

# Model Based Reinforcement Learning For Atari (SimPLe)

Lukasz Kaiser, Mohammad Babaeizadeh, Piotr Milos, Blazej Osinski, Roy H Campbell,  
Konrad Czechowski, Dumitru Erhan, Chelsea Finn, Piotr Kozakowski, Sergey Levine, Afroz  
Mohiuddin, Ryan Sepassi, George Tucker, Henryk Michalewski

# Objectives

“ How “ this paper appeared?

“ Which” technique this paper used?

“ What “ does this paper mean?

---



# Preliminaries

What is Model-Based?



# Preliminaries

- Model Free

algorithm which does **not** use the *transition probability distribution*

ex) DQN, A3C, TRPO, PPO, DDPG, SAC ...

$$\mathcal{P}_{ss'} = \mathbb{P}[S_{t+1} = s' \mid S_t = s]$$

- Model Based

algorithm which does use the *transition probability distribution*

ex) AlphaZero, MCTS




Actually...

How about “Getting” a model with **Trial and Error**?

Model Based 험오를 멈춰주세요..

# Preliminaries



	Model-based RL	Model-free RL
Pros	Sample efficiency & Scalability	Don't need a model
Cons	Computation cost	Training Data



# Introduction





# Introduction

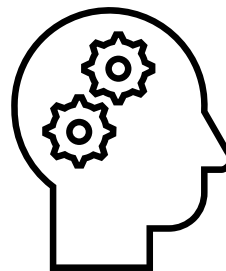
Model free Reinforcement Learning

Great performance in Atari games!

However, requires too many interactions (= Too many times) !!!



Environment



Agent

Then… How about Human?

What's Next?

# Introduction

Human possess an **intuitive understanding** of the physical processes

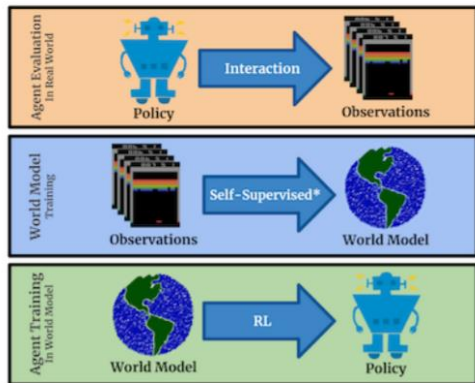
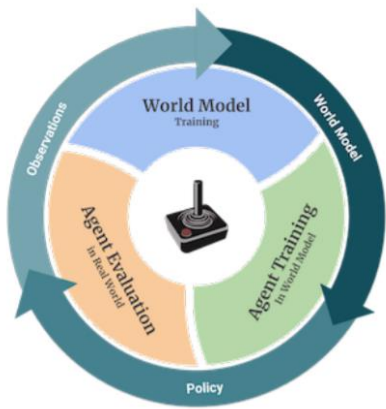
We know planes can fly, ball can roll, and bullet can destroy aliens!

Let's make our agent learns intuitive understanding!

How?

**Video Prediction!**

# Introduction



\*World Model:

우리가 생각하는 가상의 simulation 환경!  
실제 env와 interaction 없이도 가능한 Network

1. Interacting with the real environment following the latest policy
2. Collected observations will be used to train the current world model
3. Update the policy by acting inside the world model

# Introduction

SIMulated Policy Learning

SimPLe

Utilize video prediction techniques and trains a policy to play the game within the learned model

여기서 Simulation은 “Neural Network” 라는 것이 포인트!



# Related Work



# Related Work

\_Model free

DQN

A3C

Double DQN

REINFORCE

Dueling DQN

A2C

Prioritized DQN

TRPO

Rainbow DQN

PPO

DDPG

Good Performance in Atari,  
but remains far higher than the **amount of experience required** for human players to learn each game

# Related Work\_Video prediction

(2015) Action conditional video prediction using deep networks in atari games (Oh et al)

(2017) Recurrent environment simulators (Chiappa et al)

(2016) A deep learning approach for joint video frame and reward prediction in Atari games (Leibfried et al)

This paper focus on using video prediction in the context of learning how to play the game well and positively verify that learned simulators can be used to train a policy useful in original environments



# Related Work

Other works

(2017) Value Prediction Network (Oh et al)

→ use a model of reward to augment model-free learning, but not aim to model or predict future

(2019) Recurrent world models facilitate policy evaluation (Sodhani et al)

→ present a way to compose VAE with RNN, but need enough exploration

MCTS, Dyna-DQN, Generative Adversarial Tree Search (GATS), ...

