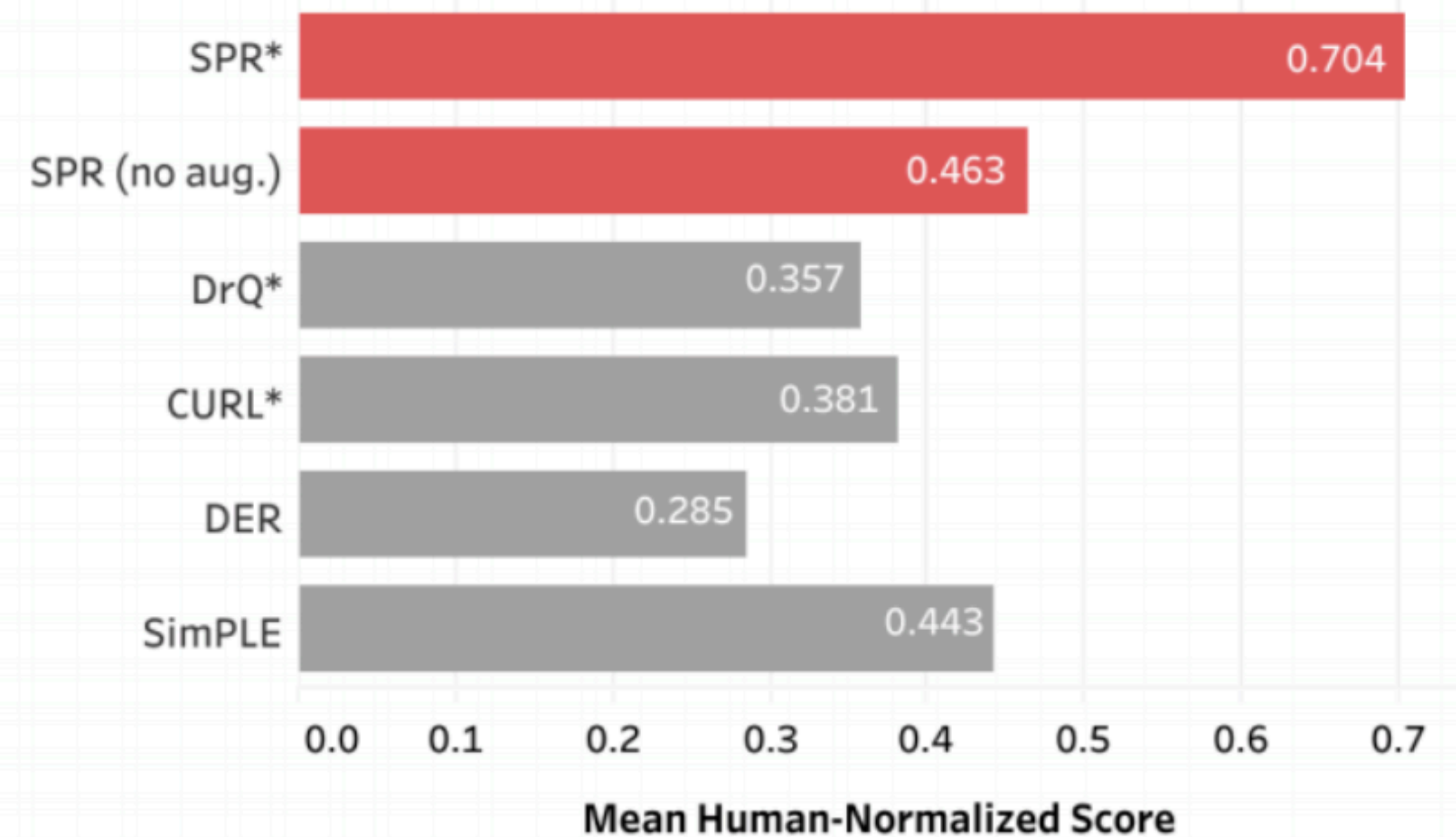
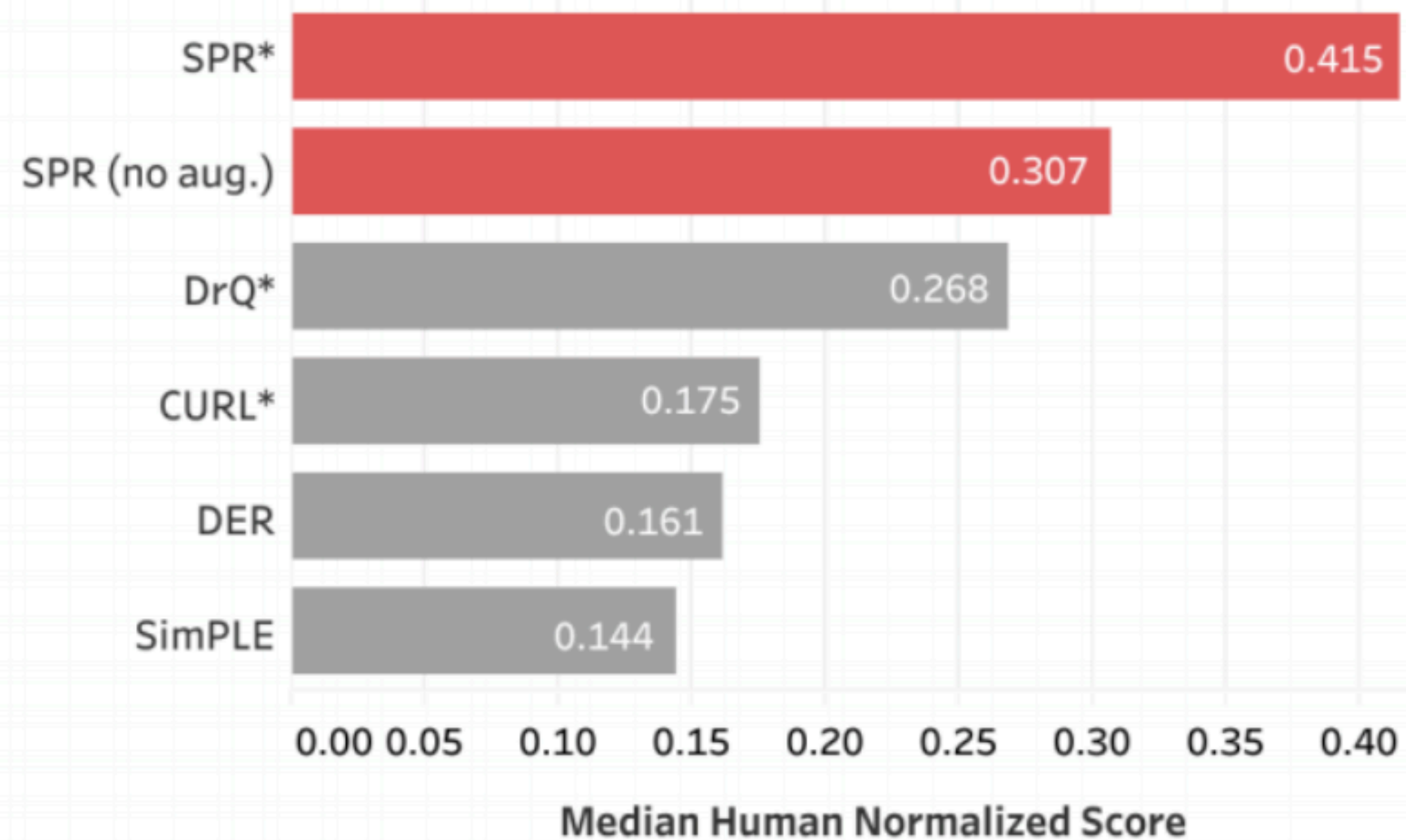
# Reinforcement Learning with Augmented Data



# Reinforcement Learning with Augmented Data



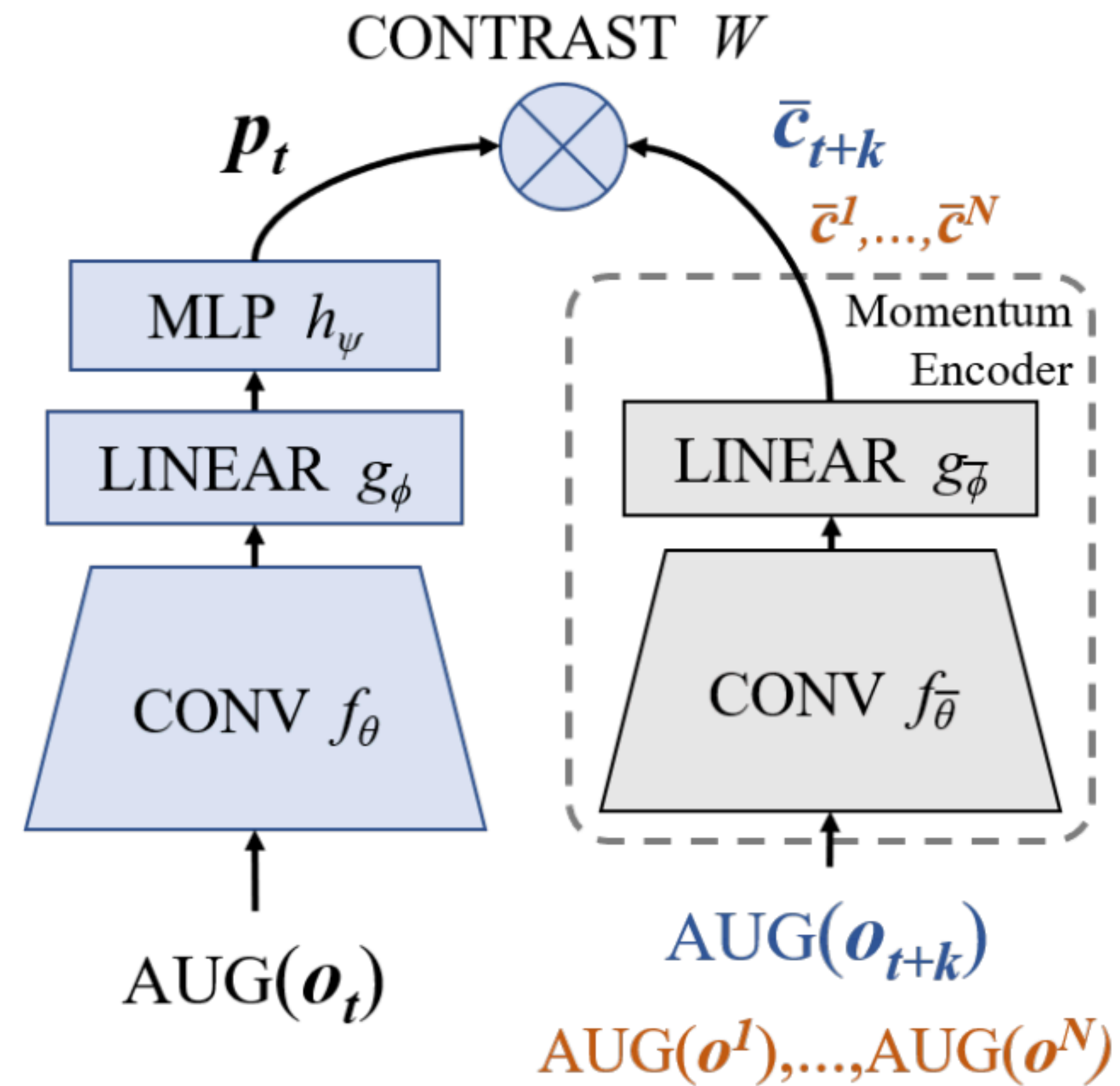
# Data Augmentation vs Auxiliary Losses



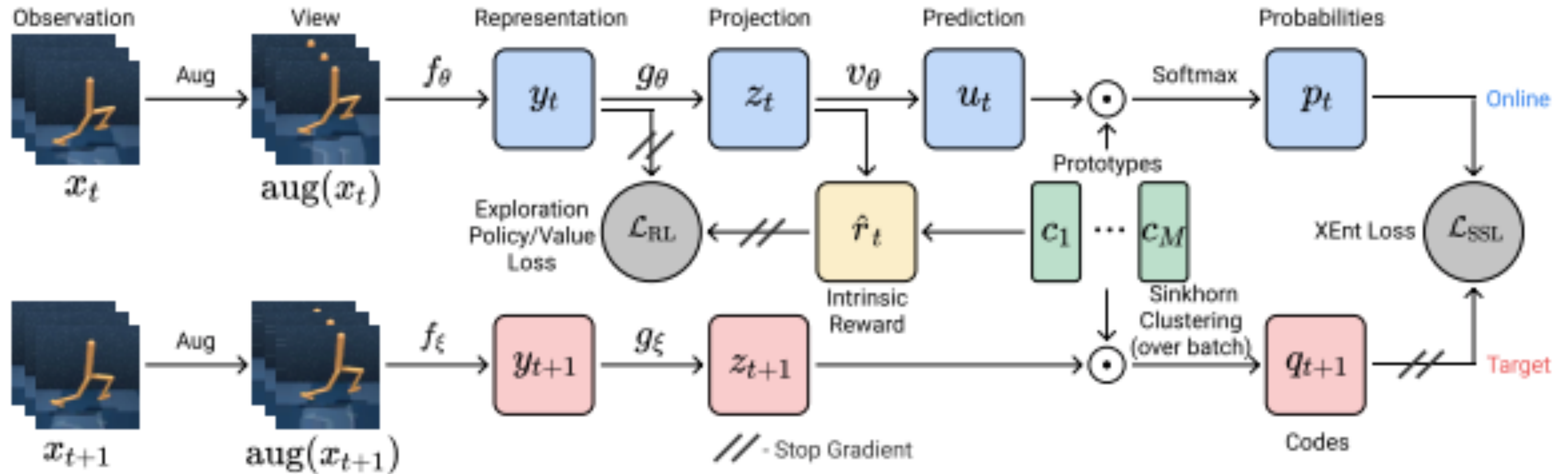
# Interpretation of augmentations in auxiliary and direct form

- CURL = Un-annotated (can be done w/o rewards or tasks; flexible with respect to positive & anchors)
- RAD = Annotated (only for a specific reward or task)
- CURL can help you learn one general purpose encoder for many tasks in the decoupled setup.
- Lots of followup work combining both the ideas now.

# Decoupling Representation and Reinforcement Learning



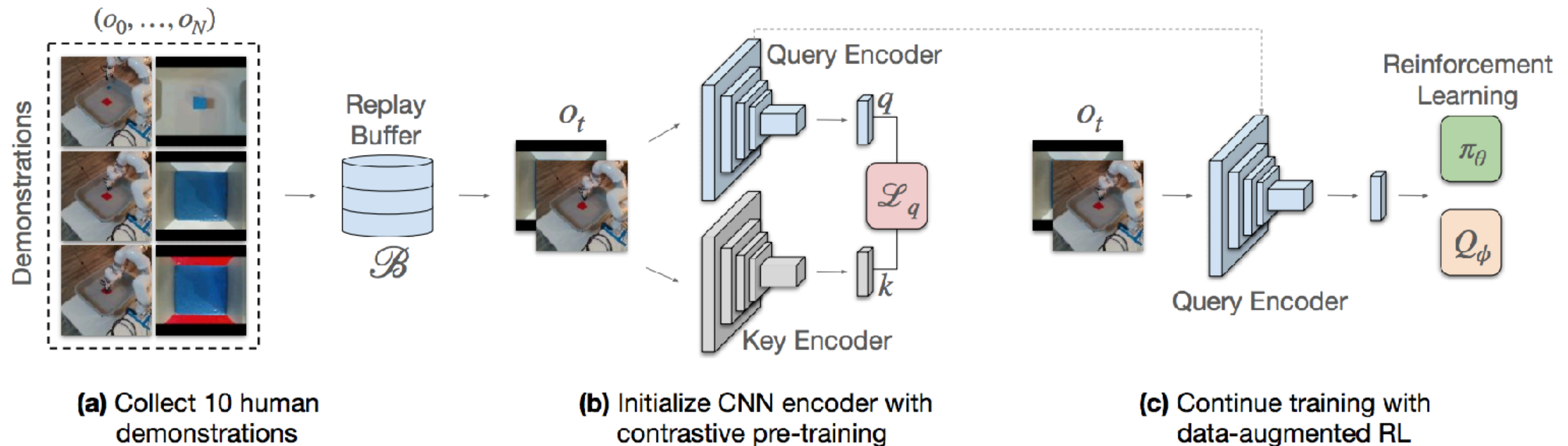
# SwaV + DrQ (Proto-RL) - Yarats



Proto-RL, Yarats et al 2021



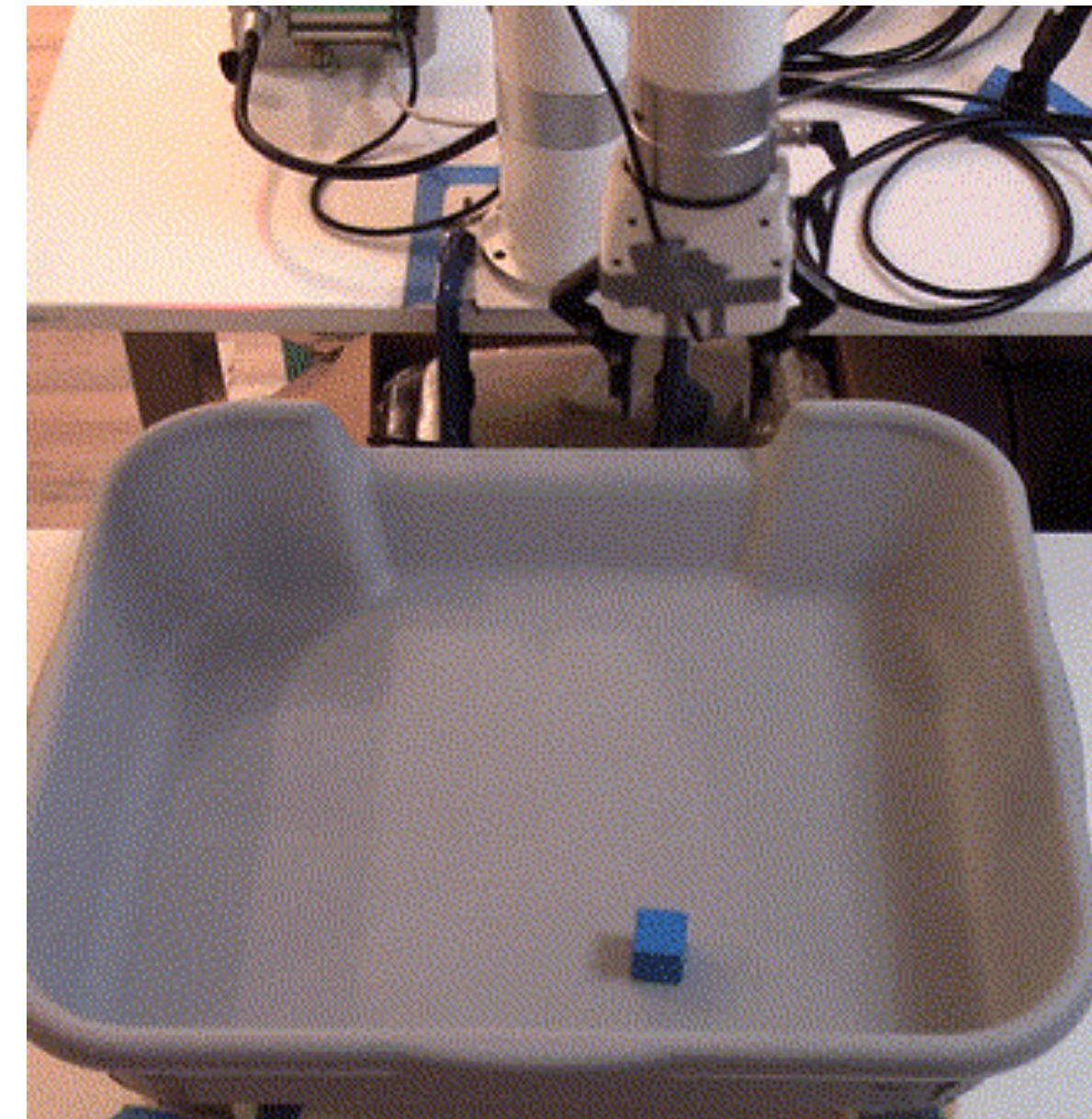
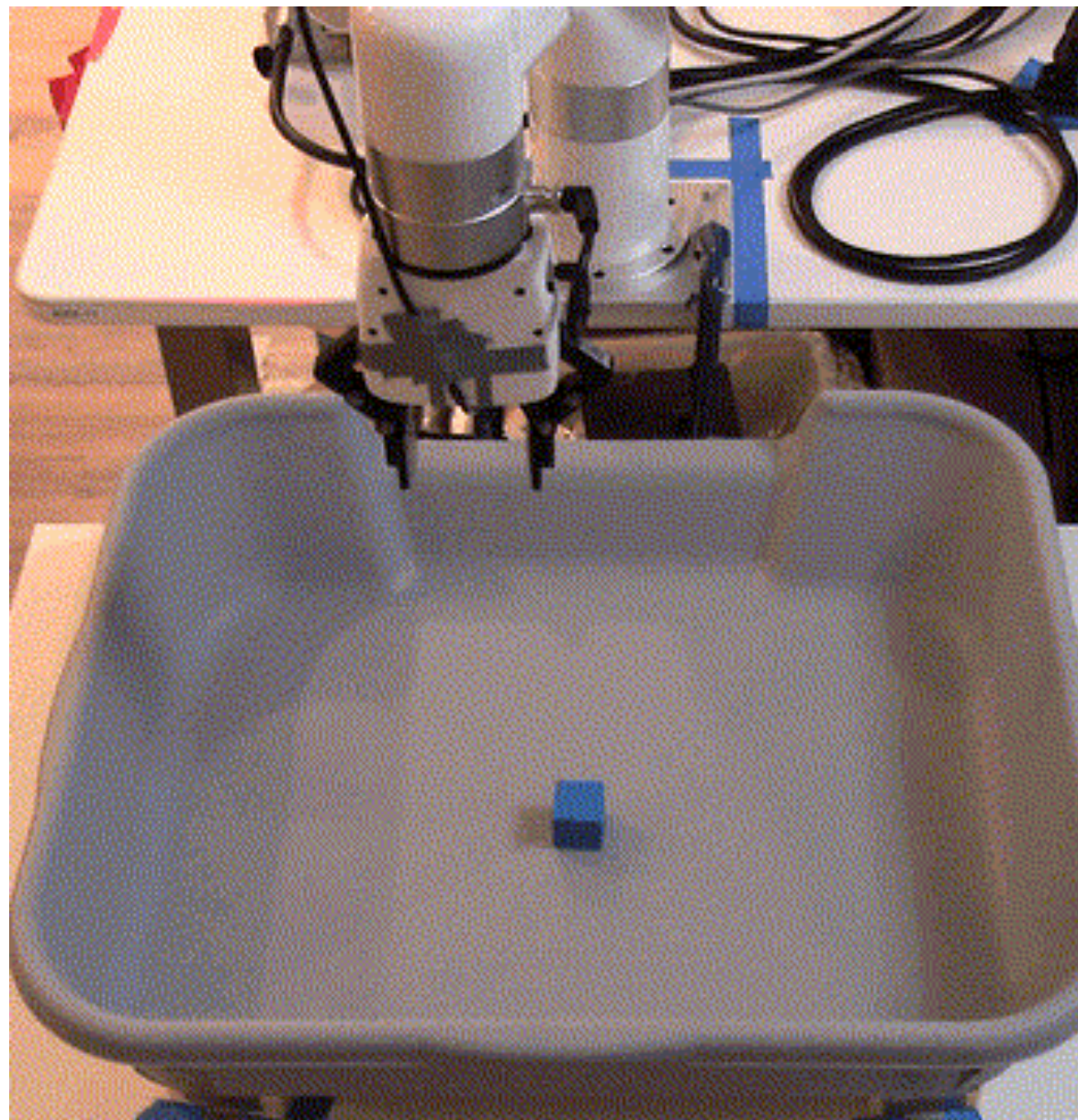
# FERM: CURL + RAD for Robotics



Robot Learning from Pixels for Sparse Reward Tasks, within 30 minutes of real time training

Framework for Efficient Robot Manipulation - Zhan et al, 2020

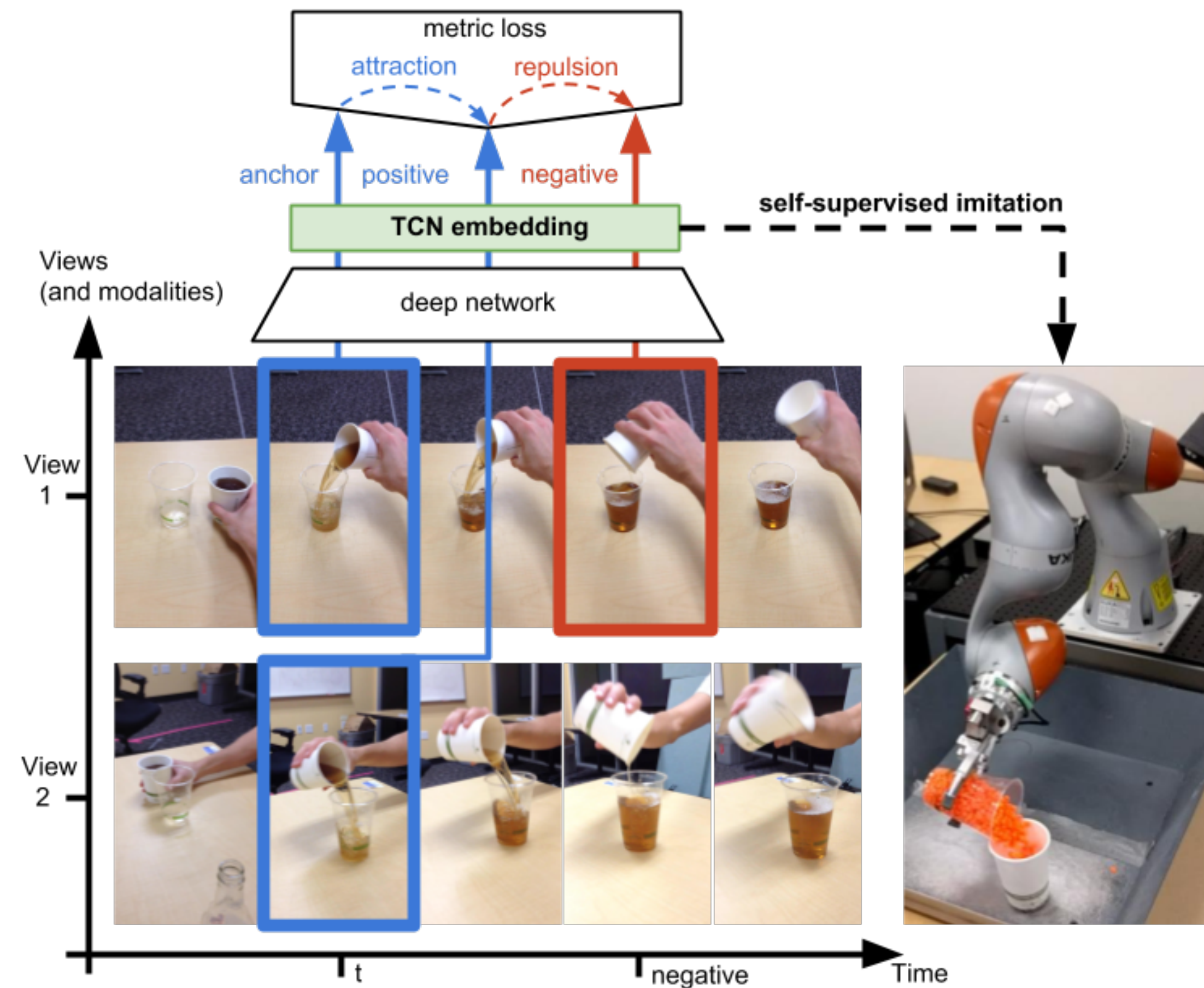
# FERM: CURL + RAD for Robotics



Robot Learning from Pixels for Sparse Reward Tasks, within 30 minutes of real time training