No prior work has successfully demonstrated Model-based control via predictive models that achieve competitive results with model-free RL!



# Model

Deterministic Model / Stochastic Model



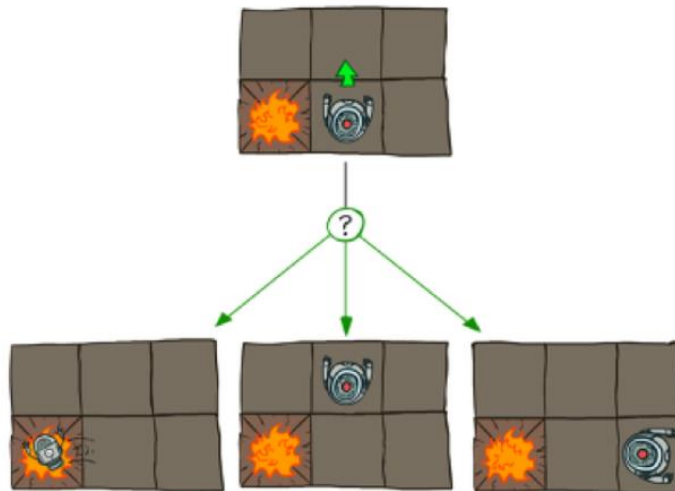
# Model

## Stochastic Model & Deterministic Model

Deterministic Grid World

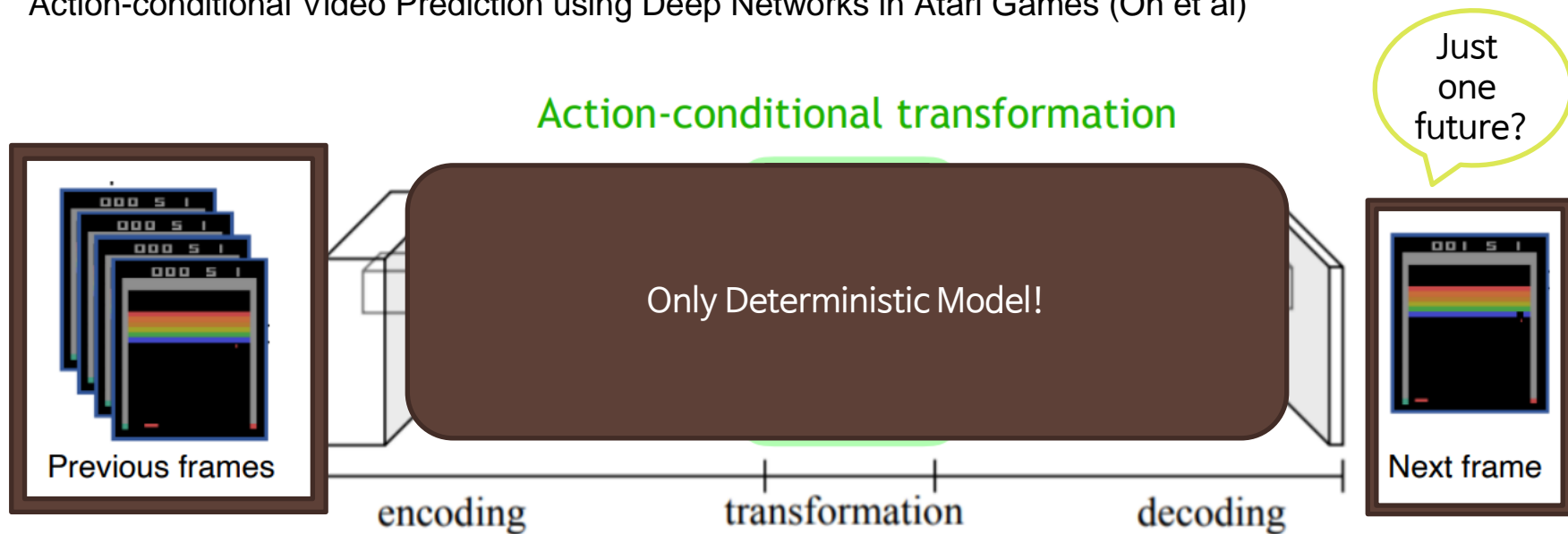


Stochastic Grid World

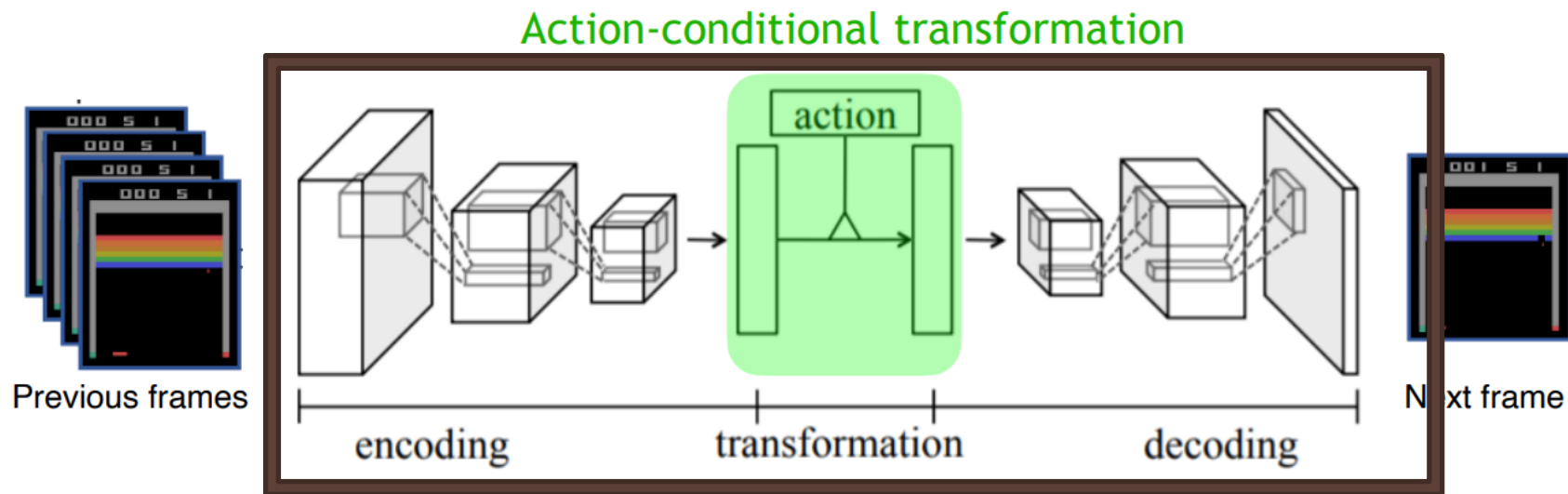


# Model

Action-conditional Video Prediction using Deep Networks in Atari Games (Oh et al)



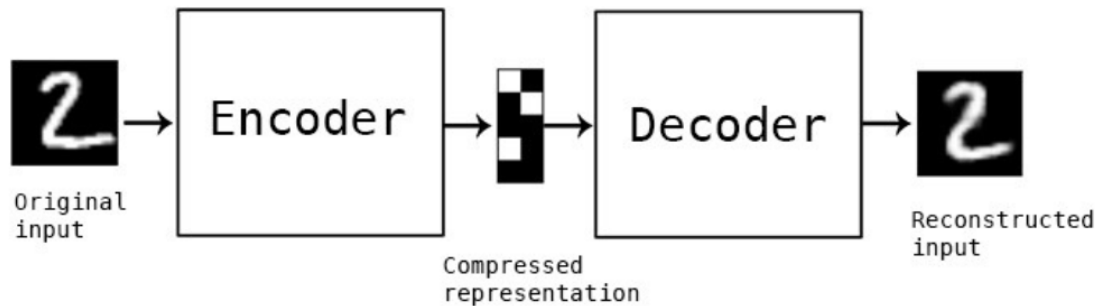
# Model



Wait.. Have you seen this structure before?