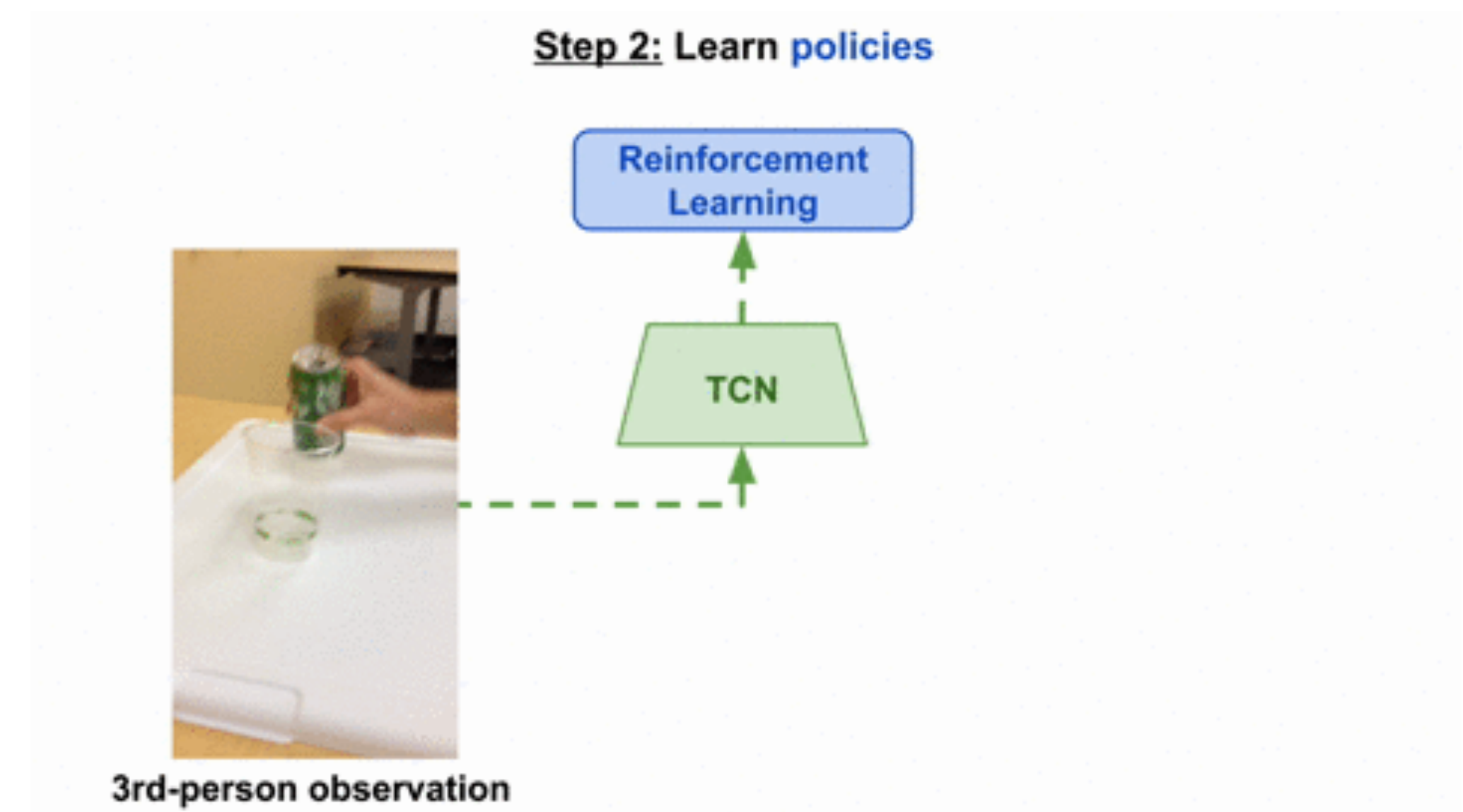
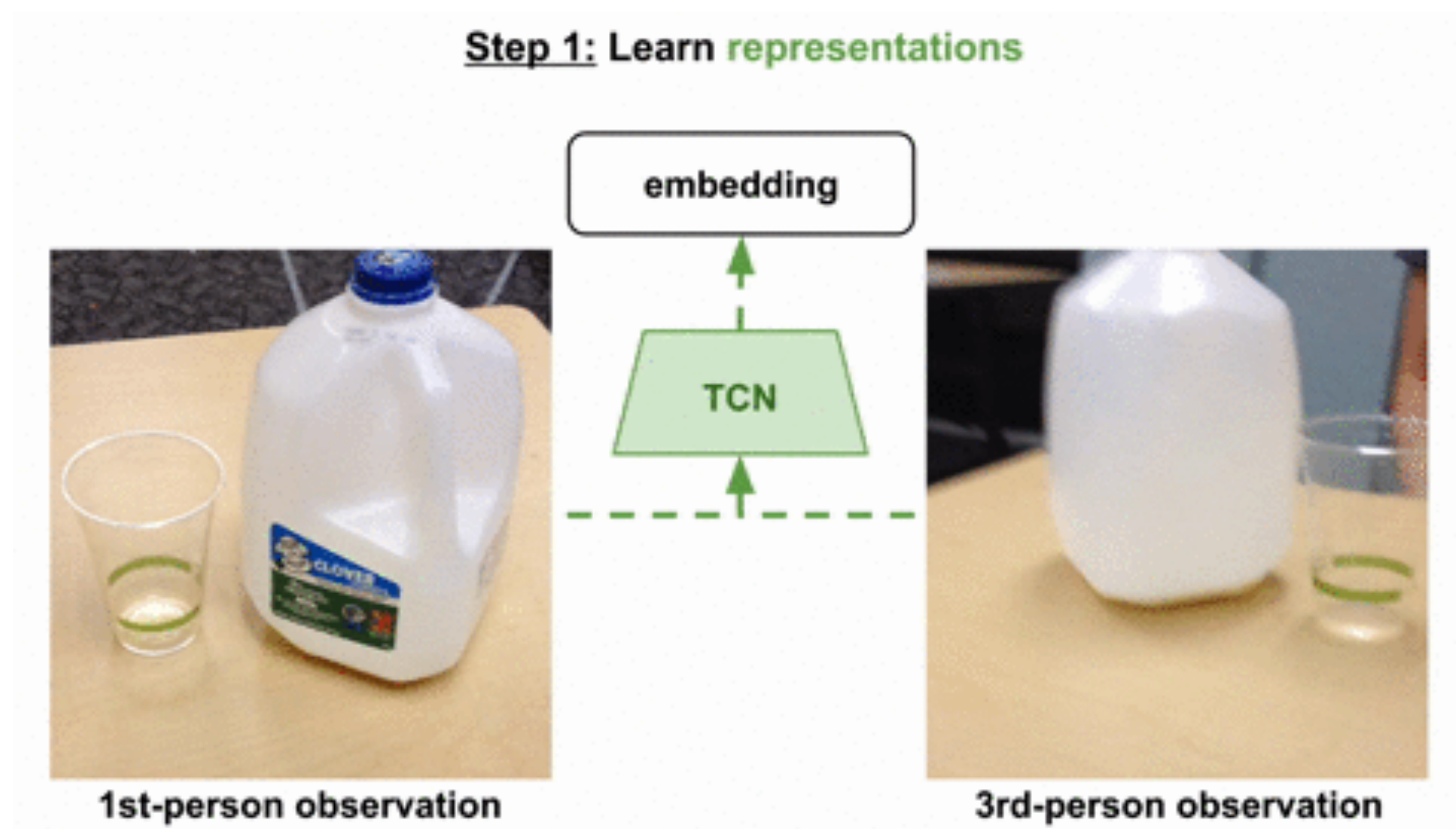
Framework for Efficient Robot Manipulation - Zhan et al, 2020

# Time Contrastive Networks (TCN)



Time Contrastive Networks: Self-Supervised Learning from Video, Sermanet et al 2018

# Time Contrastive Networks (TCN)



Time Contrastive Networks: Self-Supervised Learning from Video, Sermanet et al 2018

# Representative Work (not exhaustive)

1. UNREAL: Reinforcement Learning with Unsupervised Auxiliary Tasks (Jaderberg et al 2016)
2. DARLA: Improving Zero-Shot Transfer in Reinforcement Learning (Higgins et al 2017)
3. World Models (Ha & Schmidhuber, 2018)
4. Learning to summarize from human feedback (Stiennon et al 2020)

**Generative UL**

5. CPC: Contrastive Predictive Coding (Van den Oord et al 2018)
6. CURL: Contrastive Unsupervised Representations for RL (Srinivas & Laskin et al 2020)
7. DrQ: Image Augmentation is all you need - (Kostrikov & Yarats et al 2020),  
RAD: Reinforcement Learning with Augmented Data (Laskin et al 2020)
8. SPR: Self-Predictive Representations (Schwarzer & Anand et al 2020)
9. Decoupling Representation Learning from Reinforcement Learning (Stooke et al 2021)
10. Reinforcement Learning with Prototypical Representations (Yarats et al 2021)
11. Pre-training representations for data-efficient RL (Schwarzer & Rajkumar et al 2021)

**Contrastive-like UL**

12. Deep Spatial Autoencoders for Visuomotor Learning (Finn et al 2015)
13. TCN: Time Contrastive Networks (Sermanet et al 2018)
14. FERM: A Framework for Efficient Robotic Manipulation (Zhan et al 2020)

**Robotics Applications**

# Representative Work (not exhaustive)

1. UNREAL: Reinforcement Learning with Unsupervised Auxiliary Tasks (Jaderberg et al 2016)
2. DARLA: Improving Zero-Shot Transfer in Reinforcement Learning (Higgins et al 2017)
3. World Models (Ha & Schmidhuber, 2018)
4. Learning to summarize from human feedback (Stiennon et al 2020)

5. CPC: Contrastive Predictive Coding (Van den Oord et al 2018)
6. CURL: Contrastive Unsupervised Representations for RL (Srinivas & Laskin et al 2020)
7. DrQ: Image Augmentation is all you need - (Kostrikov & Yarats et al 2020),  
RAD: Reinforcement Learning with Augmented Data (Laskin et al 2020)

**8. SPR: Self-Predictive Representations (Schwarzer & Anand et al 2020)**

**Significant improvements**

9. Decoupling Representation Learning from Reinforcement Learning (Stooke et al 2021)
10. Reinforcement Learning with Prototypical Representations (Yarats et al 2021)
- 11. Pre-training representations for data-efficient RL (Schwarzer & Rajkumar et al 2021)**

12. Deep Spatial Autoencoders for Visuomotor Learning (Finn et al 2015)
13. TCN: Time Contrastive Networks (Sermanet et al 2018)
14. FERM: A Framework for Efficient Robotic Manipulation (Zhan et al 2020)

# Research Ideas

1. So many opportunities for future research in this space
2. Exploring objectives like DINO, Barlow Twins for RL - connections to Policy Distillation
3. Applying these ideas to indoor navigation in Habitat
4. Real-world results in robot manipulation
5. Observational Imitation (Revisiting TCN with the modern tools)
6. Learning world models in latent spaces discovered by Siamese Nets
7. Larger networks pre-training
8. Multi-task learning with shared unsupervised pre-trained backbones
9. ....
10. ....

**Thanks!**

[aravindsrinivas@gmail.com](mailto:aravindsrinivas@gmail.com)





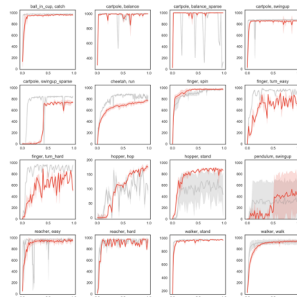
# ICML 2021 Tutorial on Unsupervised Learning for RL: Part II: Reward-Free RL

Pieter Abbeel & Aravind Srinivas  
UC Berkeley

Thanks to Marc Bellemare, Yang Gao, Misha Laskin, Sergey Levine, Hao Liu, Vlad Mnih, Pierre-Yves Oudeyer, Deepak Pathak, Lerrel Pinto, Aravind Rajeswaran for valuable feedback and suggestions.

# Tutorial Overview

- Part I: Representation Learning in RL



- Part II: Reward-Free RL

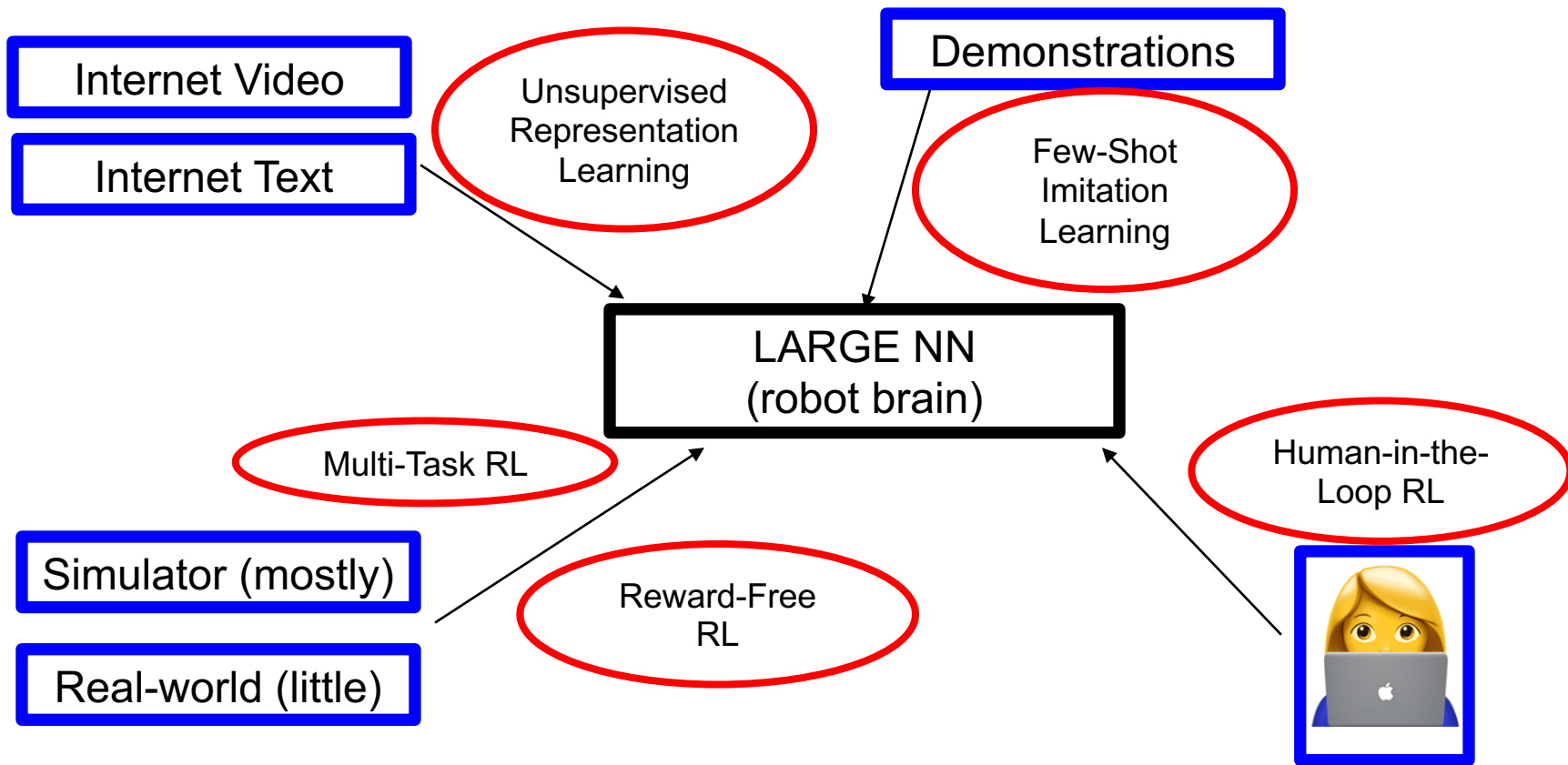


# But Wait: Lots of Passive Data Out There...

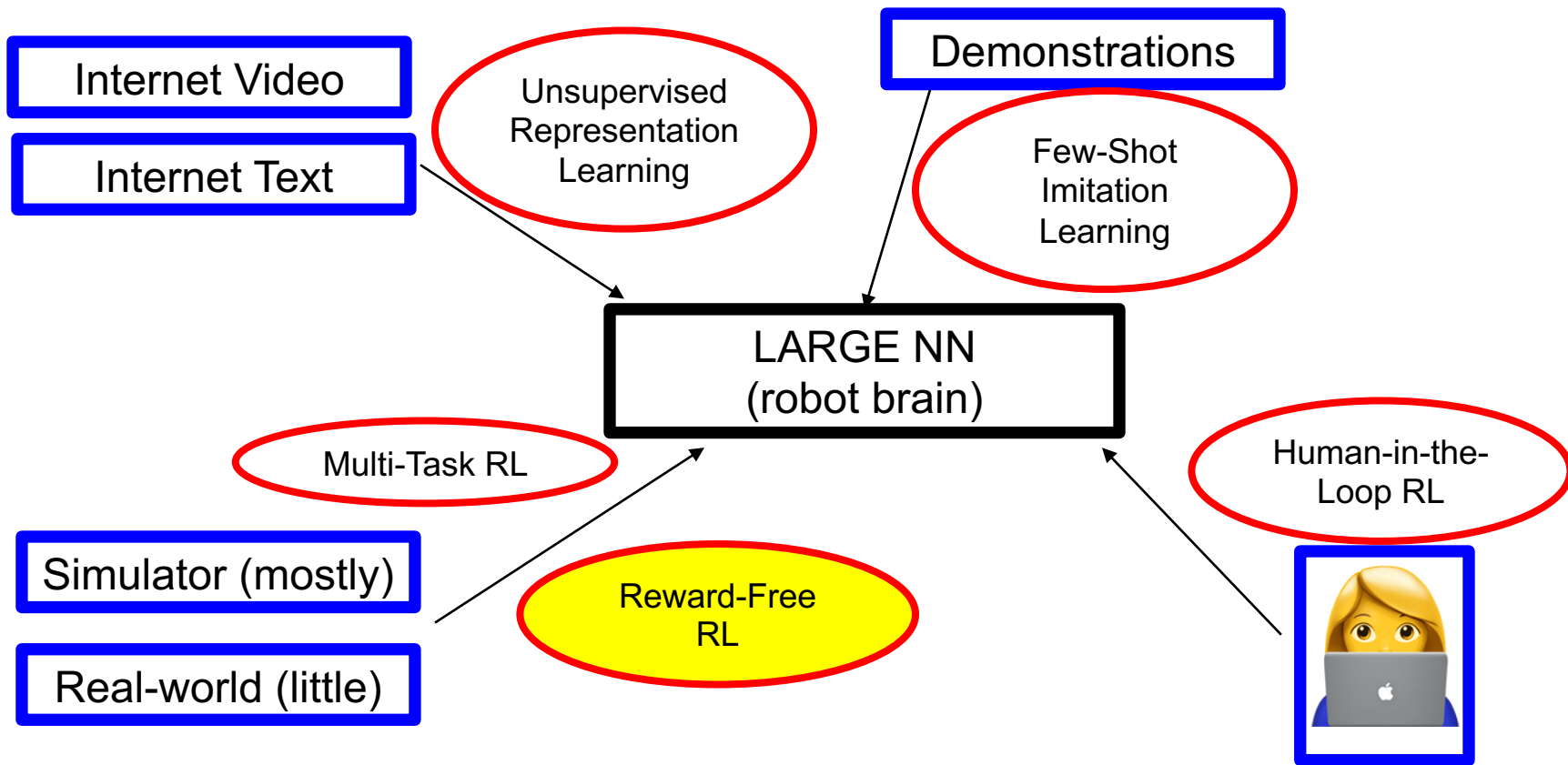
---

- So, shouldn't we leverage that?
  - YES, absolutely
  - But, IN ADDITION, we want our agents to be able to *intelligently collect their own data*, and that's what *this tutorial* will cover

# Interlude: An Attempt at a Complete Picture



# Interlude: An Attempt at a Complete Picture



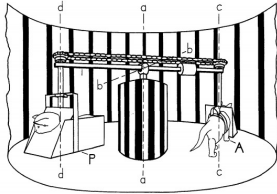
# Why Study Agents Collecting Their Own Data (1/2)

---

- In real world, reward can be tedious to provide
- Alleviate need for bootstrapping from demonstrations
- Superhuman solutions
  - Debug a new game for unexpected issues
  - Circuit design
  - Scientific discovery RL
  - ...

# Why Study Agents Collecting Their Own Data (2/2)

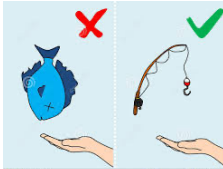
- Might lead to more robust learning (or not)



[Held & Hein, 1969]

also: [Walk, Shepherd, Miller, 1988]

- “general brain” for learning decision making



“Give an AI agent a task-reward, they can learn for a day,  
give an AI agent intrinsic reward, they can learn for a lifetime.”

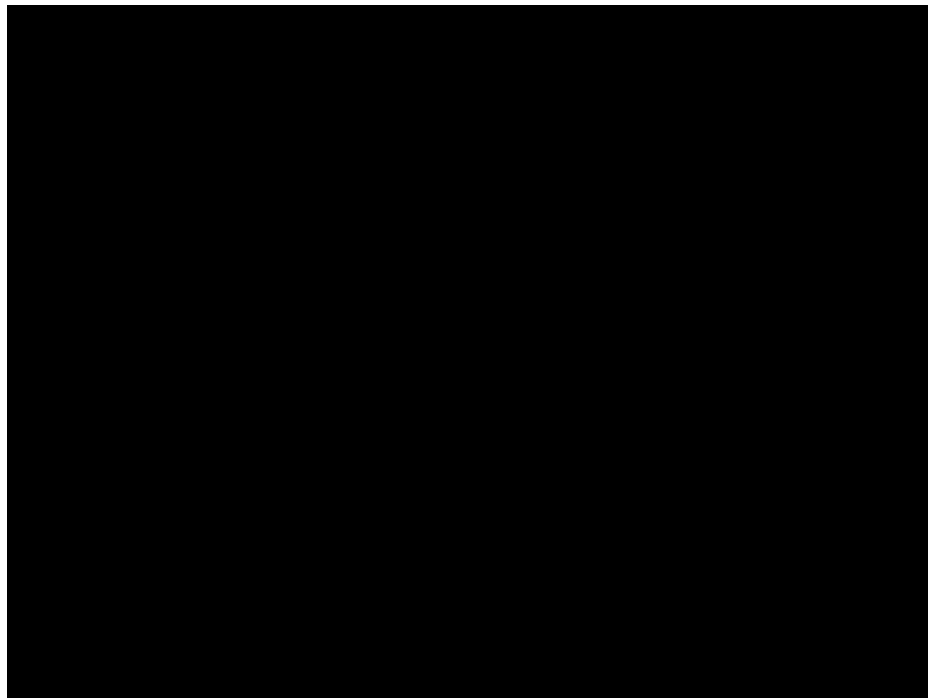
- AI motivation



[image source: Francis Vachon]

It's how humans/babies learn  
***lots of play*** (+ some supervision)

# Play



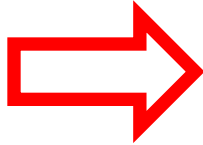
- What's going on here?  
[see, e.g., Gopnik (Berkeley), Schulz (MIT)]
- How to formalize into machine learning / silicon compute?
- What does success even mean?  
Robot rolling on the floor with toys??



# Formalizing “Play”: “Reward-Free Pre-Training (RFPT)”

Reward-free  
Pre-Training

Sample complexity  
typically not too heavily  
scrutinized here



## Possible measures of success:

- $\pi(a|s)$ : 0-shot
- $\pi(a|s)$ : fine-tune
- $\pi(a|s, g)$ : Goal-conditioned policy
- $\pi(a|s, z)$ : Latent-conditioned policy
- $\pi(a|s, z)$ : Diverse skill set
- $\{(s, a, s')\}$ : Data set
- $f_{\theta}(s'|s, a)$ : Dynamics model

Further challenge: some works are in image space, some use hardcoded image encodings, some in state space, some use hand-selected subsets of state variables, etc..  
Same video demo could have very different assumptions...