# Model

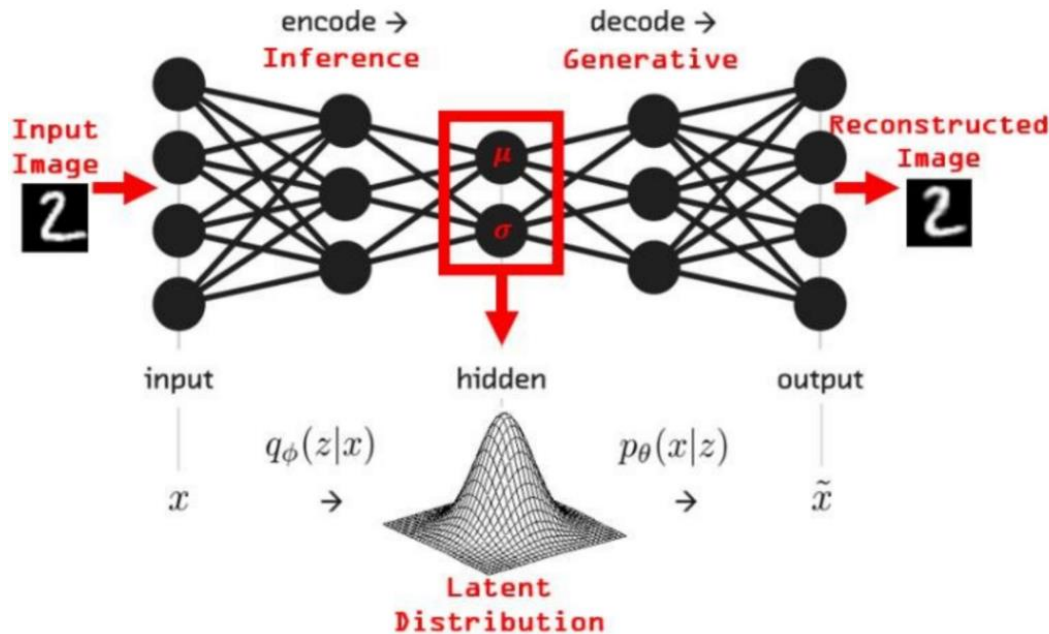
## Autoencoder



Only one output image(frame)...

# Model

## Variational Autoencoder



Distribution 형태의 latent variable  $z$ 를 통해  $x$ 를 generate 하고 싶다!

$$x = g_\theta(z)$$

$$p(x) = p(x|z)$$

$$p(x) = E_{z \sim p_\theta(z)}[p(x|z)]$$

$$p(x) = E_{z \sim p_\theta(z|x)}[p(x|z)] \\ \approx E_{z \sim q_\theta(z|x)}[p(x|z)]$$

자세한 증명은 생략

# Model

---

## Algorithm 1: Pseudocode for SimPLe

---

Initialize policy  $\pi$

Initialize model parameters  $\theta$  of  $env'$

Initialize empty set  $\mathbf{D}$

**while** not done **do**

▷ collect observations from real env.

$\mathbf{D} \leftarrow \mathbf{D} \cup \text{COLLECT}(env, \pi)$

▷ update model using collected data.

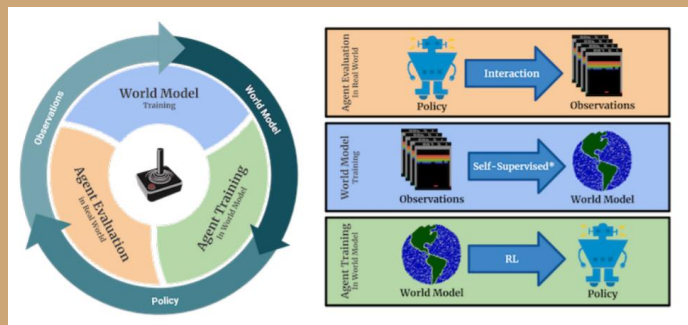
$\theta \leftarrow \text{TRAIN\_SUPERVISED}(env', \mathbf{D})$

▷ update policy using world model.

$\pi \leftarrow \text{TRAIN\_RL}(\pi, env')$

**end while**

---



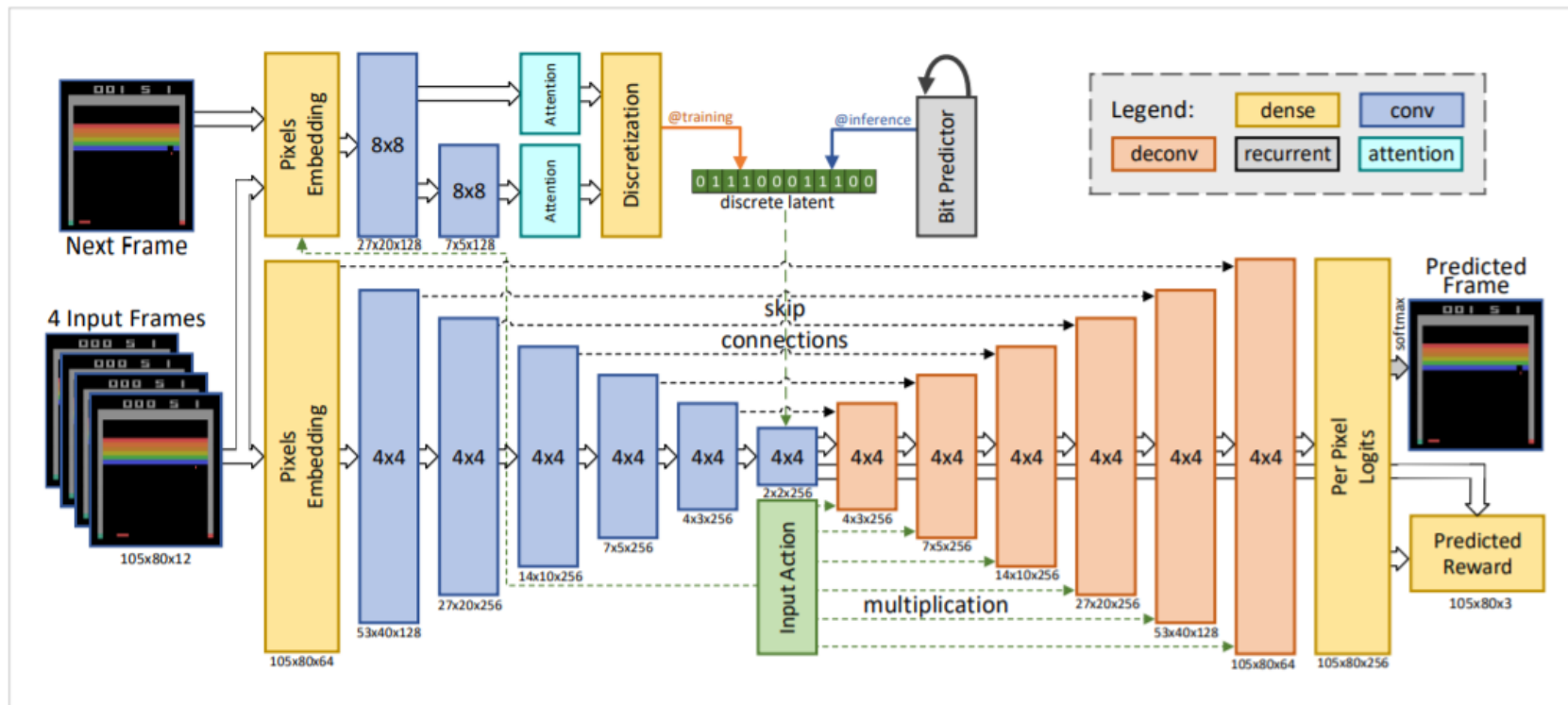
Neural Network simulated environment,  $env'$

Our purpose is to train a policy  $\pi$  using  $env'$  so that  $\pi$  achieves good performance in the original environment  $env$

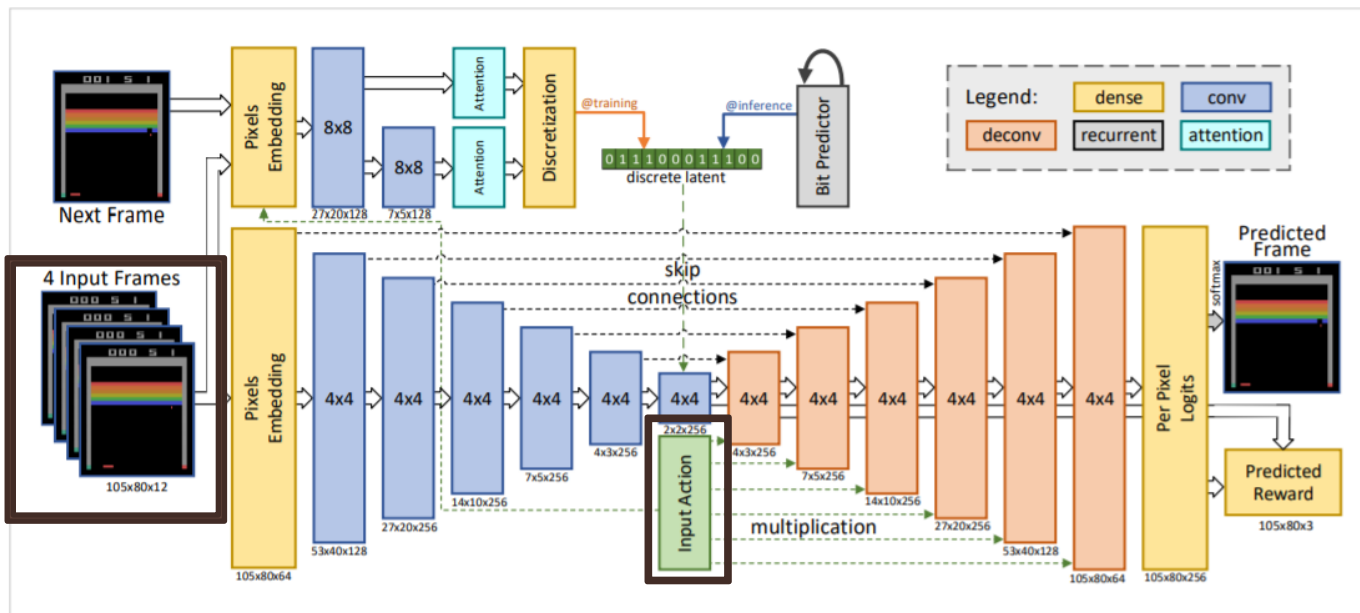
We aim to use as few interactions with  $env$  as possible