# Is Reward-Free Pre-Training Realistic?

- We don't simply leave a baby/child on its own for 18 years, then expect them to quickly fine-tune into adult life
- Similarly, it's likely our AI agents will benefit from a more interleaved approach between training on their own *and* getting access to side-information, e.g.
  - human providing some extrinsic rewards
  - human providing some rewards for good exploration
  - passive data being leveraged for representation learning
  - passive data being leveraged for guiding what's "interesting" to explore
  - etc...
- Such variants are beyond the scope of this tutorial
- Such variants will likely still greatly benefit from progress in "pure" RFPT regime

## Reward-Free Pre-Training

- First:  
pre-train without task reward available
- Then:  
leverage pre-training to learn faster  
once task reward is available

- Yes, very related, and so are research ideas/results for each
- Often arbitrary in which of the 2 contexts an idea ends up (and sometimes both)
- Ideas often complementary and by combining can amplify each other

## Exploration

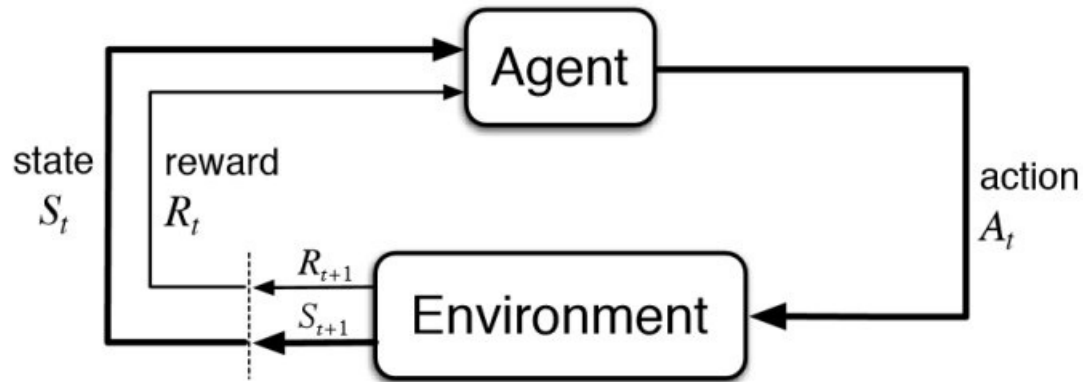
- Task reward is available from the very beginning
- But: often task reward provides very little signal in early learning stages, hence AI Agent needs exploration to find out how to achieve task reward
- Note: as long as no task reward has been experienced, the AI Agent is forced to learn reward-free

# Outline

---

- Problem Motivation
- ***Baseline RL Algorithms Refresher***
- Intrinsic Rewards for Reward-Free Pre-Training and Exploration
- Algorithmic Approaches to Exploration (can complement intrinsic reward RFPT!)
- Algorithmic Approaches to Reward-Free Pre-Training

# Reinforcement Learning (RL)



## ■ Popular RL Algorithms

- PPO
- SAC / TD3
- DDQN
- Model-Based RL (MBRL)

# Popular RL Algorithms are Mostly Optimizers, with Modest Exploration

- PPO, SAC, TD3, DDQN, MBRL in their standard forms
  - use experience so far to optimize expected reward in next roll-out
  - yet also have a *very basic exploration mechanism* built-in:
    - Epsilon-greedy
    - Boltzmann exploration (direct in DDQN, indirect in SAC)
    - Parameter perturbation [Plappert et al, 2017]
    - Noisy-nets [Fortunato et al, 2017]

→ Since so simple, remains very popular to just rely on these, even if clearly not very exploratory

→ But, *hopefully*, this statement about popularity and about superior simplicity can be untrue a few years from now 😊

# Outline

---

- Problem Motivation
- Baseline RL Algorithms Refresher
- ***Reward-Free Pre-Training and Exploration through Baseline RL Algorithms + Intrinsic Rewards***
- Algorithmic Approaches to Exploration (can complement intrinsic reward RFPT!)
- Algorithmic Approaches to Reward-Free Pre-Training

# Main Idea behind Intrinsic Rewards

---

- Intrinsic Reward = generic reward signal that encourages experiencing diversity
- IF we can design such Intrinsic Reward
- THEN we can simply use our existing RL Algorithms as Optimizers to yield exploratory / pre-training behaviors



# Reward-Free Pre-Training with Intrinsic Rewards

Run PPO / SAC / TD3 / DDQN / etc.:

$$\max E [ r_{\{INTRINSIC\}} ]$$

## Possible measures of success:

- $\pi(a|s)$ : 0-shot
- $\pi(a|s)$ : fine-tune
- $\pi(a|s, g)$ : Goal-conditioned policy
- $\pi(a|s, z)$ : Latent-conditioned policy
- $\pi(a|s, z)$ : Diverse skill set
- $\{(s, a, s')\}$ : Data set
- $f_{\theta}(s'|s, a)$ : Dynamics model

# Exploration with Intrinsic Rewards

Run PPO / SAC / TD3 / DDQN / etc.:

$$\max E [ r_{\{TASK\}} + \lambda r_{\{INTRINSIC\}} ]$$

## Measure of success:

- performance on  $r_{\{TASK\}}$

# Intrinsic Rewards

## CURIOUS MODEL-BUILDING CONTROL SYSTEMS

In Proc. International Joint Conference on Neural Networks, Singapore, volume 2, pages 1458-1463. IEEE, 1991.

Jürgen Schmidhuber\*  
Department of Computer Science  
University of Colorado  
Campus Box 430, Boulder, CO 80309, USA

[Schmidhuber, 1991]

## **What is intrinsic motivation? A typology of computational approaches**

---

**Pierre-Yves Oudeyer<sup>1,2,\*</sup> and Frederic Kaplan<sup>3</sup>**

1. Sony Computer Science Laboratory Paris, Paris, France
2. INRIA Bordeaux-Sud-Ouest, France
3. Ecole Polytechnique Federale de Lausanne, EPFL – CRAFT, Lausanne, Switzerland

[Oudeyer and Kaplan, 2007]

# Intrinsic Rewards

					Homeostatic (-) vs Heterostatic (+)	Motivation	Exploration potential	Organization potential	Computational cost	Existing models
Internal	Intrinsic	Adaptive	Knowledge-based	Information theoretic	+	UM	***	*	***	**
						IGM	***	***	***	**
					DSM	**	***	***	*	
				-	DFM	*	***	***	*	
				Predictive	+	NM	***	*	*	***
			-		ILNM	**	**	*	**	
			+		LPM	***	***	**	**	
			Competence-based	SM	**	**	**	*		
				-	FM	*	***	**	**	
				+	IM	***	*	**	*	
	CPM	***		***	**	*				
	-	CM		*	***	**	*			
	Fixed	Morphological	-	SyncM	*	***	**	**		
			StabM	*	***	*	**			
			+	VarM	***	*	*	*		
Extrinsic			-	SocM	/	/	*	***		
			-	EnerM	/	/	*	***		

# Three Main Types of Intrinsic Reward

---

- Knowledge-based: Surprise / unpredictability / how much learned about world from experience
- Competence-based: Empowerment / Skills
- Data-based: Entropy (i.e. coverage) of data collected

Note 1: Not the only way to categorize, but it's a categorization that can help us understand some of the most prominent work

Note 2: Data can give us Knowledge which can give us Competence, so one could possibly think of Data as the most fundamental, however, it's also least informed during data collection about the value of the data (for knowledge, for building competences)...

# Three Main Types of Intrinsic Reward

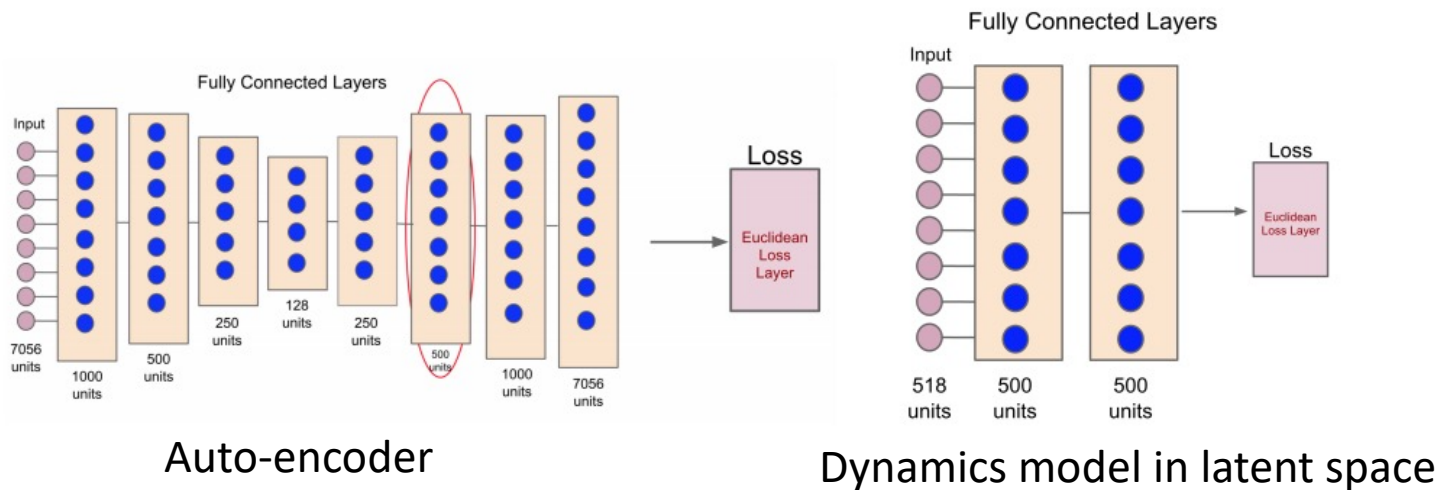
---

- **Knowledge-based: Surprise / unpredictability / how much learned about world from experience**
- Competence-based: Empowerment / Skills
- Data-based: Entropy (i.e. coverage) of data collected

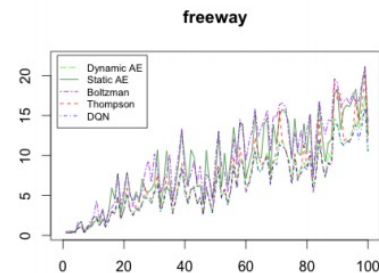
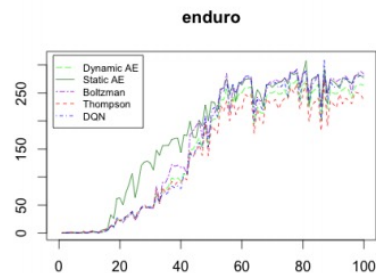
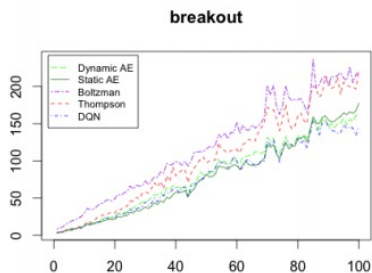
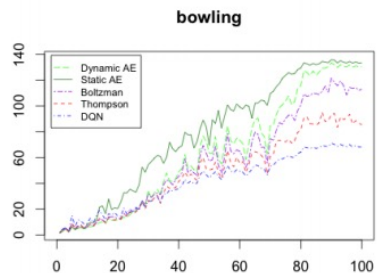
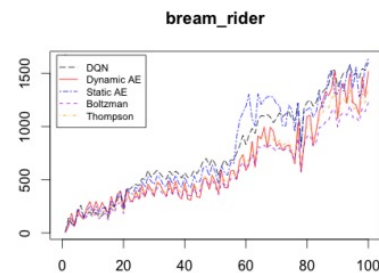
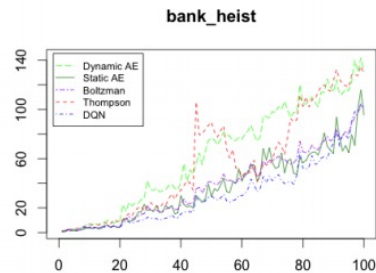
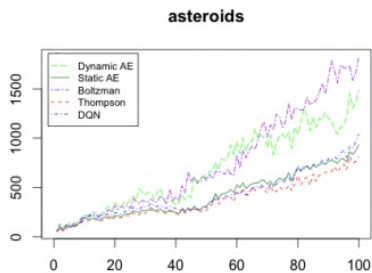
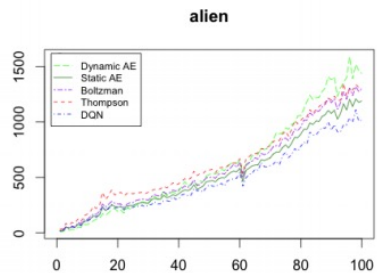
# Intrinsic Reward: Error in Learned Dynamics Model

- Key Idea:

- Train dynamics model on data collected by agent
- Intrinsic reward = prediction error of the learned dynamics model



# Intrinsic Reward: Error in Learned Dynamics Model



Experimental Findings: tends to often aid exploration / learning speed in Atari

# Why not doing even better?

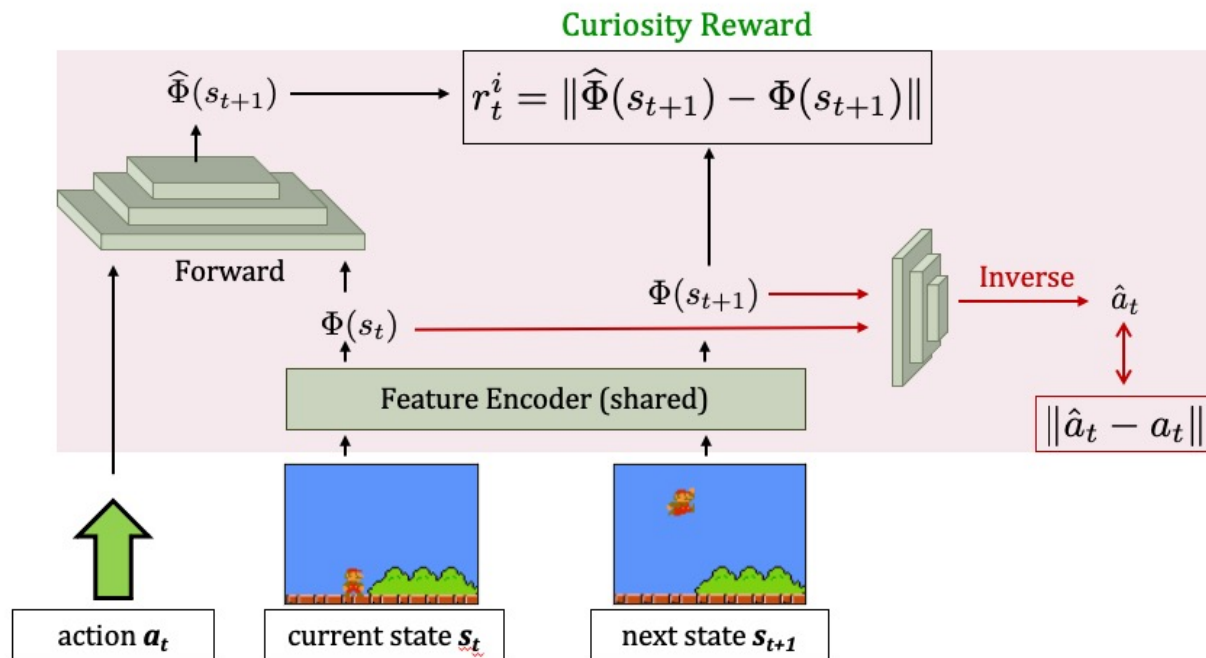
---

- Maybe:
  - VAE spends time modeling things the agent doesn't control
  - Dynamics model makes prediction errors on things the agent doesn't control (“noisy TV”)



# Inverse Dynamics Model for Encoder

Intrinsic Curiosity Module (ICM)



# Inverse Dynamics Model for Encoder



Key findings:

- Inverse dynamics helps focus learning
- 0-shot performance in many games

# Deeper dive on feature space?

---

## Large-Scale Study of Curiosity-Driven Learning

---

**Yuri Burda\***  
OpenAI

**Harri Edwards\***  
OpenAI

**Deepak Pathak\***  
UC Berkeley

**Amos Storkey**  
Univ. of Edinburgh

**Trevor Darrell**  
UC Berkeley

**Alexei A. Efros**  
UC Berkeley

- Pixels
- Random features
- VAE
- Inverse Dynamics

### Findings:

- Random features do well
- Learned features generalize better

# Current Status

- Intrinsic Reward:

$$r_t^i = \|\widehat{\Phi}(s_{t+1}) - \Phi(s_{t+1})\|$$

(with features trained with inverse dynamics)

- How can it still break down?
  - What if agent encounters dice and keeps rolling them?
    - It'll keep getting intrinsic reward, b/c the forward model keeps failing
  - More generally, need to distinguish:
    - Epistemic Uncertainty
    - Aleatoric Uncertainty

# Intrinsic Reward for Reduction in Epistemic Uncertainty

- Planning to be Surprised [Sun, Gomez, Schmidhuber, 2011]

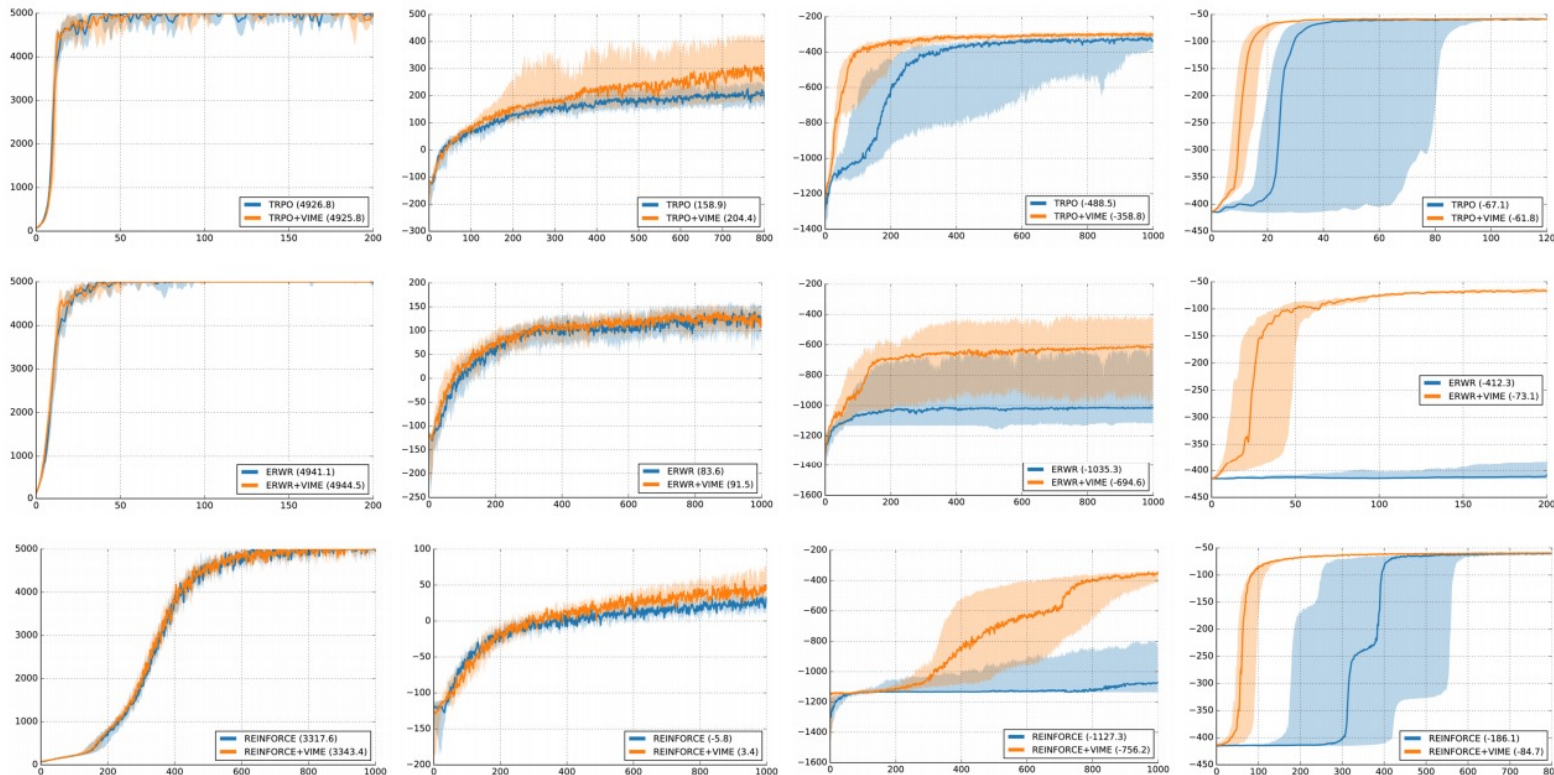
$$\sum_t (H(\Theta|\xi_t, a_t) - H(\Theta|S_{t+1}, \xi_t, a_t))$$

Reward for reducing entropy of posterior distribution over possible dynamics models

- VIME: Variational Information Maximization Exploration [Houthoofd, Chen, Duan, Schulman, De Turck, Abbeel, 2017]
  - Variational Bayes NN approximation to above (intractable objective)
- Self-Supervised Exploration via Disagreement [Pathak, Gandhi, Gupta, 2019]
  - Learn ensemble of NN dynamics models → disagreement signals epistemic uncertainty

# Intrinsic Reward for Reduction in Epistemic Uncertainty

VIME



(a) CartPole

(b) CartPoleSwingup

(c) DoublePendulum

(d) MountainCar

# Intrinsic Reward for Reduction in Epistemic Uncertainty

VIME

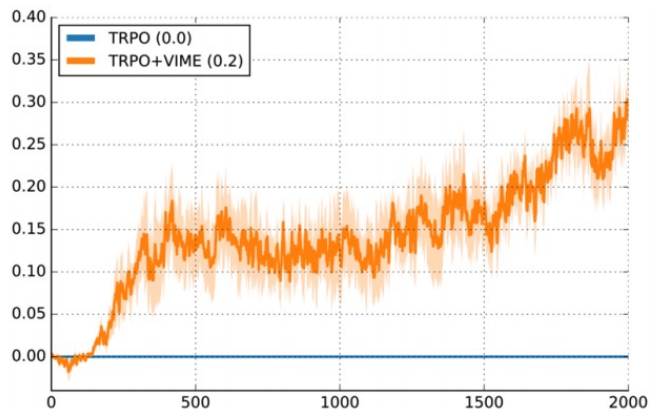
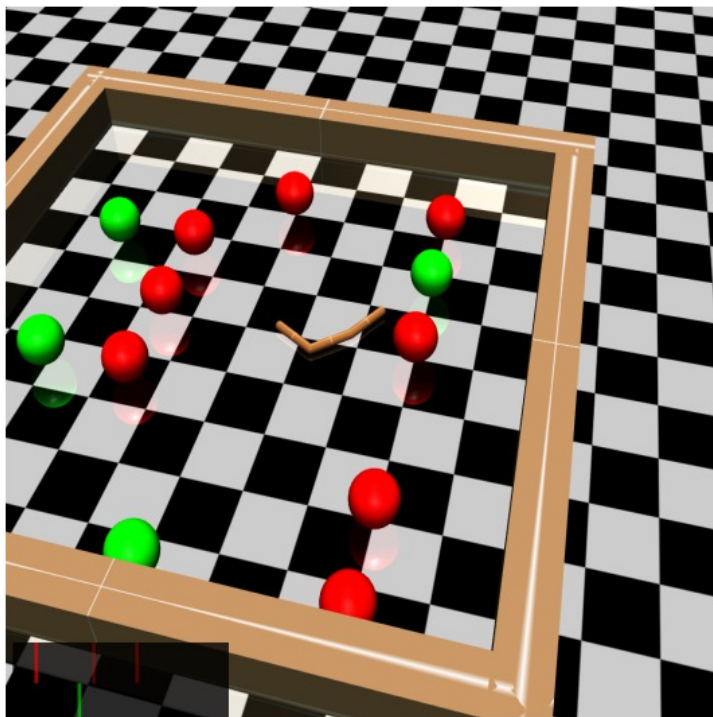
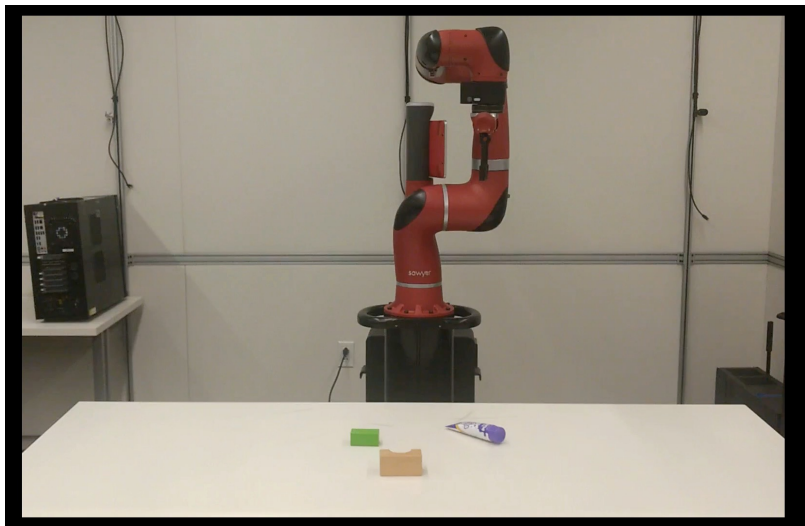


Figure 5: Performance of TRPO with and without VIME on the challenging hierarchical task SwimmerGather.