# Model

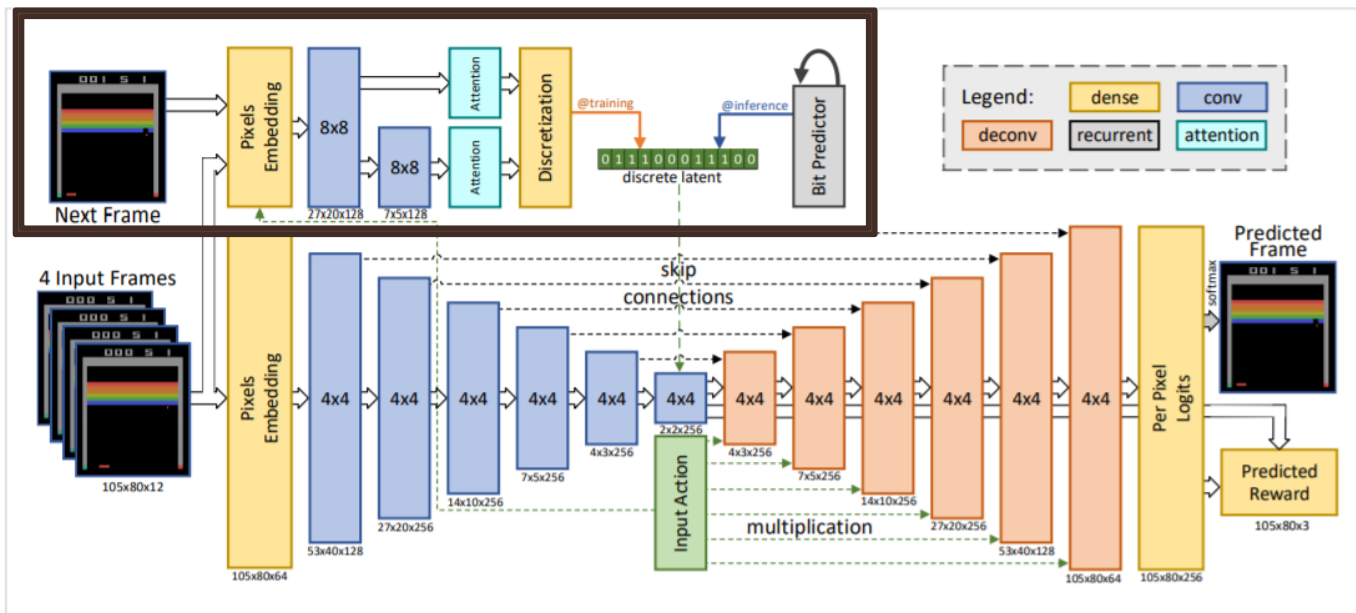


# Model



Input : four consecutive game frames and action  $a$

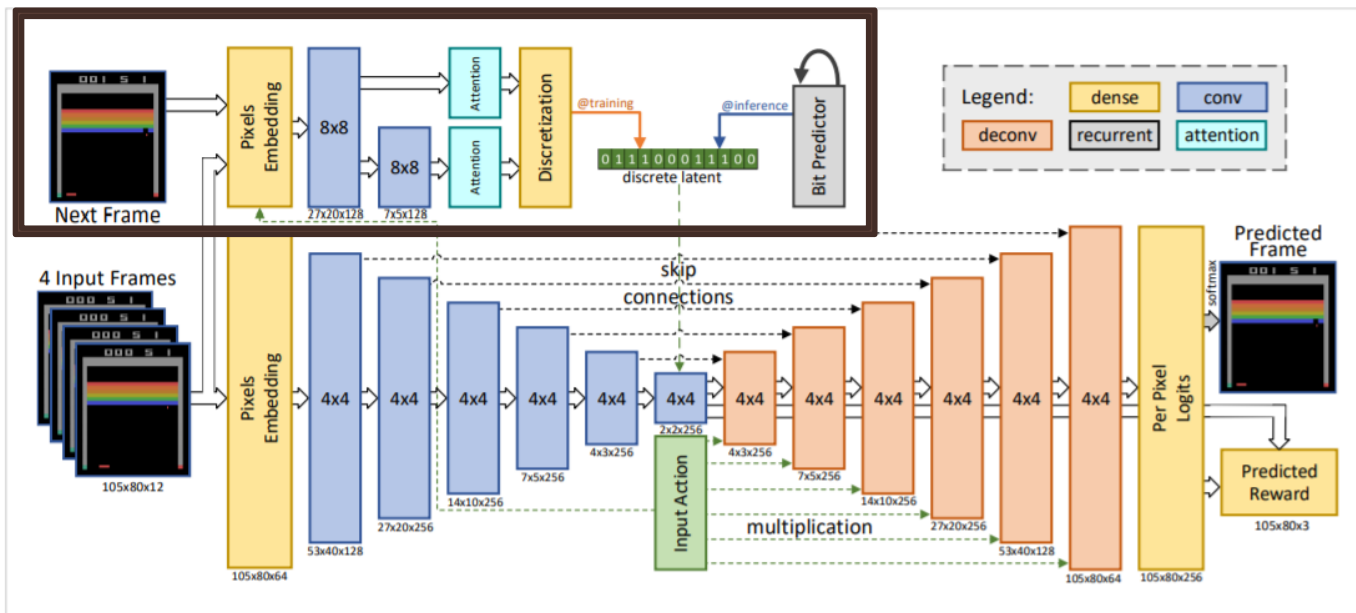
# Model



Input : receives the future target frame