# Intrinsic Reward for Reduction in Epistemic Uncertainty

## Disagreement

Baseline: random exploration



Disagreement



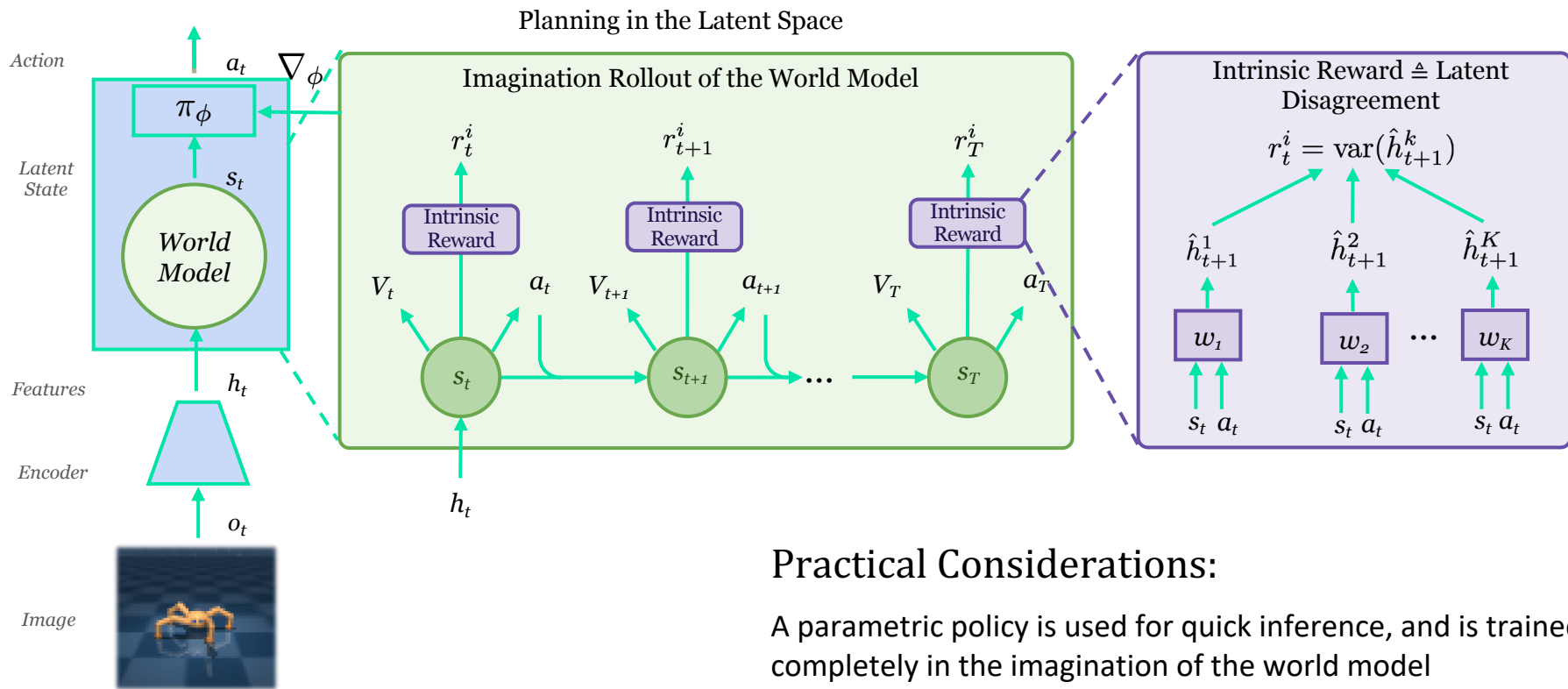
Setup: Overhead Camera; Action space = (position, direction, gripper angle, gripper opening)



# Current Status

- Intrinsic Reward = (reduction in) epistemic uncertainty
- As agent collects high intrinsic reward data, the epistemic uncertainty gets driven down
- What might still cause agent inefficiency?
  - Agent only gets rewarded for (reduction in) epistemic uncertainty \*after\* it's achieved it. I.e. agent needs to be lucky, and will then be encouraged to learn from that luck.
  - More generally, can distinguish:
    - Retrospective signals
    - Prospectively seeking out signals

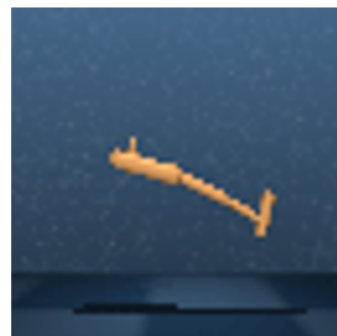
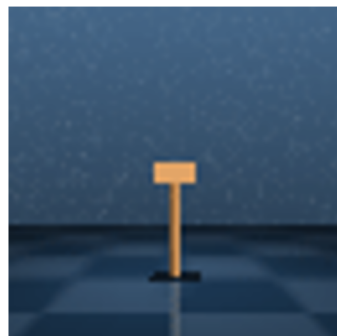
# Planning to Achieve Disagreement



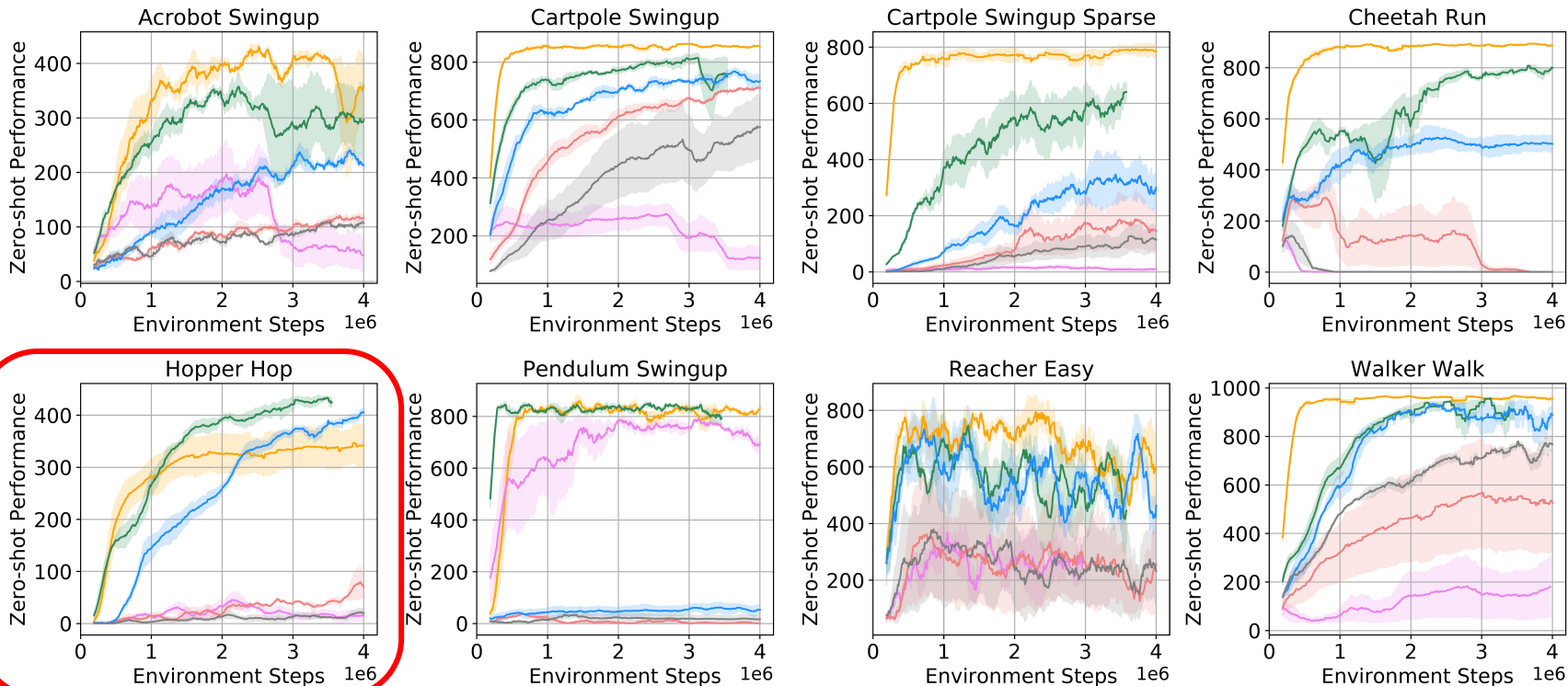
## Practical Considerations:

A parametric policy is used for quick inference, and is trained completely in the imagination of the world model

# Self-Supervised Exploration Results



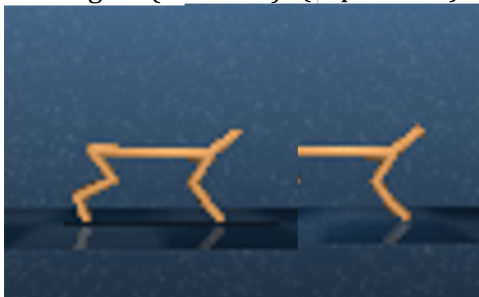
# Zero-Shot Reinforcement Learning



— Dreamer (Hafner et al., 2020) [sup] — Plan2Explore (Ours) [unsup] — Curiosity (Pathak et al., 2017) [unsup]  
— MAX (Shyam et al., 2019) [unsup] — Retrospective (Pathak et al., 2019) [unsup] — Random [unsup]

# Few-Shot Adaptation

Our Agent (few-shot) Oracle (supervised)

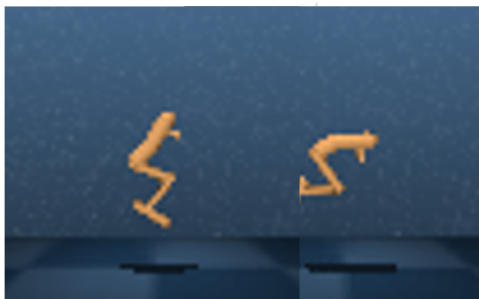


Cheetah Run

Our Agent (few-shot) Oracle (supervised)



Reacher Easy

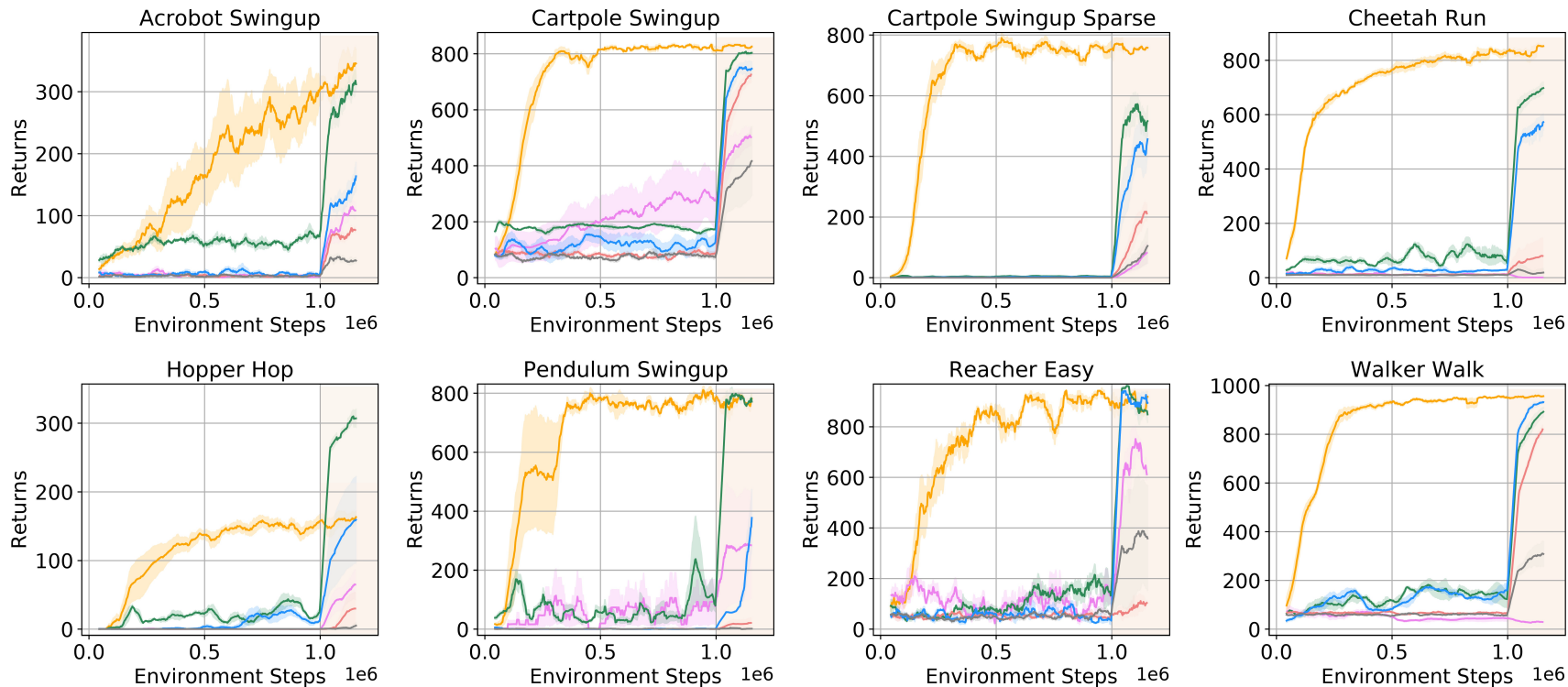


Hopper Hop



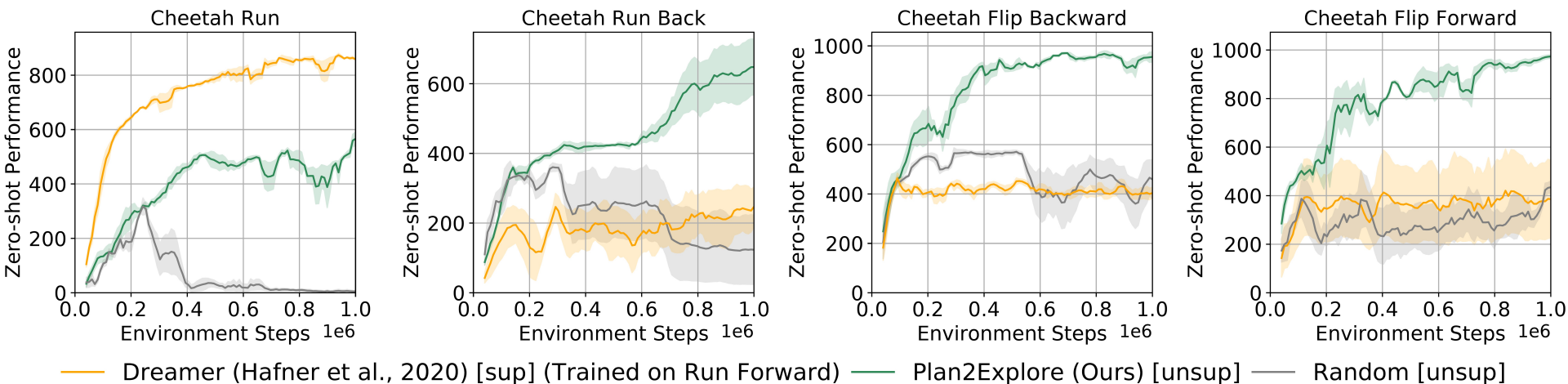
Walker Walk

# Few-Shot Adaptation



— Dreamer (Hafner et al., 2020) [sup] — Plan2Explore (Ours) [unsup] — Curiosity (Pathak et al., 2017) [unsup]  
— MAX (Shyam et al., 2019) [unsup] — Retrospective (Pathak et al., 2019) [unsup] — Random [unsup]

# Can one model be used for multiple tasks?



# Three Main Types of Intrinsic Reward

---

- Knowledge-based: Surprise / unpredictability / how much learned about world from experience
- **Competence-based: Empowerment / Skills**
- Data-based: Entropy (i.e. coverage) of data collected



# Competences

- Why care about the agent's competences?
  - High competence is a natural desideratum after pre-training
  - Competence-based RFPT can be seen as end-to-end optimizing for the competence desideratum during RFPT
  - Other desiderata are possible, of course, e.g. 0-shot, fast fine-tuning, dataset coverage, high quality dynamics model
  - High competence might be achievable in “Phase 2” assuming great dataset coverage or great learned dynamics model

# Formalizing Competence: Empowerment

- Formalizes and quantifies the degrees of freedom (or options) that an organism or agent has as a proxy for “preparedness”
- Behavioral Empowerment Hypothesis
  - “The adaptation brought about by natural evolution produced organisms that in absence of specific goals behave as if they were maximising their empowerment.”
- Evolutionary Empowerment Hypothesis
  - “The adaptation brought about by natural evolution increases the empowerment of the resulting organism.”
- AI Empowerment Hypothesis
  - “Empowerment provides a task-independent motivation that generate AI behaviour which is beneficial for a range of goal-oriented behaviour.”

# Empowerment

$$\mathfrak{E}_t = C\left(p(s_{t+n}|a_t^n)\right) = \max_{p(a_t^n)} I(A_t^n; S_{t+n})$$

$A$ : the agent's actuator<sup>3</sup> which takes values  $a \in \mathcal{A}$

$S$ : the agent's sensor which takes values  $s \in \mathcal{S}$

$M$ : the agent's internal state (or memory) which takes values  $m \in \mathcal{M}$

$R$ : the state of the environment which takes values  $r \in \mathcal{R}$

- Measures the optionality a (perfect) agent can create over its future
- As phrased, is a number associated with a state + environment
- Empowerment can be increased by:
  - (1) going to more empowered states
  - (2) changing the agent/environment

# Empowerment

$$\mathfrak{E}_t = C\left(p(s_{t+n}|a_t^n)\right) = \max_{p(a_t^n)} I(A_t^n; S_{t+n})$$

So, what does this have to do with Exploration / RFPT / Competences?

An agent has maximal competence if it has figured out a policy that achieves the "empowerment / channel capacity" available in its environment

→ Intrinsic Reward that measures the ***empowerment level of the agent's policy***

- (1) going to more empowered states
- (2) changing the agent/environment

# Empowerment through Mutual Information Maximization

- Find policy  $\pi_{\theta}(a | s, z)$  such that  $MI(z ; \tau)$  is maximized
- How to compute MI?
- $MI(z ; \tau) = H(z) - H(z | \tau) = H(\tau) - H(\tau | z)$
- Both decompositions have been investigated (though mostly the first one)

# Variational Approximation to Mutual Information

$$MI(z ; \tau) = H(z) - H(z | \tau)$$

$$\rightarrow r_{\{intrinsic\}} = \log q_{\phi}(z | \tau) - \log p(z)$$

- Intuition: need to be able to recover the “intent”/“latent”  $z$  from the trajectory; i.e. different  $z$  results in distinctly different trajectories
- For  $z$  discrete: this means training a classifier, and can set  $p(z)$  equal to uniform to maximize entropy

# Variational Approximation to Mutual Information:

$$MI(z ; \tau) = H(z) - H(z | \tau)$$

- SSN4HRL:  $r_{\{intrinsic\}} = \log q_{\phi}(z | s_t) - \log p(z)$
- VIC:  $r_{\{intrinsic\}} = \log q_{\phi}(z | s_H) - \log p(z)$
- DIAYN / VISR:  $r_{\{intrinsic\}} = \log q_{\phi}(z | s_t) - \log p(z)$
- VALOR:  $r_{\{intrinsic\}} = \log q_{\phi}(z | s_{1:H}) - \log p(z)$
- DISCERN:  $r_{\{intrinsic\}} = \log q_{\phi}(z | s_{1:H}) - \log p(z)$

DISCERN is the only one to use continuous  $z$ , which doesn't allow for a classifier, so it uses a contrastive approximation to the first term

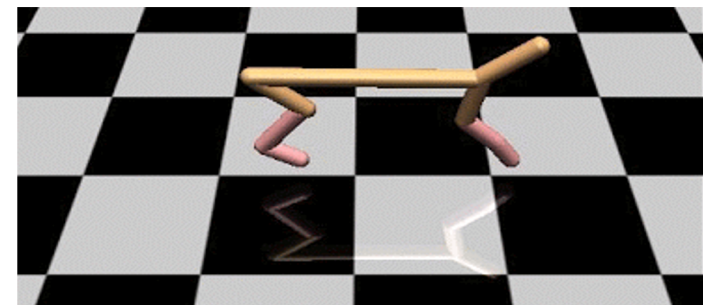
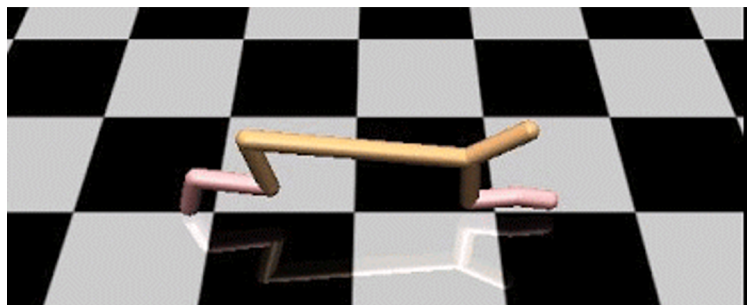
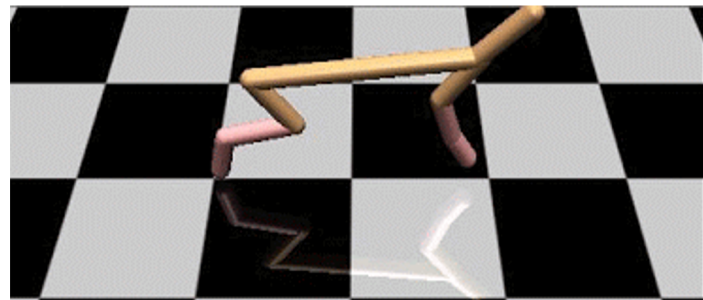
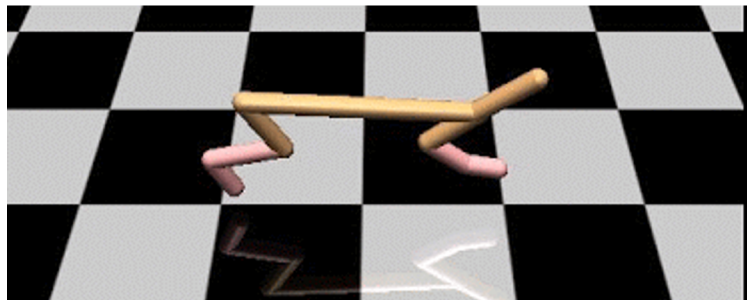
# References for Previous Slide

---

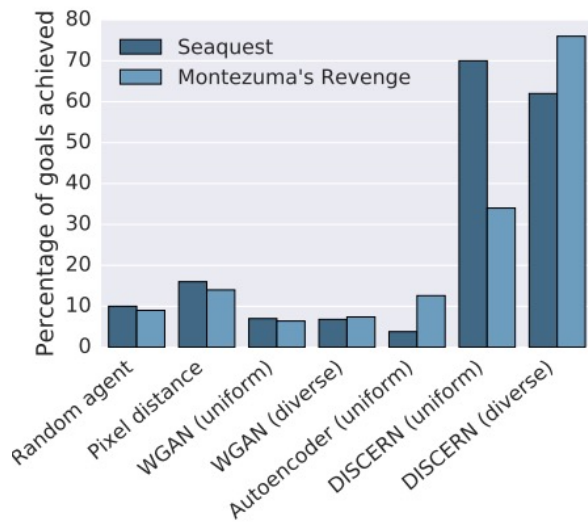
- SSN4HRL: Florensa, Duan, Abbeel, 2017
- VIC: Gregor, Rezende, Wierstra, 2016
- DIAYN: Eysenbach, Gupta, Ibarz, Levine, 2018
- VALOR: Achiam, Edwards, Amodei, Abbeel, 2018
- DISCERN: Warde-Farley, Van de Wiele, Kulkarni, Ionesu, Hansen, Mnih, 2018
- VISR: Hansen, Dabney, Barreto, Warde-Farley, Van de Wiele, Mnih, 2019



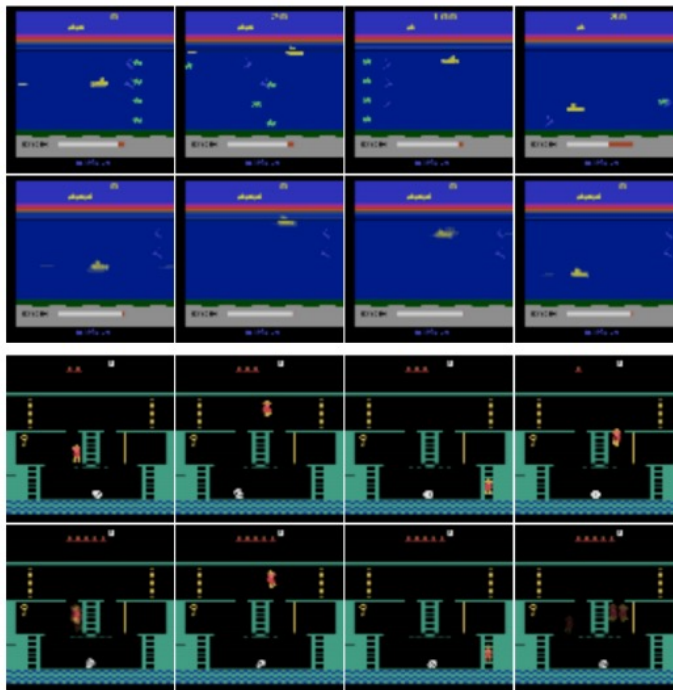
# Results Highlights: DIAYN



# Results Highlights: DISCERN



(a)



(b)

# Results Highlights: DISCERN

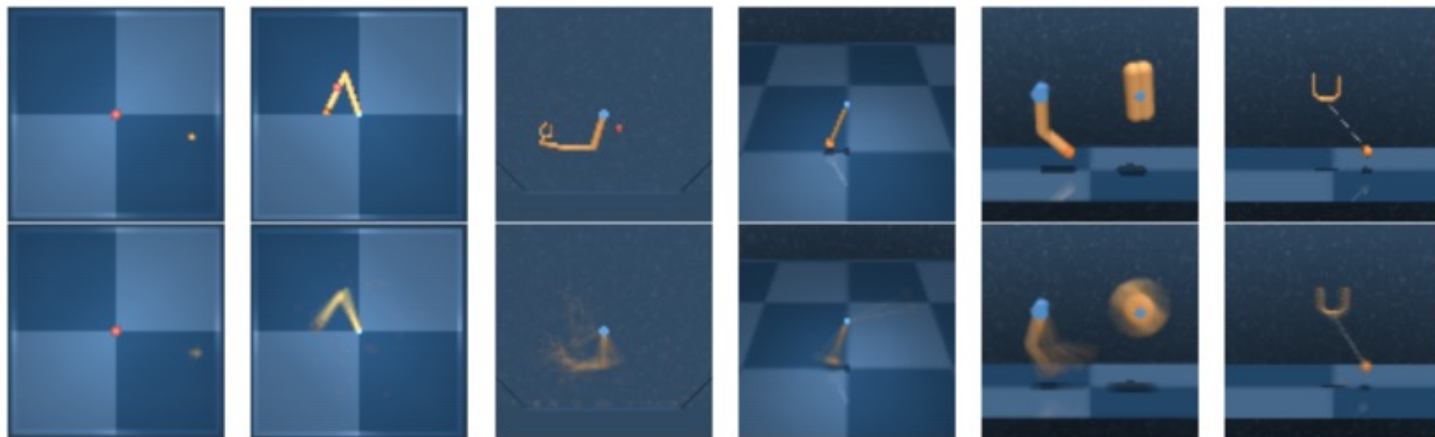


Figure 2: Average achieved frames for point mass (task easy), reacher (task hard), manipulator (task bring ball), pendulum (task swingup), finger (task spin) and ball in cup (task catch) environments. The goal is shown in the top row and the achieved frame is shown in the bottom row.