as input and approximates the distribution of the posterior  $p(z|x)$

At test time, the latent values are sampled from an assumed prior

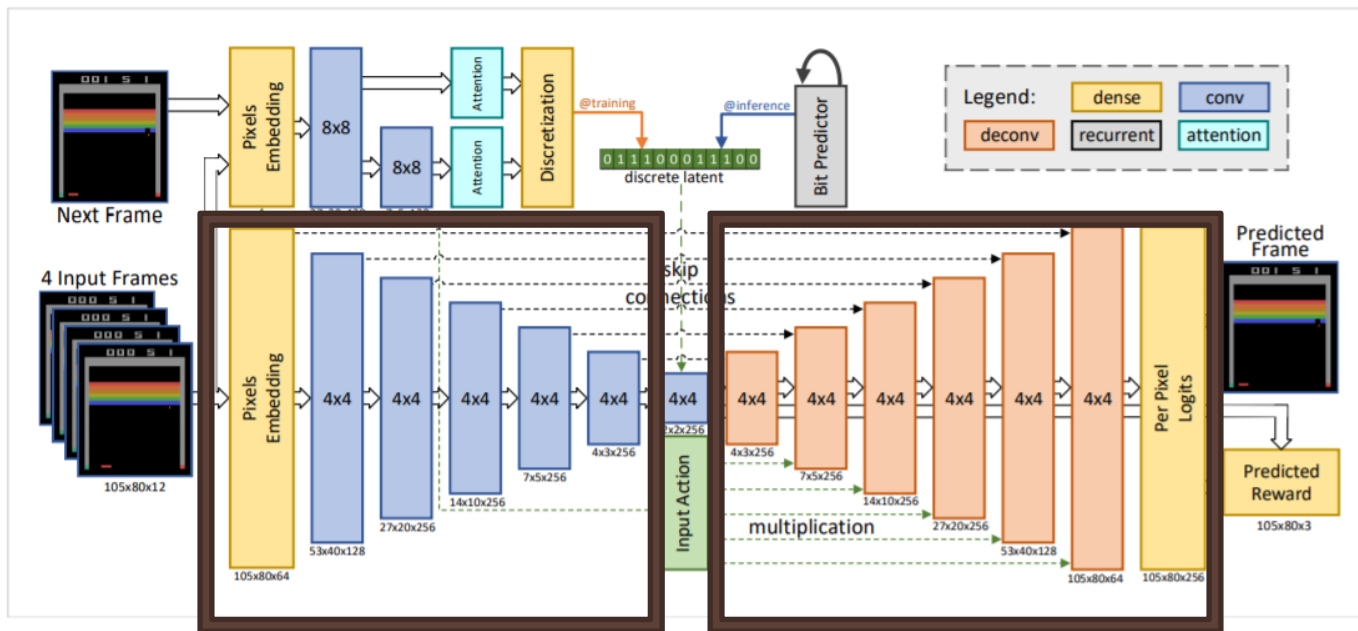
# Model



Why discrete latent?