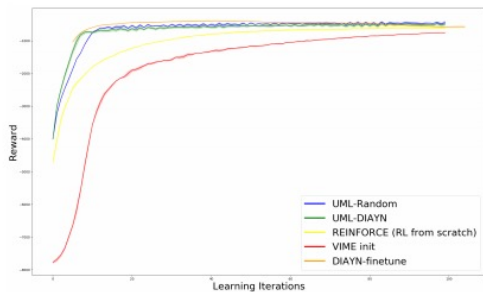
# VISR

---

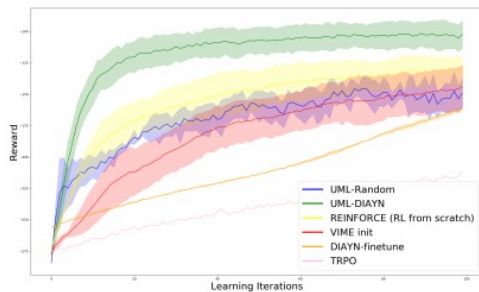
- Addresses slowness of fine-tuning
- Main ideas:
  - Same MI objective as DIAYN
  - Add success features [Barreto et al, 2017] to the agent

# Unsupervised Meta-RL

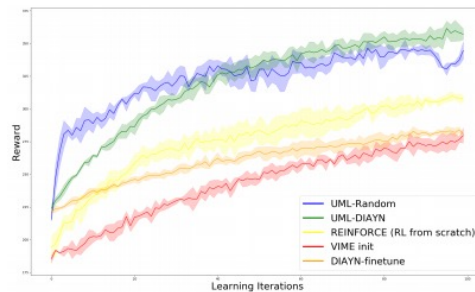
- Challenge: how to quickly adapt to new rewards?
  - Key Ideas:
    - Train with DIAYN (or similar)
    - After training, we can consider  $z$  as indexing into tasks, with reward  $\log q(z|s)$
- use this as tasks for Meta-RL training



2D navigation



Half-Cheetah

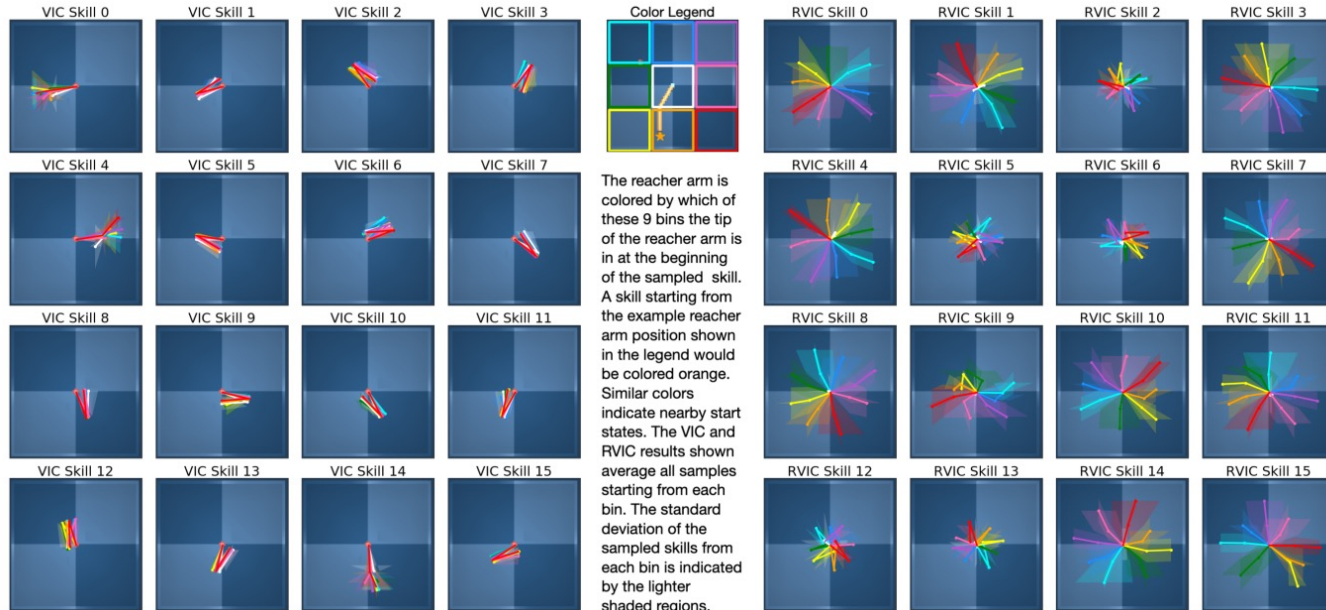


Ant

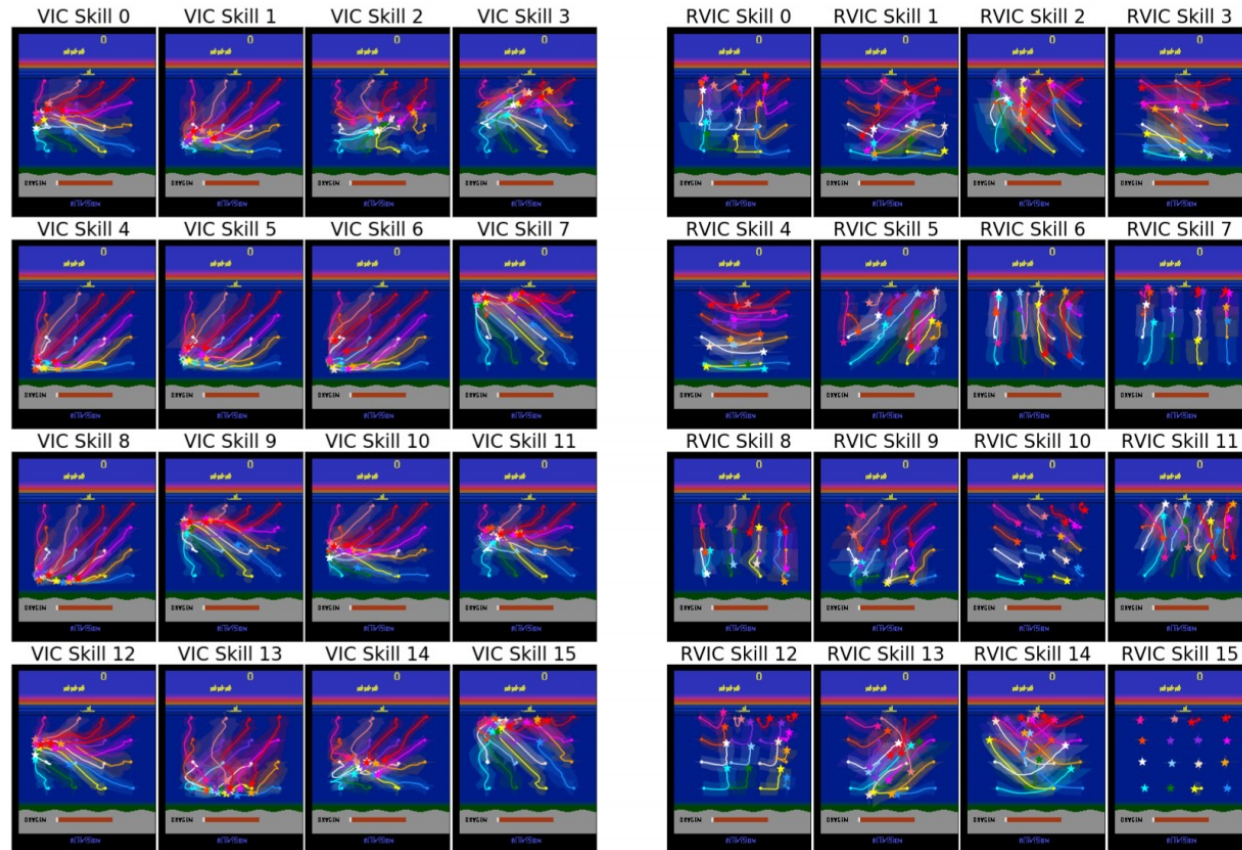
# Relative VIC

- Main idea:  $z$  ( $\omega$ ) more readily recoverable from  $s_0$  and  $s_T$  than just  $s_T$

$$\log q_\phi(\Omega|s_T, s_0) - \log q_\psi^{\text{abs}}(\Omega|s_T)$$



# Relative VIC



# DADS

$$r_{\{intrinsic\}} = \log q_{\phi}(s_{t+1}|s_t, z) - \log q_{\phi}(s_{t+1} | s_t)$$





# Some Extra Perspective on $MI(z; \tau) = H(z) - H(z | \tau)$

- Different papers treat  $\tau$  slightly differently, but the bigger difference might be in other experiment design factors:
  - Which RL algorithm is used?
  - State-based or image-based (only DISCERN is image based)
  - Termination conditions of episodes affect learning signal
  - When state-based, which variables to include in  $\tau$ ?
    - E.g. all state variables, vs. just x-y coordinates of ant robot?

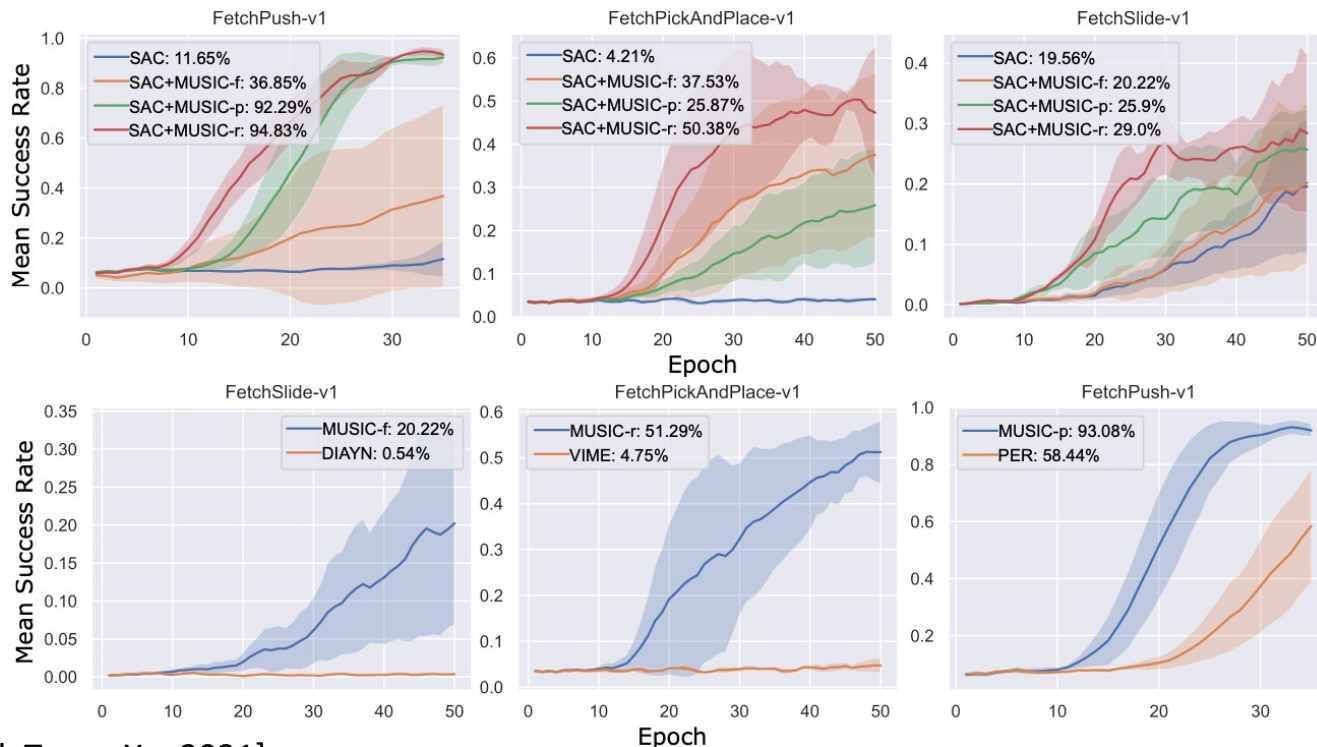
# MUSIC: Mutual Info State Intrinsic Control

---

- Embraces there is a natural split of the state:
  - State of the agent (body)
  - State of the environment around the agent
- Optimizes  $MI(\text{agent state ; env state})$ 
  - i.e. incentivize the agent to visit agent states that ***affect*** the environment

# MUSIC: Mutual Info State Intrinsic Control

- The proposed intrinsic reward help the agent to quickly learn to solve different downstream tasks.



# More Perspective on $MI(z; \tau) = H(z) - H(z | \tau)$

- Limited exploration signal:
  - $H(z)$  is maximized by simply sampling the latent from a high entropy distribution at the start of the trajectory
  - $H(z | \tau)$  encourages re-visiting states for which the classifier can recover  $z$ , which in turn likely leads to also visiting nearby states and gradual expansion, but not a very strong signal
- Might want to combine with other exploration

# How about the other decomposition of MI

- $MI(z ; \tau) = H(\tau) - H(\tau | z)$ 
  - $H(\tau)$  directly encourages coverage / exploration
  - $-H(\tau | z)$  encourages predictability of the trajectory once we have decided on the latent  $z$
- How to estimate entropy and conditional entropy over trajectories?
  - $\rightarrow$  recent work showing promising results, we'll look at this a bit later, when studying entropy based exploration/RFPT

# Reverse Question: Can there be too much MI available?

---

- In case of very open-ended environments, info-theoretic approaches presented so far might keep busy forever
  - Showcase humanoid on floor in many configurations ,but never bothers to stand up...
- How to know what data is relevant to collect?
  - Option 1: passive data / human feedback
  - Option 2: try to more directly measure “skill” in terms of what the neural net policy could learn
    - ASP

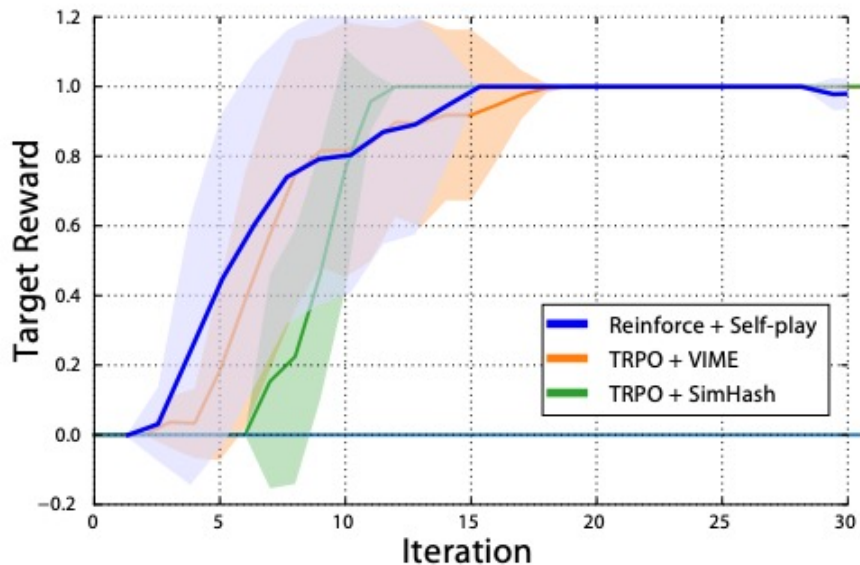
# Asymmetric Self-Play (ASP)

- Key Idea: two-player game, Alice challenges Bob, which encourages Alice to try new things (explore) and Bob to acquire skill
- Algorithm:
  - Alice does roll-out, gives final state as goal to Bob
  - Bob learns goal-conditioned policy with goals set by Alice
  - Alice intrinsic reward =  $\max(0, t_B - t_A)$
  - Bob reward =  $-t_B$

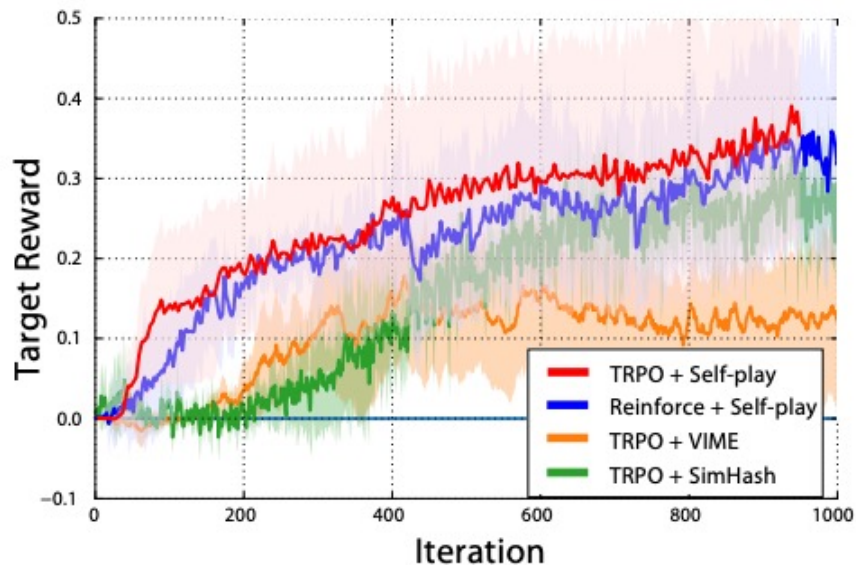
$t_A$  and  $t_B$  are completion times of Alice and Bob

# Asymmetric Self-Play (ASP)

Mountain car

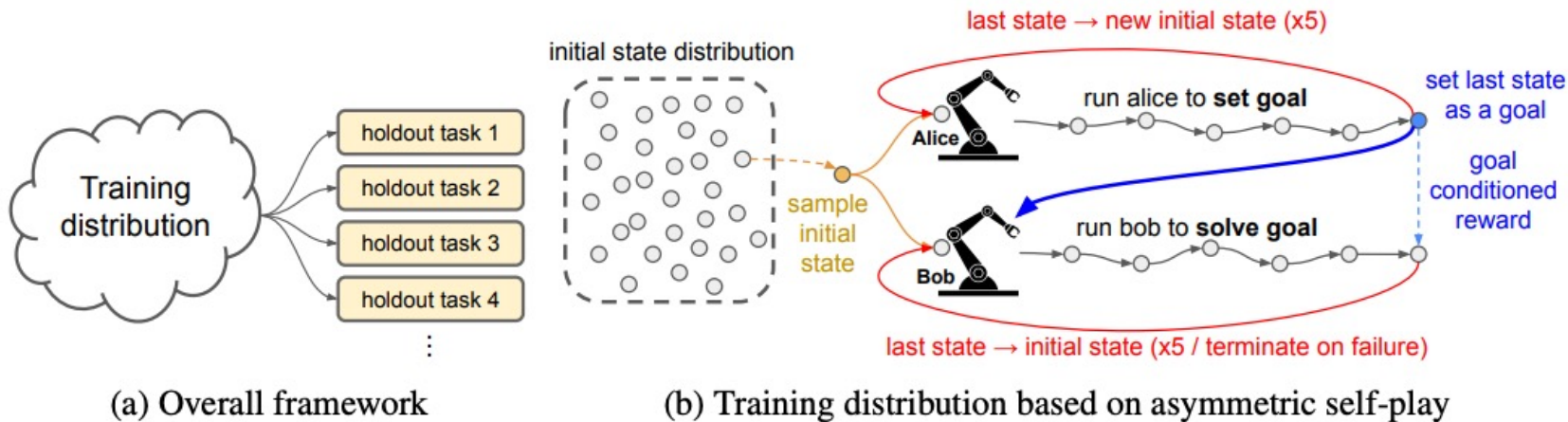


Swimmer Gather



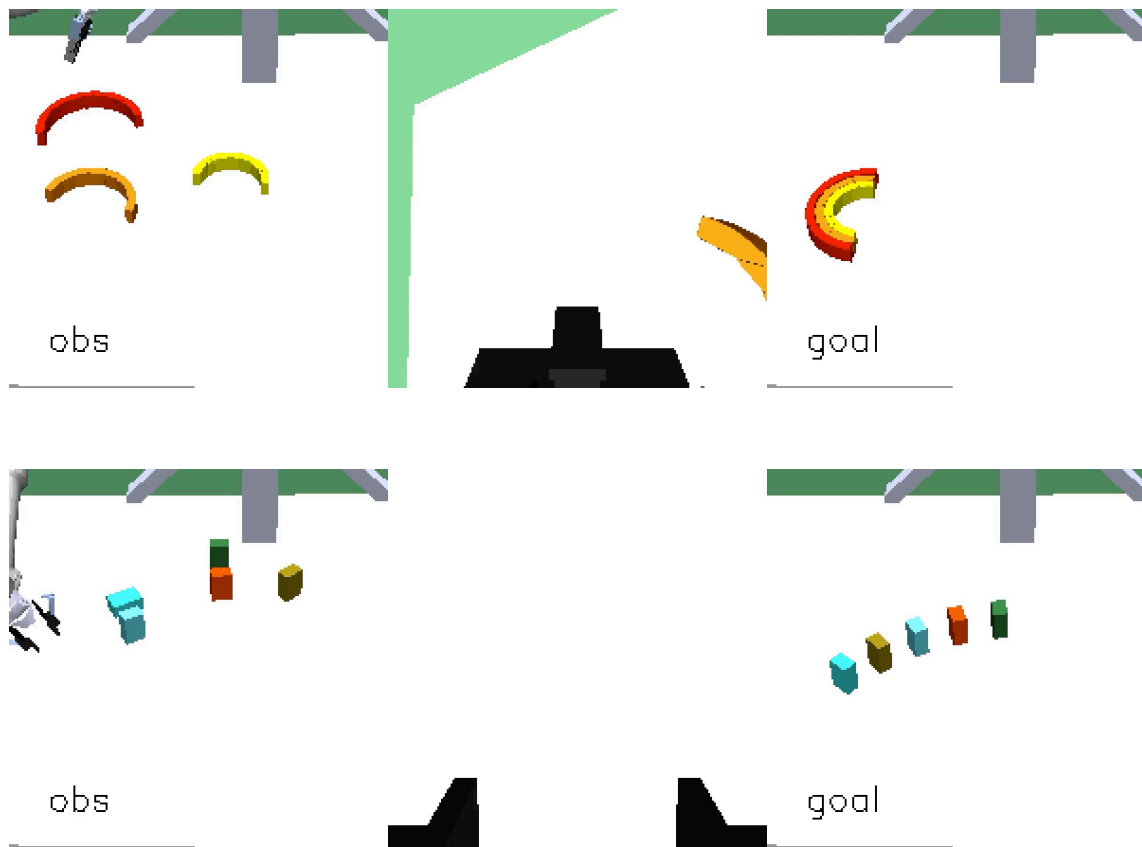


# Scaled-up ASP

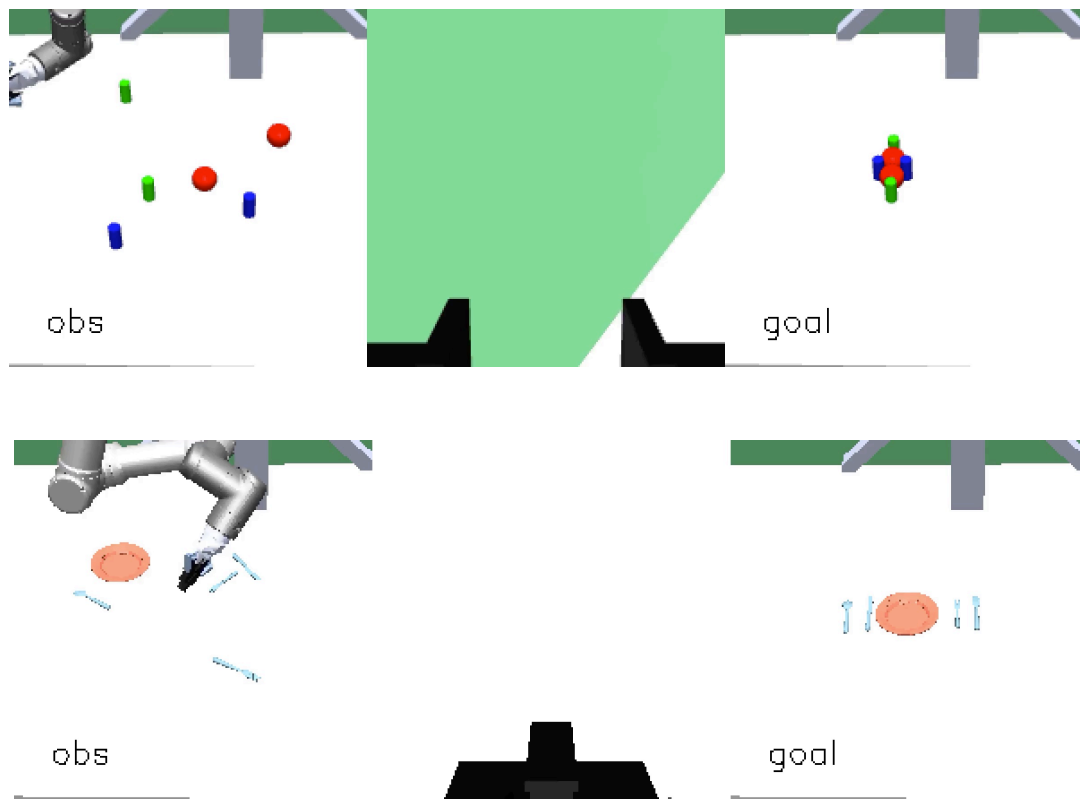


- Goal setting filter: Goals only valid if an object was moved
- 1 billion steps

# Scaled-up ASP: 0-shot generalization



# Scaled-up ASP: 0-shot generalization



# Three Main Types of Intrinsic Reward

---

- Knowledge-based: Surprise / unpredictability / how much learned about world from experience
- Competence-based: Empowerment / Skills
- **Data-based: Entropy (i.e. coverage) of data collected**

# Tabular MDPs: Count-based Exploration Bonus

Model-Based Interval Estimation Exploration Bonus (MBIE-EB)

$$V(x) = \max_{a \in \mathcal{A}} \left[ \hat{R}(x, a) + \gamma \mathbb{E}_{\hat{P}} [V(x')] + \beta N(x, a)^{-1/2} \right]$$

# High-D Spaces: Pseudo-Counts

- Bellemare et al 2016: Unifying Count-based Exploration and Intrinsic Motivation
  - learn  $P(s)$ , change in  $P(s)$  can be connected to a pseudo-count
- Ostrovski et al 2017: Count-Based Exploration with Neural Density Models
  - Similar to above, significantly improved performance thanks to better density models
- Tang et al 2017: #exploration
  - NN hashes high-D observation into lower-dimensional space, use counts there
- Burda et al 2018: RND
  - Train a neural net to match the classification of a random neural net on observations encountered; as long as discrepancy on a new observation, count considered low

# Density-based Pseudo-Counts

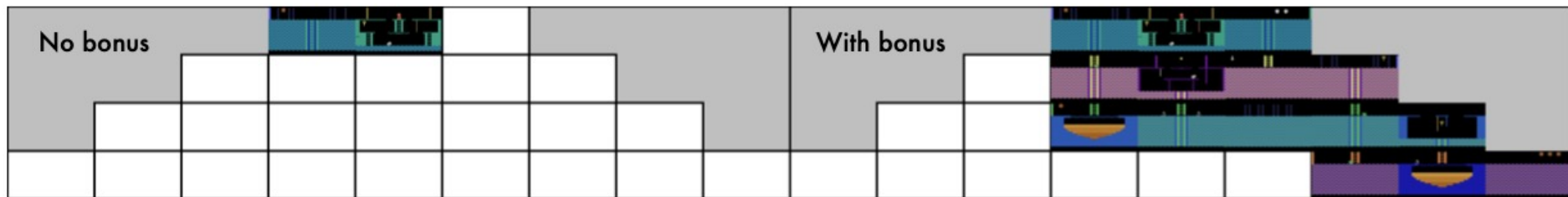


Figure 3: “Known world” of a DQN agent trained for 50 million frames with (**right**) and without (**left**) count-based exploration bonuses, in **MONTEZUMA’S REVENGE**.

# Density-based Pseudo-Counts

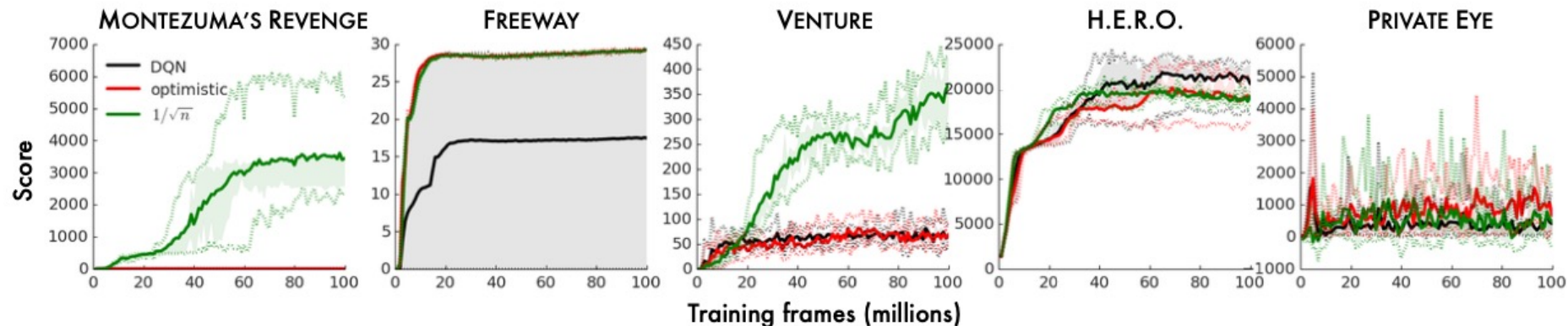


Figure 2: Average training score with and without exploration bonus or optimistic initialization in 5 Atari 2600 games. Shaded areas denote inter-quartile range, dotted lines show min/max scores.



# Directly Optimizing Entropy of Data Collected

- Incentivizing exploration by introducing intrinsic rewards based on a measure of state novelty
- **State entropy** as intrinsic reward

$$r^{\text{intrinsic}} = \mathcal{H}(s) = -\mathbb{E}_{s \sim p(s)} [\log p(s)]$$

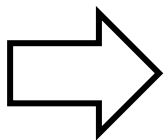
- Maximizing state entropy  $\approx$  good state coverage

# Directly Optimizing Entropy of Data Collected

- Incentivizing exploration by introducing intrinsic rewards based on a measure of state novelty
- **State entropy** as intrinsic reward

$$r^{\text{intrinsic}} = \mathcal{H}(s) = -\mathbb{E}_{s \sim p(s)} [\log p(s)]$$

- Maximizing state entropy  $\approx$  good state coverage

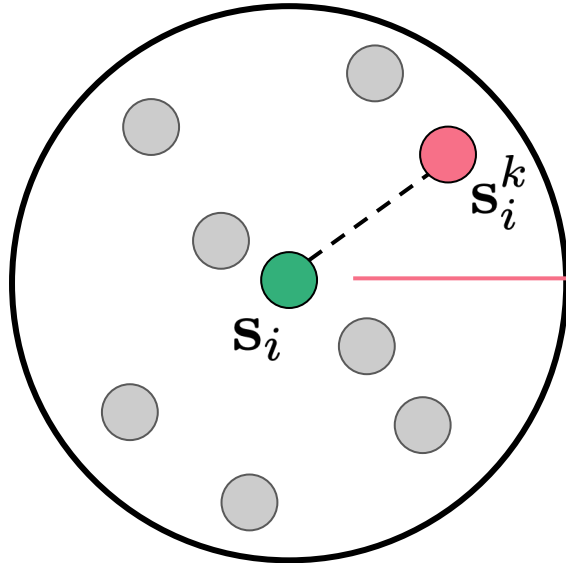


**Measuring state entropy is intractable** to compute in most setting

# K-Nearest-Neighbor Entropy Estimator

- K-nearest entropy estimator

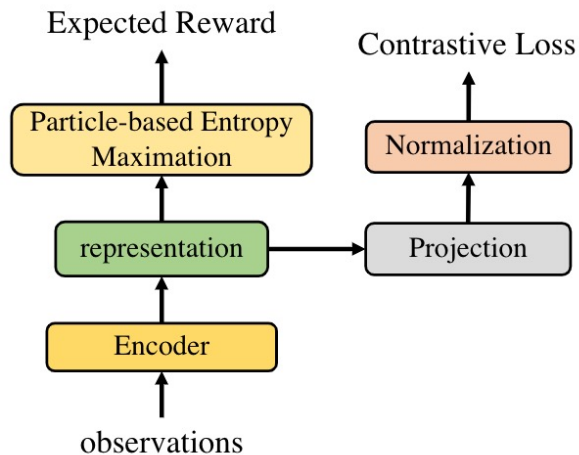
$$\mathcal{H}(s) = -\mathbb{E}_{s \sim p(s)} [\log p(s)]$$



$$\hat{\mathcal{H}}(s) \propto \sum_i \log(\|s_i - s_i^k\|)$$

- Distribution  $\rightarrow$  Store N number of visited states
- Compute the distance between each state and its K-NN

# APT: Active Pre-Training



---

**Algorithm 1:** Training APT

---

Randomly Initialize  $f$  encoder

Randomly Initialize  $\pi$  and  $Q$  networks

**for**  $e := 1, \infty$  **do**

**for**  $t := 1, T$  **do**

    Receive observation  $s_t$  from environment

    Take action  $a_t \sim \pi(\cdot | s_t)$ , receive observation  $s_{t+1}$  and  $\mathcal{R}_t$  from environment

$\mathcal{D} \leftarrow \mathcal{D} \cup (s_t, a_t, \mathcal{R}_t, s'_t)$

$\{(s_i, a_i, \mathcal{R}_i, s'_i)\}_{i=1}^N \sim \mathcal{D}$

    // sample a mini batch

    Train neural encoder  $f$  on mini batch

    // representation learning

**for each**  $i = 1..N$  **do**

$a'_i \sim \pi(\cdot | s'_i)$

$\hat{Q}_i = Q_{\theta'}(s'_i, a'_i)$

      Compute  $r_{\text{APT}}$  with equation (5)

      // particle-based entropy reward

$y_i \leftarrow r_{\text{APT}} + \gamma \hat{Q}_i$

**end**

$loss_Q = \sum_i (Q(s_i, a_i) - y_i)^2$

    Gradient descent step on  $Q$  and  $\pi$

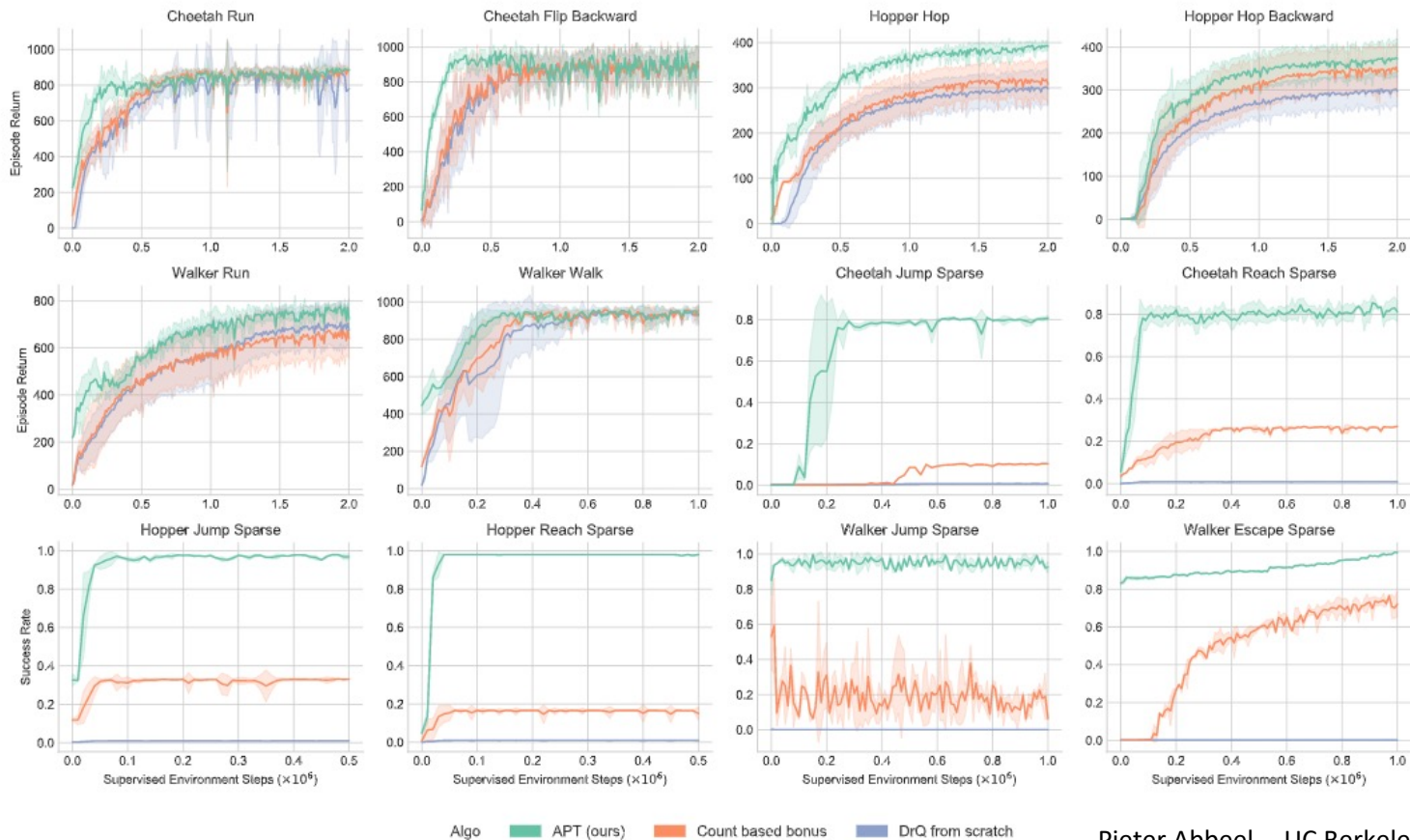
    // standard Q-learning

**end**

**end**

---

# Experiments: DM Control Suite

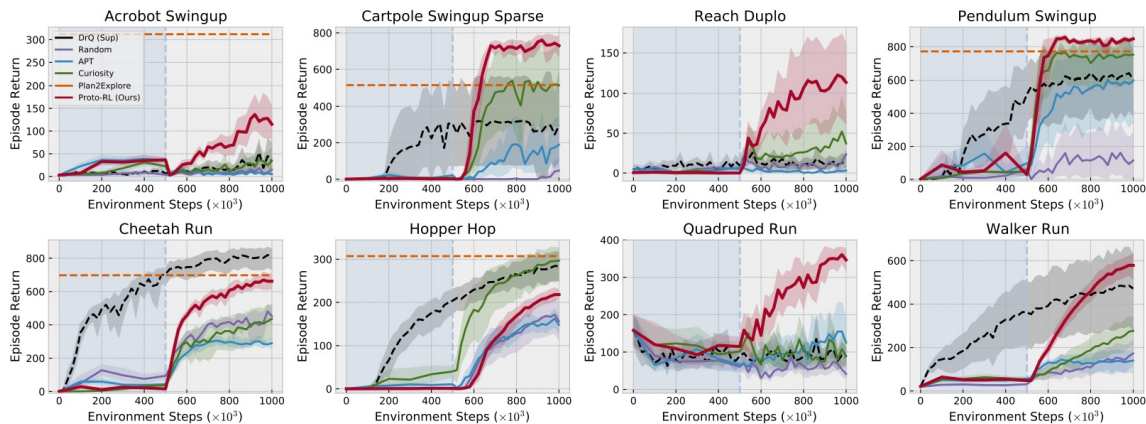
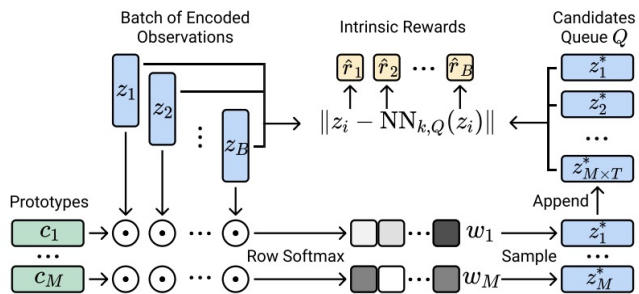


# Experiments: Atari

Game	Random	Human	SimPLe	DER	CURL	DrQ	SPR	VISR	APT (ours)
Alien	227.8	7127.7	616.9	739.9	558.2	771.2	801.5	364.4	<b>2614.8</b>
Amidar	5.8	1719.5	88.0	188.6	142.1	102.8	176.3	186.0	<b>211.5</b>
Assault	222.4	742.0	527.2	431.2	600.6	452.4	571.0	<b>12091.1</b>	891.5
Asterix	210.0	8503.3	1128.3	470.8	734.5	603.5	977.8	<b>6216.7</b>	185.5
Bank Heist	14.2	753.1	34.2	51.0	131.6	168.9	380.9	71.3	<b>416.7</b>
BattleZone	2360.0	37187.5	5184.4	10124.6	14870.0	12954.0	<b>16651.0</b>	7072.7	7065.1
Boxing	0.1	12.1	9.1	0.2	1.2	6.0	<b>35.8</b>	13.4	21.3
Breakout	1.7	30.5	16.4	1.9	4.9	16.1	17.1	<b>17.9</b>	10.9
ChopperCommand	811.0	7387.8	<b>1246.9</b>	861.8	1058.5	780.3	974.8	800.8	317.0
Crazy Climber	10780.5	23829.4	<b>62583.6</b>	16185.2	12146.5	20516.5	42923.6	49373.9	44128.0
Demon Attack	107805	35829.4	62583.6	16185.3	12146.5	20516.5	42923.6	<b>8994.9</b>	5071.8
Freeway	0.0	29.6	20.3	27.9	26.7	9.8	24.4	-12.1	<b>29.9</b>
Frostbite	65.2	4334.7	254.7	866.8	1181.3	331.1	<b>1821.5</b>	230.9	1796.1
Gopher	257.6	2412.5	771.0	349.5	669.3	636.3	715.2	498.6	<b>2590.4</b>
Hero	1027.0	30826.4	2656.6	6857.0	6279.3	3736.3	<b>7019.2</b>	663.5	6789.1
Jamesbond	29.0	302.8	125.3	301.6	471.0	236.0	365.4	<b>484.4</b>	356.1
Kangaroo	52.0	3035.0	323.1	779.3	872.5	940.6	<b>3276.4</b>	1761.9	412.0
Krull	1598.0	2665.5	<b>4539.9</b>	2851.5	4229.6	4018.1	2688.9	3142.5	2312.0
Kung Fu Master	258.5	22736.3	17257.2	14346.1	14307.8	9111.0	13192.7	16754.9	<b>17357.0</b>
Ms Pacman	307.3	6951.6	1480.0	1204.1	1465.5	960.5	1313.2	558.5	<b>2827.1</b>
Pong	-20.7	14.6	<b>12.8</b>	-19.3	-16.5	-8.5	-5.9	-26.2	-8.0
Private Eye	24.9	69571.3	58.3	97.8	218.4	-13.6	<b>124.0</b>	98.3	96.1
Qbert	163.9	13455.0	1288.8	1152.9	1042.4	854.4	669.1	666.3	<b>17671.2</b>
Road Runner	11.5	7845.0	5640.6	9600.0	5661.0	8895.1	<b>14220.5</b>	6146.7	4782.1
Seaquest	68.4	42054.7	683.3	354.1	384.5	301.2	583.1	706.6	<b>2116.7</b>
Up N Down	533.4	11693.2	3350.3	2877.4	2955.2	3180.8	<b>28138.5</b>	10037.6	8289.4
Mean HNS	0.000	1.000	44.3	28.5	38.1	35.7	<b>70.4</b>	64.31	69.55
Median HNS	0.000	1.000	14.4	16.1	17.5	26.8	41.5	12.36	<b>47.50</b>
# Superhuman	0	N/A	2	2	2	2	<b>7</b>	6	<b>7</b>

# How about size of replay buffer for entropy estimates?

→ Keep around cluster representatives for entropy estimation



# Active Pre-Training with Successor Features (APS)

---

- Similar to how VISR added Successor Features to DIAYN for faster adaptation
- APS add Successor Features to APT



# APS on Atari

Game	Random	Human	SimPLe	DER	CURL	DrQ	SPR	VISR	APT	APS (ours)
Alien	227.8	7127.7	616.9	739.9	558.2	771.2	801.5	364.4	<b>2614.8</b>	934.9
Amidar	5.8	1719.5	88.0	188.6	142.1	102.8	176.3	186.0	<b>211.5</b>	178.4
Assault	222.4	742.0	527.2	431.2	600.6	452.4	571.0	<b>12091.1</b>	891.5	413.3
Asterix	210.0	8503.3	1128.3	470.8	734.5	603.5	977.8	<b>6216.7</b>	185.5	1159.7
Bank Heist	14.2	753.1	34.2	51.0	131.6	168.9	380.9	71.3	<b>416.7</b>	262.7
BattleZone	2360.0	37187.5	5184.4	10124.6	14870.0	12954.0	16651.0	7072.7	7065.1	<b>26920.1</b>
Boxing	0.1	12.1	9.1	0.2	1.2	6.0	35.8	13.4	21.3	<b>36.3</b>
Breakout	1.7	30.5	16.4	1.9	4.9	16.1	17.1	17.9	10.9	<b>19.1</b>
ChopperCommand	811.0	7387.8	1246.9	861.8	1058.5	780.3	974.8	800.8	317.0	<b>2517.0</b>
Crazy Climber	10780.5	23829.4	62583.6	16185.2	12146.5	20516.5	42923.6	49373.9	44128.0	<b>67328.1</b>
Demon Attack	107805	35829.4	62583.6	16185.3	12146.5	20516.5	42923.6	<b>8994.9</b>	5071.8	7989.0
Freeway	0.0	29.6	20.3	27.9	26.7	9.8	24.4	-12.1	<b>29.9</b>	27.1
Frostbite	65.2	4334.7	254.7	866.8	1181.3	331.1	<b>1821.5</b>	230.9	1796.1	496.5
Gopher	257.6	2412.5	771.0	349.5	669.3	636.3	715.2	498.6	<b>2590.4</b>	2386.5
Hero	1027.0	30826.4	2656.6	6857.0	6279.3	3736.3	7019.2	663.5	6789.1	<b>12189.3</b>
Jamesbond	29.0	302.8	125.3	301.6	471.0	236.0	365.4	484.4	356.1	<b>622.3</b>
Kangaroo	52.0	3035.0	323.1	779.3	872.5	940.6	3276.4	1761.9	412.0	<b>5280.1</b>
Krull	1598.0	2665.5	<b>4539.9</b>	2851.5	4229.6	4018.1	2688.9	3142.5	2312.0	4496.0
Kung Fu Master	258.5	22736.3	17257.2	14346.1	14307.8	9111.0	13192.7	16754.9	17357.0	<b>22412.0</b>
Ms Pacman	307.3	6951.6	1480.0	1204.1	1465.5	960.5	1313.2	558.5	<b>2827.1</b>	2092.3
Pong	-20.7	14.6	<b>12.8</b>	-19.3	-16.5	-8.5	-5.9	-26.2	-8.0	12.5
Private Eye	24.9	69571.3	58.3	97.8	218.4	-13.6	<b>124.0</b>	98.3	96.1	117.9
Qbert	163.9	13455.0	1288.8	1152.9	1042.4	854.4	669.1	666.3	17671.2	<b>19271.4</b>
Road Runner	11.5	7845.0	5640.6	9600.0	5661.0	8895.1	<b>14220.5</b>	6146.7	4782.1	5919.0
Seaquest	68.4	42054.7	683.3	354.1	384.5	301.2	583.1	706.6	2116.7	<b>4209.7</b>
Up N Down	533.4	11693.2	3350.3	2877.4	2955.2	3180.8	<b>28138.5</b>	10037.6	8289.4	4911.9
Mean Human-Norm'd	0.000	1.000	44.3	28.5	38.1	35.7	70.4	64.31	69.55	<b>99.04</b>
Median Human-Norm'd	0.000	1.000	14.4	16.1	17.5	26.8	41.5	12.36	47.50	<b>58.80</b>
# Superhuman	0	N/A	2	2	2	2	7	6	7	<b>8</b>

# Never Give Up

- Main idea:
  - Short-term exploration through particle-based entropy
  - Long-term exploration through RND

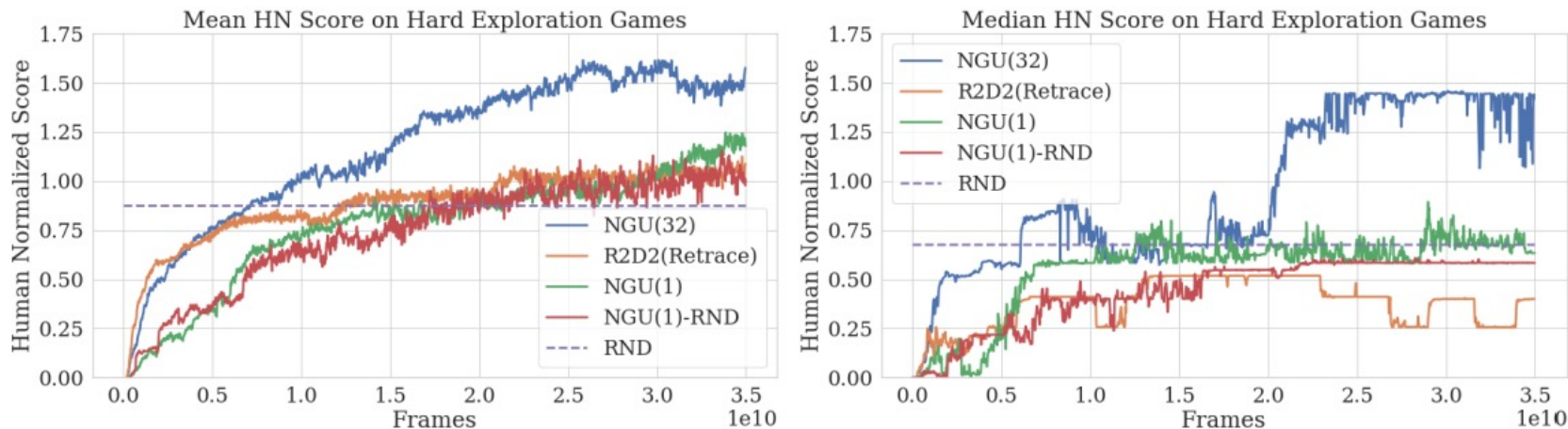


Figure 4: Human Normalized Scores on the 6 hard exploration games.

# Some Theory

## Provably Efficient Maximum Entropy Exploration

Elad Hazan<sup>1,2</sup>   Sham M. Kakade<sup>1,3,4</sup>   Karan Singh<sup>1,2</sup>   Abby Van Soest<sup>2</sup>

<sup>1</sup> Google AI Princeton

<sup>2</sup> Department of Computer Science, Princeton University

<sup>3</sup> Allen School of Computer Science and Engineering, University of Washington

<sup>4</sup> Department of Statistics, University of Washington

{ehazan,karans,asoest}@princeton.edu, sham@cs.washington.edu

## Kinematic State Abstraction and Provably Efficient Rich-Observation Reinforcement Learning

Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford

Microsoft Research, New York, NY

# Outline

---

- Problem Motivation
- Baseline RL Algorithms Refresher
- Intrinsic Rewards for Reward-Free Pre-Training and Exploration
- ***Algorithmic Approaches to Exploration*** (*can complement intrinsic reward RFPT!*)
- Algorithmic Approaches to Reward-Free Pre-Training

# Key Idea: Optimism in Face of Uncertainty

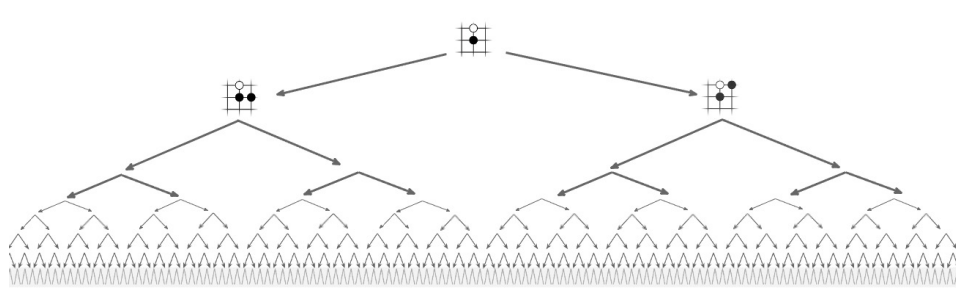
- Key Idea permeating most works: Optimism in Face of Uncertainty
  - If haven't visited a state, let's assume it might have high reward, until we experience otherwise
- Classic references (experiments largely tabular / linear)

Lai and Robbins' Upper Confidence Bounds (UCB) 1985;  
Sutton's Dyna, 1990;  
Schmidhuber's Curiosity 1991;  
Kaelbling's Interval Exploration 1993;  
Moore & Atkeson's Prioritized Sweeping 1993;

Kaelbling, et al, RL Survey 1996;  
Kearns and Singh, E3 2002;  
Brafman and Tennenholtz, RMax 2002  
Peter Auer's UCB regret bounds, 2002;  
Strehl and Littman, Model-based Interval Estimation (MBIE), 2008

# Algorithmic Exploration: MCTS

AlphaGo



- MCTS Assumes: (i) model-based RL OR (ii) sim with resets to any previously experienced state

- PUCB value:  $u(a) = v(a) + p(a) \cdot pb_c$

$$pb_c = \frac{\sqrt{\text{visit}_{\text{parent}}}}{\text{visit}_{\text{child}} + 1} \cdot \left( \log \left( \frac{\text{visit}_{\text{parent}} + pb_{\text{base}} + 1}{pb_{\text{base}}} \right) + pb_{\text{init}} \right)$$

→ Favor less frequently visited children in the tree

# Algorithmic Exploration: Q Ensembles

---

Q: How to more directly represent posterior over value functions?

A: DQN with an ensemble of Q functions

Q: How to use this posterior?

A1: Posterior Sampling / Thompson Sampling / “Bootstrap”

A2: Create Upper Confidence Bound (UCB)

# Algorithmic Exploration: Q Ensembles: Thompson Sampling / Bootstrap

Key idea: Rather than independent action selection in each step, use the same member of the Q-ensemble for the entire roll-out

→ more consistent behavior, which aids exploration



Figure 3: Scalable environments that requires deep exploration.



# Algorithmic Exploration: Q Ensembles: Thompson Sampling / Bootstrap

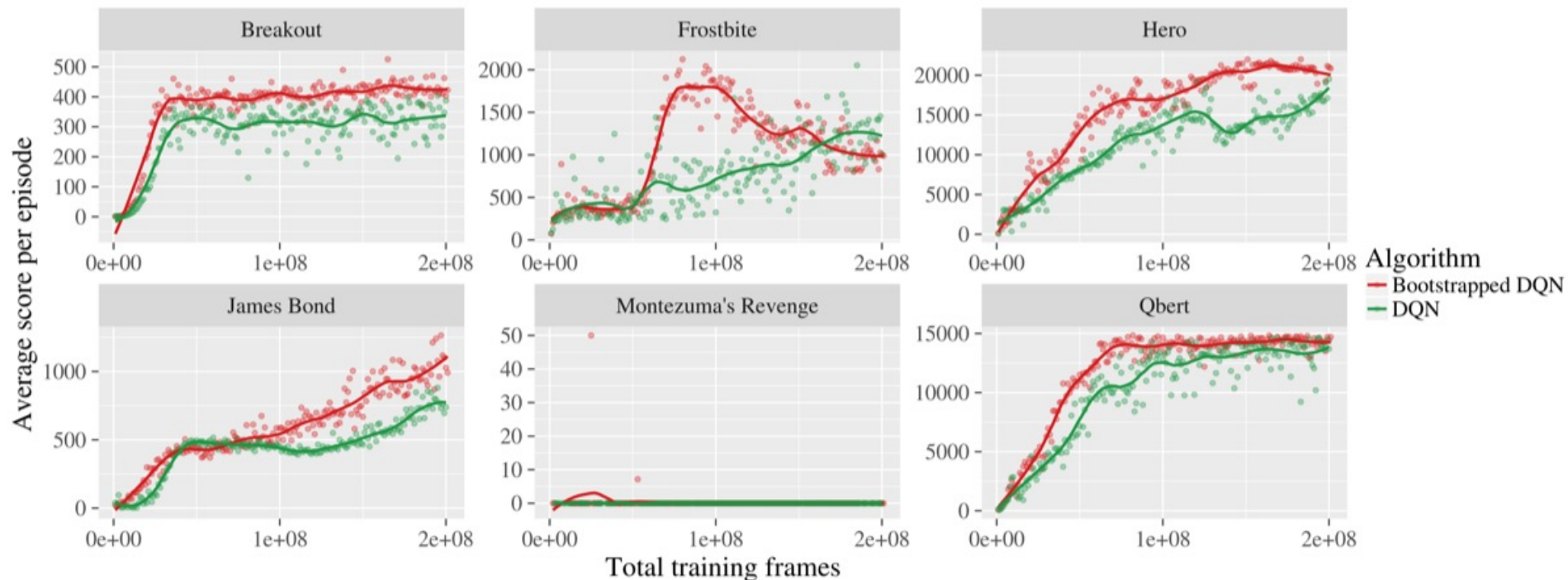


Figure 6: Bootstrapped DQN drives more efficient exploration.

# Algorithmic Exploration: Q-Ensembles: Q-UCB

- How to add UCB to deep Q-learning (DQN)?

---

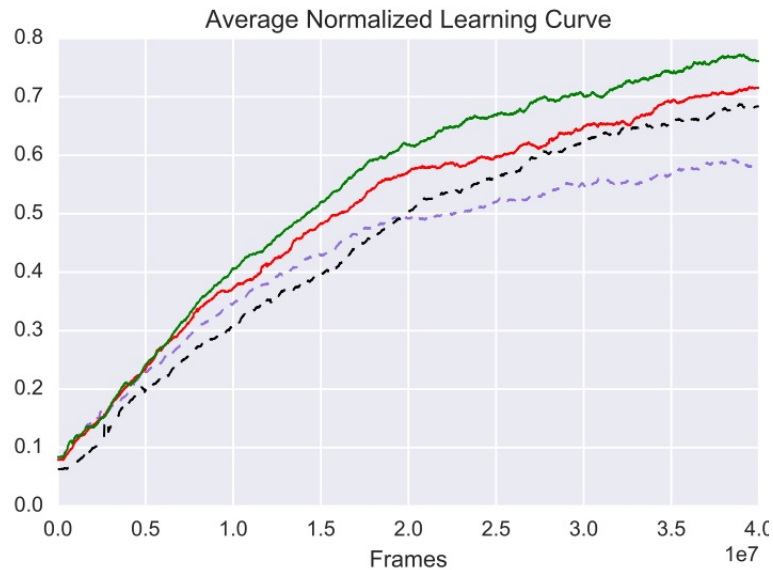
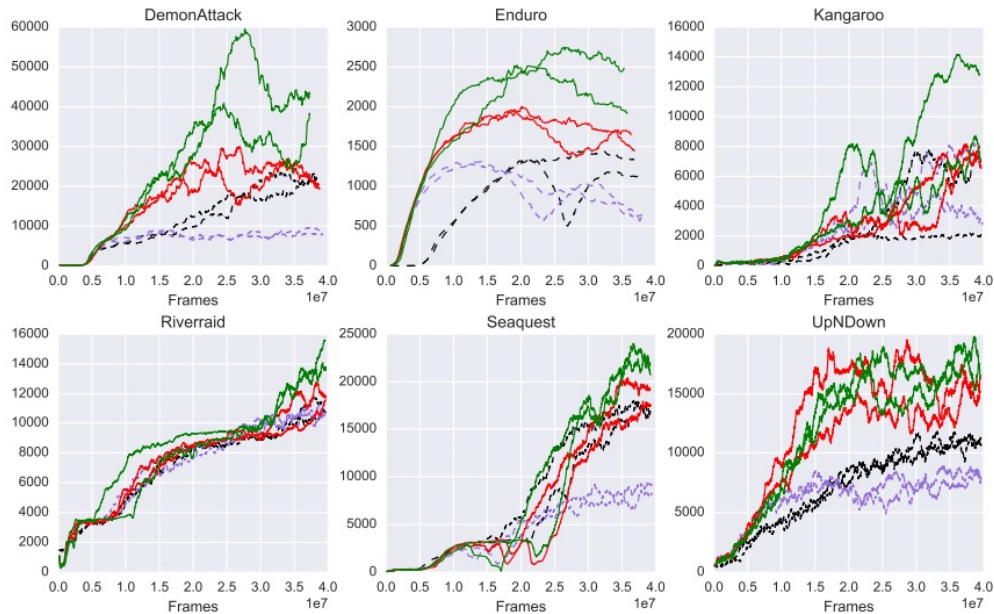
**Algorithm 2** UCB Exploration with  $Q$ -Ensembles

---

- 1: **Input:** Value function networks  $Q$  with  $K$  outputs  $\{Q_k\}_{k=1}^K$ . Hyperparameter  $\lambda$ .
  - 2: Let  $B$  be a replay buffer storing experience for training.
  - 3: **for** each episode **do**
  - 4:     Obtain initial state from environment  $s_0$
  - 5:     **for** step  $t = 1, \dots$  until end of episode **do**
  - 6:         Pick an action according to  $a_t \in \operatorname{argmax}_a \{\tilde{\mu}(s_t, a) + \lambda \cdot \tilde{\sigma}(s_t, a)\}$
  - 7:         Receive state  $s_{t+1}$  and reward  $r_t$  from environment, having taken action  $a_t$
  - 8:         Add  $(s_t, a_t, r_t, s_{t+1})$  to replay buffer  $B$
  - 9:         At learning interval, sample random minibatch and update  $\{Q_k\}$  according to (12)
  - 10:     **end for**
  - 11: **end for**
- 

mean and std across  
ensemble of Q-functions

# Algorithmic Exploration: Q-Ensembles: Q-UCB



- - bootstrapped dqn
- ucb exploration
- ensemble voting
- - double dqn

# Algorithmic Exploration: V-Ensemble + MPC: POLO

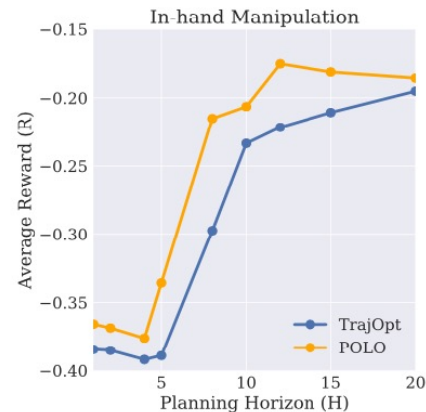
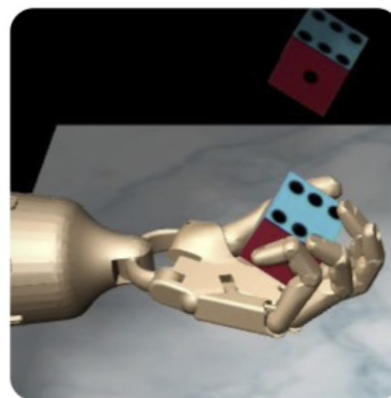
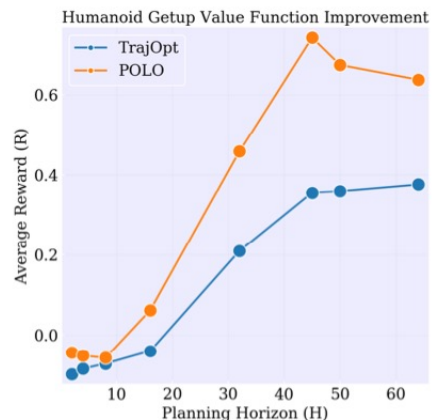
- Key idea: leverage value function ensemble to *plan* for *exploration*

$$\hat{\pi}_{MPC}(s) := \arg \max_{\pi_{0:H-1}} \mathbb{E} \left[ \sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) + \gamma^H r_f(s_H) \right]$$

$$\hat{V}(s) = \log \left( \sum_{k=1}^K \exp \left( \kappa \hat{V}_{\theta_k}(s) \right) \right)$$

MPC with optimistic final state value through softmax over ensemble

# Algorithmic Exploration: V-Ensemble + MPC: POLO



# Algorithmic Exploration: Q-Ensembles: SUNRISE

## Key Contributions beyond Q-UCB:

- Q: If we use Q-UCB, can the off-policy action selection affect learning stability?
- A: Yes. (see also recent work on stabilizing offline RL, investigating same)
  
- Q: If we have a Q-ensemble, can we leverage it to stabilize/improve the Q-learning update?
- A: Yes, and it can more generally mitigate noise propagation that happens in regular DQN/Q-learning (which is especially present when Q-learning is run off-policy\*)


\*Related: Chen, Schulman, Abbeel, Boltzmann Exploration Q-learning is secretly on-policy actor critic...

# Algorithmic Exploration: Q-Ensembles: SUNRISE

- Error propagation issue in Q-learning

$$Q(s_t, a_t) \leftarrow r_t + \gamma \max_a Q(s_{t+1}, a)$$

unseen (s,a) → high error



error propagates

- Reweight each term in Bellman backup loss

$$w(s, a) \left( Q(s, a) - [r(s, a) + \gamma \hat{Q}(s', a')] \right)^2$$

confidence score about target value based on variance of Q-ensemble

# Algorithmic Exploration: Q-Ensembles: SUNRISE

- Performance on DeepMind Control Suite at 100K and 500K environment steps

500K step	PlaNet [16]	Dreamer [17]	SLAC [31]	CURL [41]	DrQ [25]	RAD [30]	SUNRISE
Finger-spin	561 $\pm$ 284	796 $\pm$ 183	673 $\pm$ 92	926 $\pm$ 45	938 $\pm$ 103	975 $\pm$ 16	<b>983</b> $\pm$ 1
Cartpole-swing	475 $\pm$ 71	762 $\pm$ 27	-	845 $\pm$ 45	868 $\pm$ 10	873 $\pm$ 3	<b>876</b> $\pm$ 4
Reacher-easy	210 $\pm$ 44	793 $\pm$ 164	-	929 $\pm$ 44	942 $\pm$ 71	916 $\pm$ 49	<b>982</b> $\pm$ 3
Cheetah-run	305 $\pm$ 131	570 $\pm$ 253	640 $\pm$ 19	518 $\pm$ 28	660 $\pm$ 96	624 $\pm$ 10	<b>678</b> $\pm$ 46
Walker-walk	351 $\pm$ 58	897 $\pm$ 49	842 $\pm$ 51	902 $\pm$ 43	921 $\pm$ 45	938 $\pm$ 9	<b>953</b> $\pm$ 13
Cup-catch	460 $\pm$ 380	879 $\pm$ 87	852 $\pm$ 71	959 $\pm$ 27	963 $\pm$ 9	966 $\pm$ 9	<b>969</b> $\pm$ 5
100K step							
Finger-spin	136 $\pm$ 216	341 $\pm$ 70	693 $\pm$ 141	767 $\pm$ 56	901 $\pm$ 104	811 $\pm$ 146	<b>905</b> $\pm$ 57
Cartpole-swing	297 $\pm$ 39	326 $\pm$ 27	-	582 $\pm$ 146	<b>759</b> $\pm$ 92	373 $\pm$ 90	591 $\pm$ 55
Reacher-easy	20 $\pm$ 50	314 $\pm$ 155	-	538 $\pm$ 233	601 $\pm$ 213	567 $\pm$ 54	<b>722</b> $\pm$ 50
Cheetah-run	138 $\pm$ 88	235 $\pm$ 137	319 $\pm$ 56	299 $\pm$ 48	344 $\pm$ 67	381 $\pm$ 79	<b>413</b> $\pm$ 35
Walker-walk	224 $\pm$ 48	277 $\pm$ 12	361 $\pm$ 73	403 $\pm$ 24	612 $\pm$ 164	641 $\pm$ 89	<b>667</b> $\pm$ 147
Cup-catch	0 $\pm$ 0	246 $\pm$ 174	512 $\pm$ 110	769 $\pm$ 43	<b>913</b> $\pm$ 53	666 $\pm$ 181	633 $\pm$ 241

- SUNRISE consistently improves the performance of RAD



# Algorithmic Exploration: Q-Ensembles: SUNRISE

- Performance on Atari games at 100K interactions

Game	Human	Random	SimPLe [23]	CURL [41]	Rainbow [47]	SUNRISE
Alien	7127.7	227.8	616.9	558.2	789.0	<b>872.0</b>
Amidar	1719.5	5.8	88.0	<b>142.1</b>	118.5	122.6
Assault	742.0	222.4	527.2	<b>600.6</b>	413.0	594.8
Asterix	8503.3	210.0	<b>1128.3</b>	734.5	533.3	755.0
BankHeist	753.1	14.2	34.2	131.6	97.7	<b>266.7</b>
BattleZone	37187.5	2360.0	5184.4	14870.0	7833.3	<b>15700.0</b>
Boxing	12.1	0.1	<b>9.1</b>	1.2	0.6	6.7
Breakout	30.5	1.7	<b>16.4</b>	4.9	2.3	1.8
ChopperCommand	7387.8	811.0	<b>1246.9</b>	1058.5	590.0	1040.0
CrazyClimber	35829.4	10780.5	<b>62583.6</b>	12146.5	25426.7	22230.0
DemonAttack	1971.0	152.1	208.1	817.6	688.2	<b>919.8</b>
Freeway	29.6	0.0	20.3	26.7	28.7	<b>30.2</b>
Frostbite	4334.7	65.2	254.7	1181.3	1478.3	<b>2026.7</b>
Gopher	2412.5	257.6	771.0	<b>669.3</b>	348.7	654.7
Hero	30826.4	1027.0	2656.6	6279.3	3675.7	<b>8072.5</b>
Jamesbond	302.8	29.0	125.3	<b>471.0</b>	300.0	390.0
Kangaroo	3035.0	52.0	323.1	872.5	1060.0	<b>2000.0</b>
Krull	2665.5	1598.0	<b>4539.9</b>	4229.6	2592.1	3087.2
KungFuMaster	22736.3	258.5	<b>17257.2</b>	14307.8	8600.0	10306.7
MsPacman	6951.6	307.3	1480.0	1465.5	1118.7	<b>1482.3</b>
Pong	14.6	-20.7	<b>12.8</b>	-16.5	-19.0	-19.3
PrivateEye	69571.3	24.9	58.3	<b>218.4</b>	97.8	100.0
Qbert	13455.0	163.9	1288.8	1042.4	646.7	<b>1830.8</b>
RoadRunner	7845.0	11.5	5640.6	5661.0	9923.3	<b>11913.3</b>
Seaquest	42054.7	68.4	<b>683.3</b>	384.5	396.0	570.7
UpNDown	11693.2	533.4	3350.3	2955.2	3816.0	<b>5074.0</b>

Consistently  
outperform Rainbow

SOTA on 13 out of 26  
environments

# Outline

---

- Problem Motivation
- Baseline RL Algorithms Refresher
- Intrinsic Rewards for Reward-Free Pre-Training and Exploration
- Algorithmic Approaches to Exploration (can complement intrinsic reward RFPT!)
- ***Algorithmic Approaches to Reward-Free Pre-Training***

# Algorithmic RFPT: Frontier Approaches

---

- Main idea:
  - keep track of frontier of where the AI Agent has been
  - then set goals near this frontier and/or explicitly revisit past state on frontier and randomly explore from there

# Algorithmic RFPT: Frontier: HER

- Train policy  $\pi(a | s, g)$
- Key idea:
  - To improve learning signal / alleviate exploration needs:  
Hindsight-relabel 50% of goals  $g$  in the replay buffer to match the final achieved state in the corresponding trajectory
  - Leads to natural expansion of goals that can be achieved

[Andrychowicz et al, 2017: Hindsight Experience Replay]

[Schaul et al, 2015: Universal Value Function Approximators]

[Kaelbling 1993: Learning to Achieve Goals]

# Algorithmic RFPT: Frontier: HER

- Train policy  $\pi(a | s, g)$
- Key idea:
  - To improve learning signal / alleviate exploration needs:  
Hindsight-relabel 50% of goals  $g$  in the replay buffer to match the final achieved state in the corresponding trajectory
  - Leads to natural expansion of goals that can be achieved

[Andrychowicz et al, 2017: Hindsight Experience Replay]

[Schaul et al, 2015: Universal Value Function Approximators]

[Kaelbling 1993: Learning to Achieve Goals]

# Quantitative Evaluation

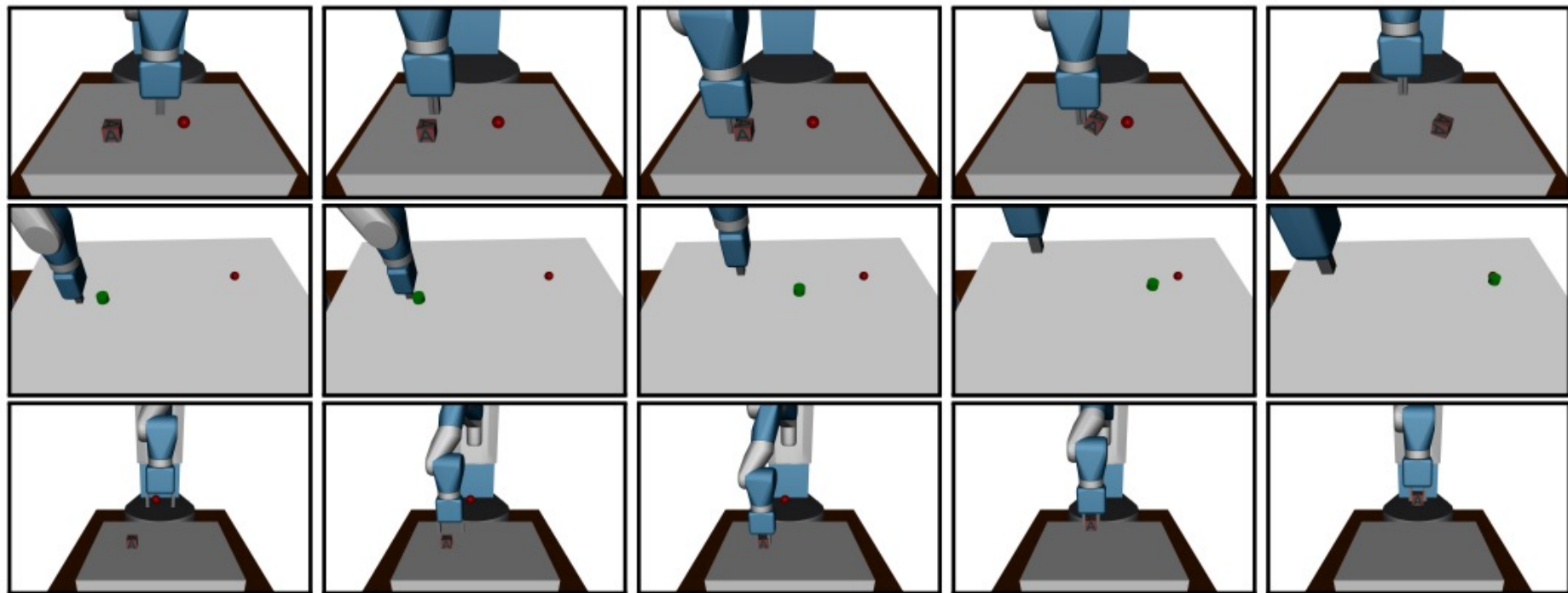
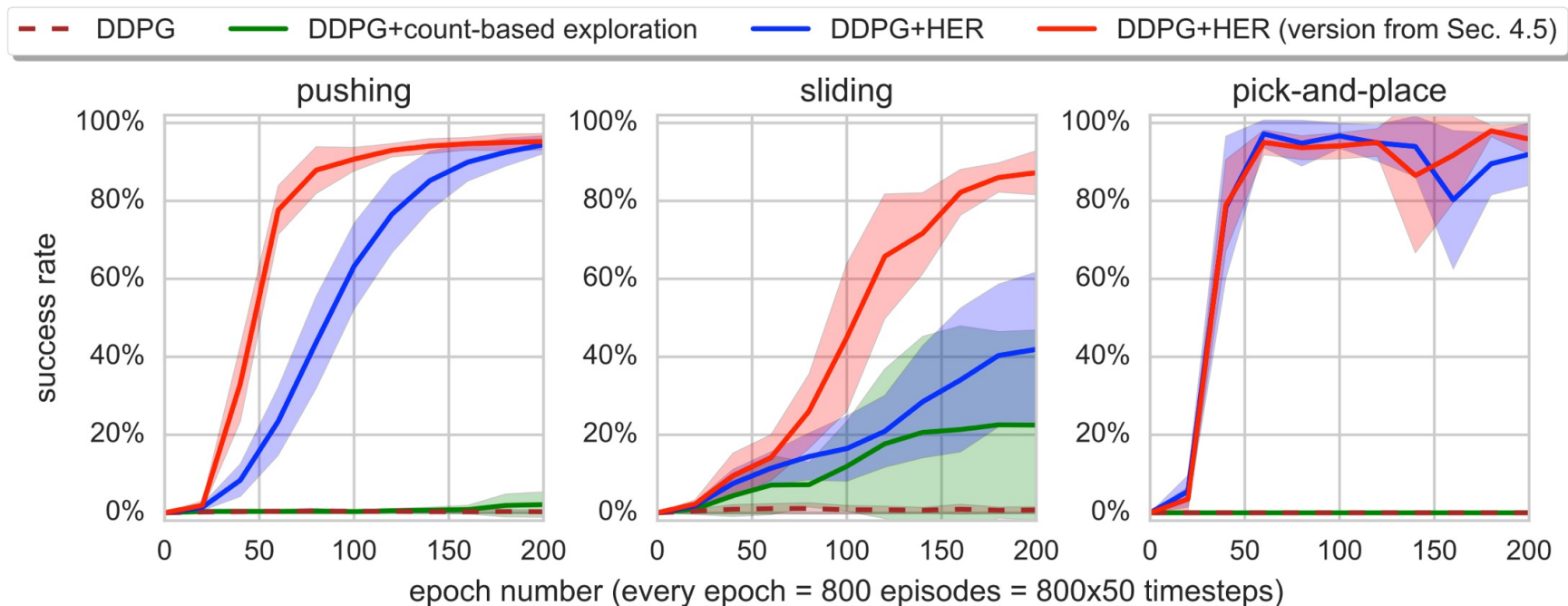


Figure 2: Different tasks: *pushing* (top row), *sliding* (middle row) and *pick-and-place* (bottom row). The red ball denotes the goal position.

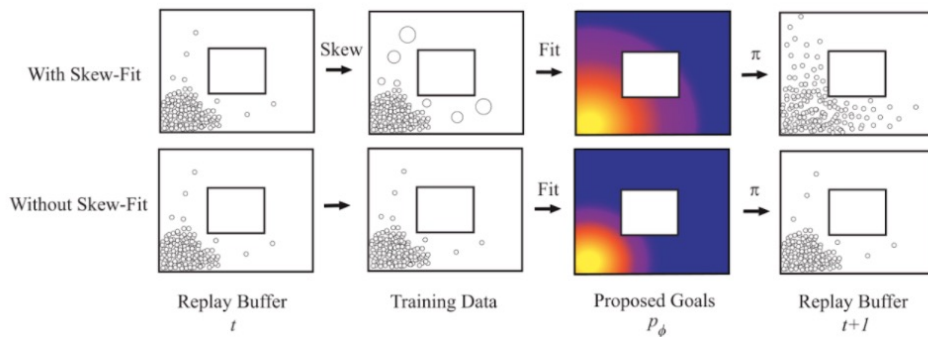
# Quantitative Evaluation



~3s/episode → 1.5 days for 50 epochs, 6 days for 200 epochs

# Algorithmic RFPT: Frontier: RIG / SkewFit

- HER assumes access to underlying state
- RIG adapts HER to Image Inputs
  - operates in latent space of a VAE
  - can use latent space for new goal setting
- SkewFit shifts RIG / HER goal sampling

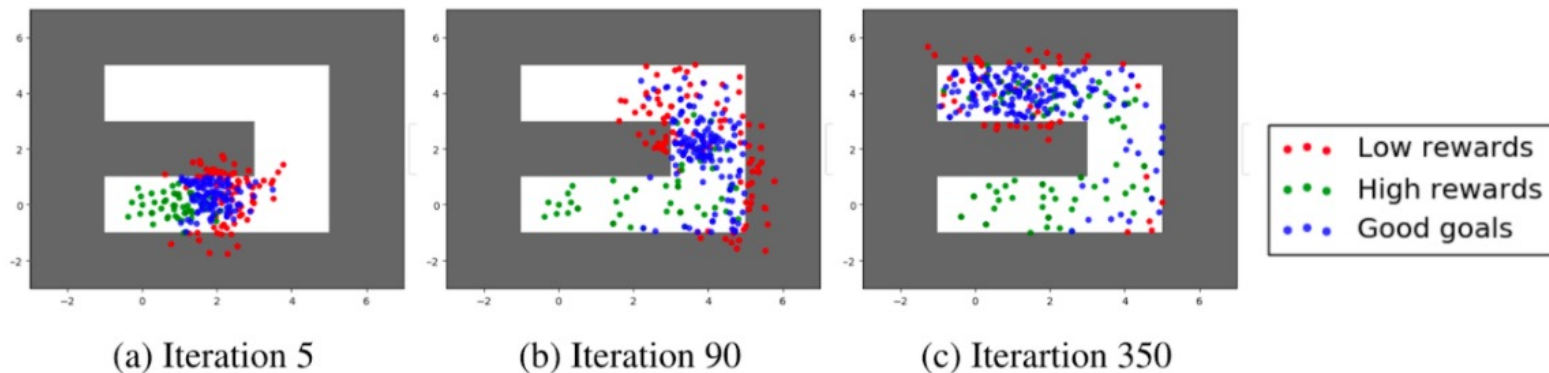




# Algorithmic RFPT: Frontier: GoalGAN

## ■ Key ideas:

- Train goal-conditioned policy / Q-function ( $\sim$ HER)
- Sample Goals of Intermediate Difficulty – goal generator = GAN



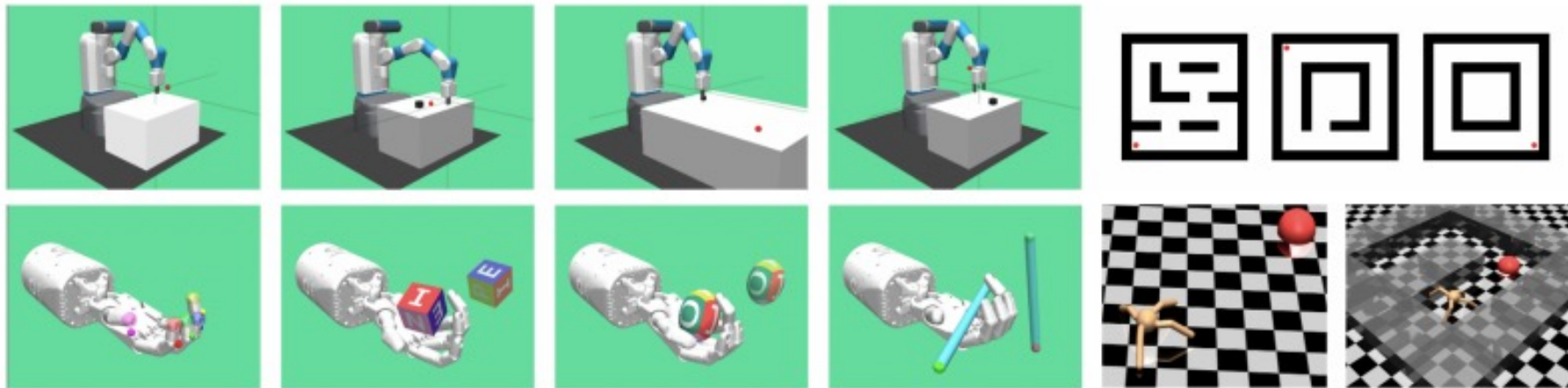
[GoalGAN: Florensa\*, Held\*, Geng\*, Abbeel, 2017: Automatic Goal Generation for RL Agents]

see also: [CURIOUS: Colas, Fournier, Sigaud, Chetouani, Oudeyer, 2019] – considers learning progress (vs. GoalGAN considers learning status)

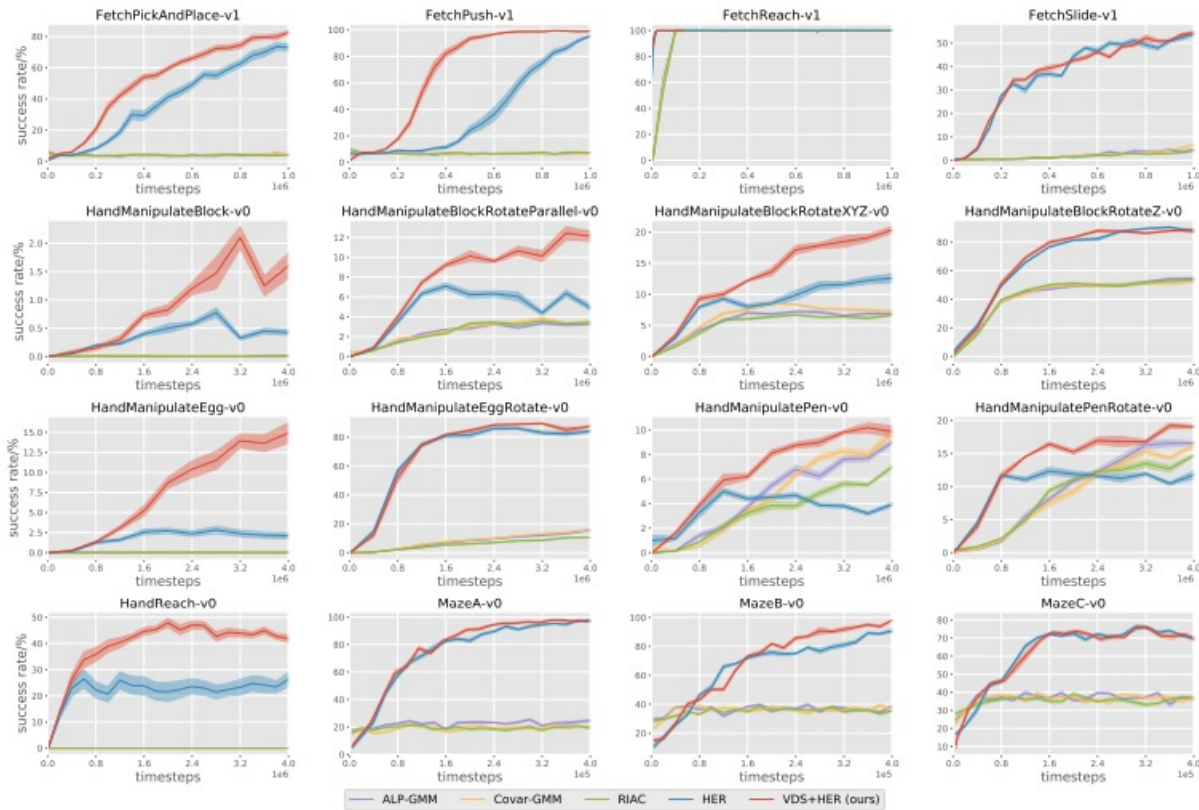
# Algorithmic RFPT: Frontier: Value Disagreement

- Key ideas:

- Train goal-conditioned policy / Q-function ( $\sim$ HER)
- Sample Goals based on Q-Value Disagreement



# Algorithmic RFPT: Frontier: Value Disagreement

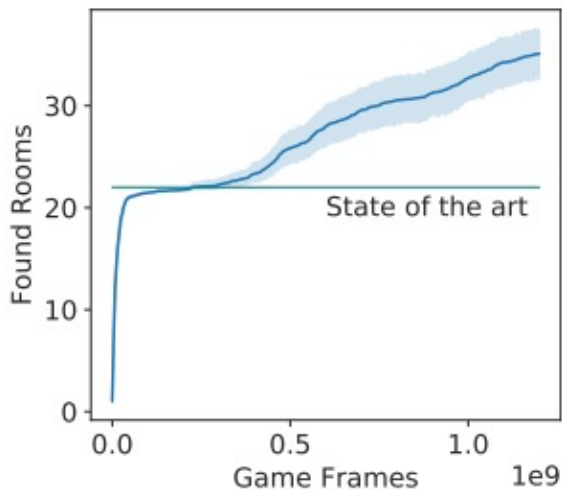


# Algorithmic RFPT: Frontier: GoExplore

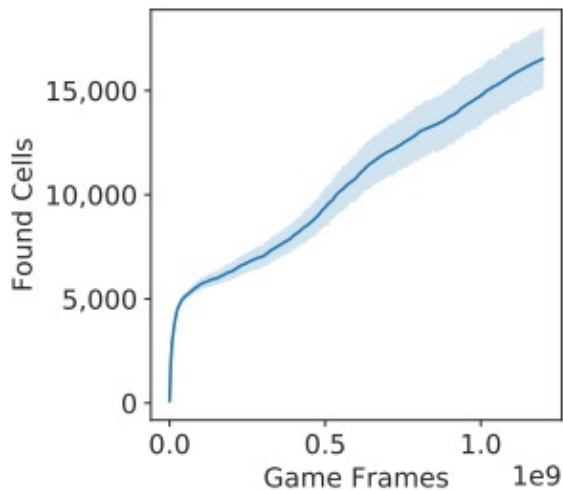
- Key idea:
  - When visiting a state with low visitation count, specifically go to that state again in future roll-outs + randomly explore around that state, which will likely yield new such exploration goal states
  - Original paper: assume ability to deterministically return or access to resets; later versions train a goal conditioned policy
  - Key assumption: reasonable way to divide state space into cells, and not too many cells to be able to explore them all --- done by low-res image in Atari
- Main result: breakthrough/”solve” Atari hard-exploration games

# Algorithmic RFPT: Frontier: GoExplore

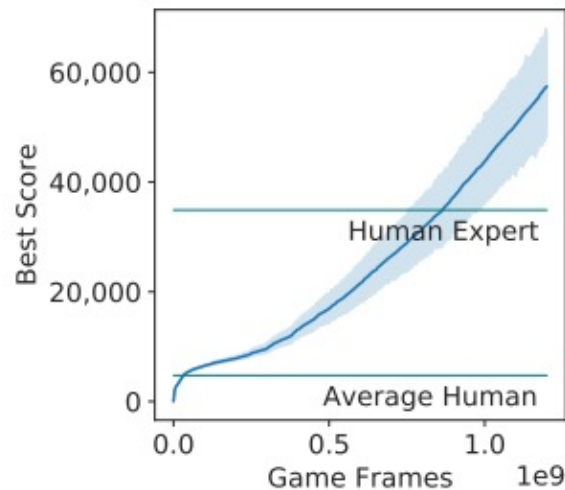
## Montezuma's Revenge Atari Game Performance



(a) Number of rooms found



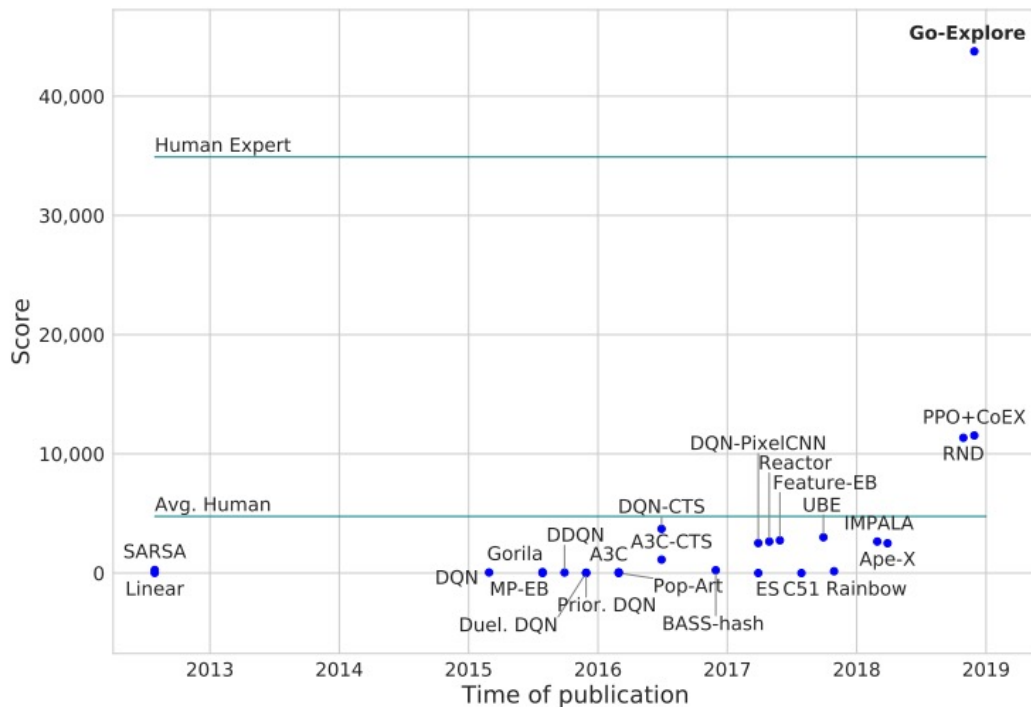
(b) Number of cells found



(c) Maximum score in archive

# Algorithmic RFPT: Frontier: GoExplore

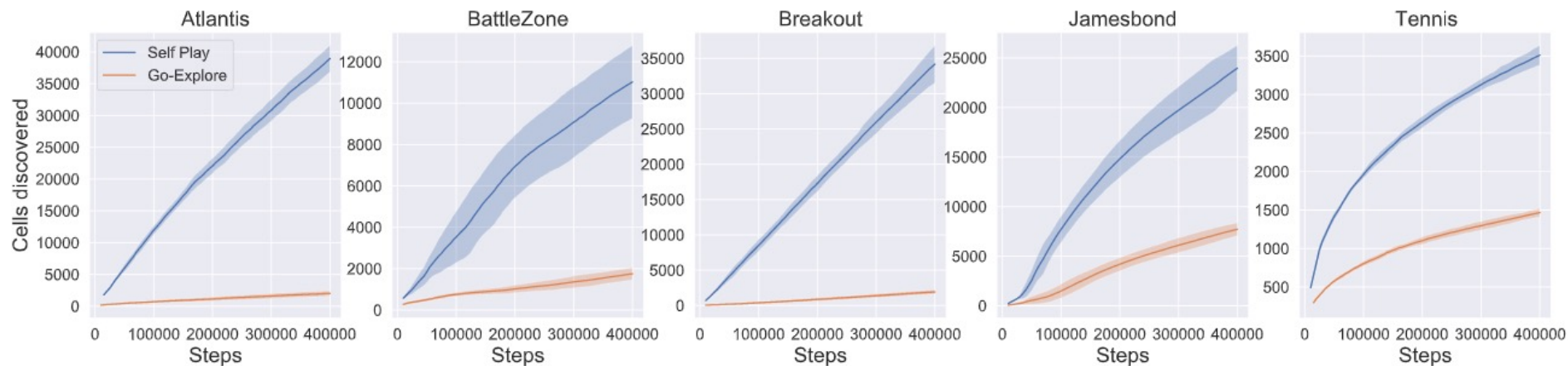
## Montezuma's Revenge Atari Game Performance



[Go-Explore: Ecoffet, Huizinga, Lehman, Stanley, Clune, 2019; also: First Return, Then Explore: Ecoffet, H, L, S, C, 2020]

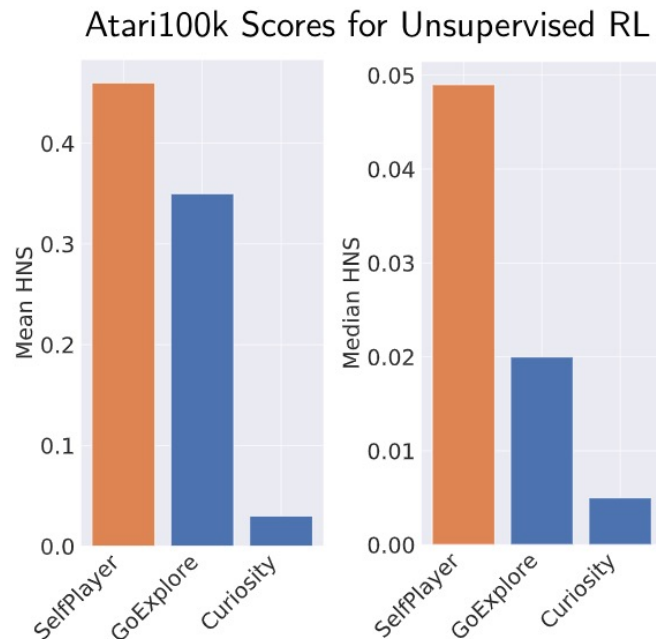
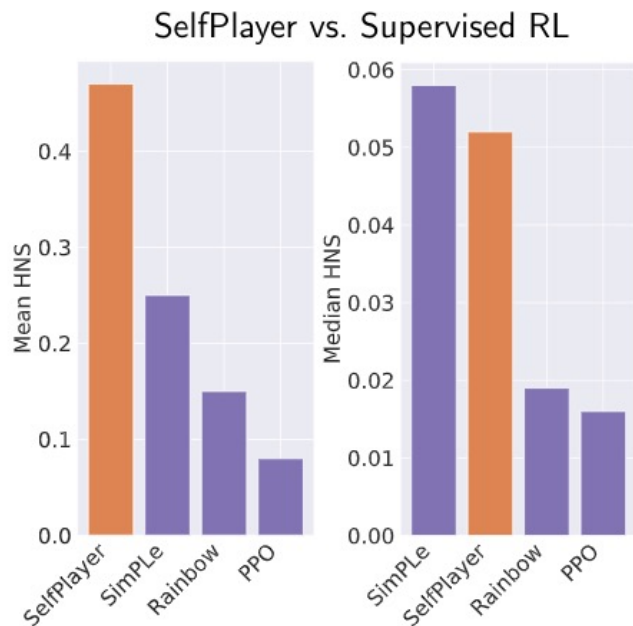
# Algorithmic RFPT: Frontier: SelfPlayer

- Brings together ideas from ASP and GoExplore



# Algorithmic RFPT: Frontier: SelfPlayer

- Brings together ideas from ASP and GoExplore





# Algorithmic RFPT: Frontier: SelfPlayer

- Brings together ideas from ASP and GoExplore

ENV	SELFPLAYER	GOEXPLORE	CURIOSITY	RANDOM	HUMAN
ALIEN	<b>551.41</b> ± 41.8	<b>550.32</b> ± 54.5	217.0	227.8	7127.7
AMIDAR	<b>270.84</b> ± 17.9	166.92 ± 11.9	18.1	5.8	1719.5
ASSAULT	<b>305.88</b> ± 17.0	180.67 ± 14.6	214.0	222.4	742.0
ASTERIX	743.43 ± 315.0	<b>1007.14</b> ± 47.8	223.0	210.0	8503.3
ASTEROIDS	718.69 ± 101.5	217.78 ± 14.0	<b>953.0</b>	719.1	47388.7
ATLANTIS	6990.91 ± 1011.0	1220.63 ± 195.3	<b>17400.0</b>	12850.0	29028.1
BANKHEIST	<b>68.99</b> ± 5.8	26.19 ± 3.9	12.9	14.2	753.1
BATTLEZONE	<b>10679.01</b> ± 604.4	6412.7 ± 636.6	2560.0	2360.0	37187.5
BEAMRIDER	294.26 ± 14.8	285.08 ± 27.1	311.0	<b>363.9</b>	16926.5
BOWLING	<b>40.53</b> ± 3.0	20.22 ± 1.2	20.2	23.1	160.7
BREAKOUT	<b>87.31</b> ± 39.2	9.9 ± 0.8	3.3	1.7	30.5
CENTIPEDE	<b>4024.76</b> ± 373.1	3308.33 ± 576.1	2100.0	2090.9	12017.0
CHOPPER	<b>850.0</b> ± 65.1	719.05 ± 51.8	597.0	<b>811.0</b>	7387.8
CRAZYCLIMBER	3574.75 ± 337.0	1326.98 ± 59.4	9690.0	<b>10780.5</b>	35829.4
DEMONATTACK	<b>455.33</b> ± 33.0	371.9 ± 26.2	161.0	152.1	1971.0
DOUBLEDUNK	-23.64 ± 0.3	-15.87 ± 0.7	-18.0	-18.6	-16.4
ENDURO	0.38 ± 0.1	<b>0.46</b> ± 0.2	0.0	0.0	860.5
FISHINGDERBY	-88.62 ± 0.5	-73.65 ± 1.9	-92.3	-91.7	-38.7
FREEWAY	4.77 ± 0.2	<b>5.16</b> ± 0.3	0.73	0.0	29.6
FROSTBITE	<b>403.09</b> ± 59.0	86.03 ± 8.9	149.0	65.2	4334.7
GOPHER	142.96 ± 14.3	53.97 ± 7.8	<b>397.0</b>	257.6	2412.5
GRAVITAR	<b>276.54</b> ± 22.0	103.97 ± 25.4	251.0	173.0	3351.4
ICEHOCKEY	-15.02 ± 0.3	-14.02 ± 0.7	-13.1	-11.2	0.9
JAMESBOND	<b>5483.89</b> ± 148.4	4663.49 ± 155.2	38.2	29.0	302.8
KANGAROO	<b>866.67</b> ± 127.5	<b>892.06</b> ± 234.8	59.3	52.0	3035.0
KRULL	83.52 ± 5.0	37.56 ± 11.2	814.0	<b>1598.0</b>	2665.5
KUNGFUMASTER	77.78 ± 33.9	17.46 ± 6.5	<b>400.0</b>	258.5	22736.3
MONTEZUMA	<b>1369.14</b> ± 140.4	26.98 ± 6.3	0.0	0.0	4753.3
MSPACMAN	<b>1381.73</b> ± 79.9	1071.11 ± 78.5	246.0	307.3	6951.6
NAMETHISGAME	1274.32 ± 64.2	1120.63 ± 62.9	2070.0	<b>2292.3</b>	8049.0
PITFALL	-446.57 ± 107.3	-430.62 ± 44.3	-39.1	-229.4	6463.7
PONG	-20.21 ± 0.2	-11.71 ± 0.7	-19.1	-20.7	14.6
PRIVATEEYE	<b>2694.78</b> ± 764.7	<b>2565.24</b> ± 893.4	-608.0	24.9	69571.3
QBERT	<b>854.01</b> ± 168.7	687.7 ± 47.7	334.0	163.9	13455.0
RIVERRAID	<b>2149.63</b> ± 203.9	1016.98 ± 54.7	1490.0	1338.5	17118.0
ROADRUNNER	<b>2271.6</b> ± 312.5	1103.17 ± 239.7	60.0	11.5	7845.0
ROBOTANK	<b>4.61</b> ± 0.3	3.62 ± 0.3	2.67	2.2	11.9
SEAQUEST	<b>252.89</b> ± 26.3	179.52 ± 45.1	237.0	68.4	42054.7
SPACEINVADERS	<b>330.06</b> ± 29.1	304.13 ± 29.2	191.0	148.0	1668.7
STARGUNNER	146.67 ± 7.6	19.05 ± 5.7	533.0	<b>664.0</b>	10250.0
TENNIS	-23.08 ± 0.1	-15.79 ± 0.5	-23.5	-23.8	-8.3
TIMEPILOT	1122.22 ± 271.7	120.63 ± 17.5	2520.0	<b>3568.0</b>	5229.2
TUTANKHAM	<b>16.2</b> ± 2.0	5.94 ± 1.1	7.41	11.4	167.6
UPNDOWN	1213.11 ± 114.5	1176.83 ± 107.3	<b>1510.0</b>	533.4	11693.2
VENTURE	23.33 ± 15.8	15.87 ± 8.9	<b>83.3</b>	0.0	1187.5
VIDEOPINBALL	10555.21 ± 1434.2	10454.87 ± 1549.6	<b>35300.0</b>	0.0	17667.9
WIZARDOFWOR	<b>7813.33</b> ± 1012.8	1206.35 ± 66.1	817.0	563.5	4756.5
ZAXXON	<b>391.11</b> ± 41.4	184.13 ± 34.8	1.68	32.5	9173.3
# ENVs BEST	<b>26</b>	6	7	7	-
AVERAGE HNS	<b>.47</b>	.35	0.03	-	-
MEDIAN HNS	<b>.05</b>	.02	0.01	-	-

# Summary

---

- Problem Motivation
- Baseline RL Algorithms Refresher
- Intrinsic Rewards for Reward-Free Pre-Training and Exploration
- Algorithmic Approaches to Exploration (can complement intrinsic reward RFPT!)
- Algorithmic Approaches to Reward-Free Pre-Training

# Many Research Opportunities

- Clearer evaluations / comparisons
- Right combination of components?
- Simpler / more robust versions → supplant existing simple SAC/TD3/PPO/DQN baselines
- Open-ended environments are the frontier; but even regular vision-based environments still challenging
- Better fine-tuning/adaptation
- Can it lead to more robust vision
- Showcase truly unexpected solutions (e.g. circuit design / ..)

# Some Additional Perspective

## Intrinsically Motivated Goal-Conditioned Reinforcement Learning: a Short Survey

**Cédric Colas**

*INRIA and Univ. de Bordeaux; Bordeaux (FR)*

CEDRIC.COLAS@INRIA.FR

**Tristan Karch**

*INRIA and Univ. de Bordeaux; Bordeaux (FR)*

TRISTAN.KARCH@INRIA.FR

**Olivier Sigaud**

*Sorbonne Université; Paris (FR)*

OLIVIER.SIGAUD@UPMC.FR

**Pierre-Yves Oudeyer**

*INRIA; Bordeaux (FR) and ENSTA Paris Tech; Paris (FR)*

PIERRE-YVES.OUDEYER@INRIA.FR

Thank you!  
pabbeel@cs.berkeley.edu