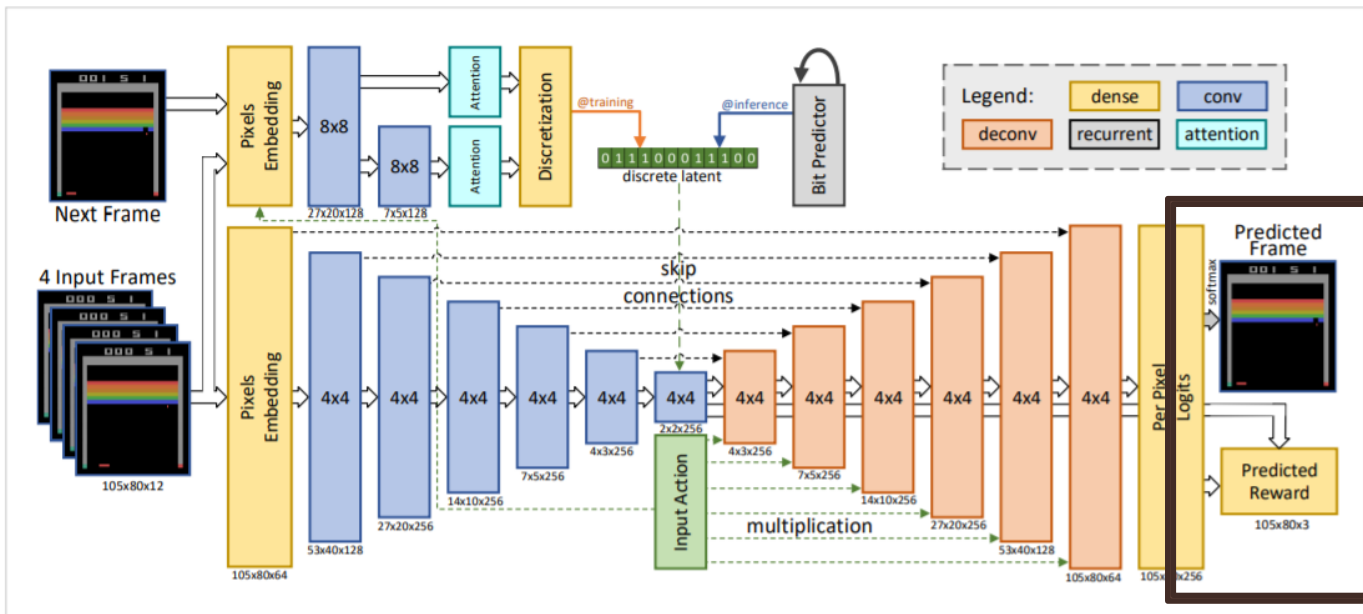
1. The weight of the KL divergence is game dependent
2. The weight is usually a very small number,  $[10^{-3}, 10^{-5}]$

# Model



Action is vector which is multiplied channel-wise with the output of the convolution layer + **Skip Connections!**

# Model



# Model

## Loss function

Use the *clipped loss*,  $\max(Loss, C)$

It may make model to concentrate on small but important areas (e.g. the ball in Pong)

$C = 10$  for L2 loss,  $0.03$  for softmax loss

## Policy Training

The algorithm generates rollouts in the *env*' and uses them to improve policy  $\pi$

Using PPO algorithm for learning policy

Imperfections of the model compounding over time

→ Short rollout ( $N=50$ )

→ degrading effect of PPO algorithm

→ add to the reward

---



# Experiment





# Experiment Setting

Main loop : 15 times

World model : 45K steps in first iteration, and 15K steps in others

PPO

- 16 parallel agents collecting different steps from the  $env'$

Training Data

- Data from real environment:

$$\begin{aligned} & 6400 \text{ interactions} + 6400 \text{ interactions} \times 15 \text{ iterations} \\ & = 102,400 \text{ interactions} = 409,600 \text{ frames} \end{aligned}$$

- Data from simulated environment:

$$1,000,000 \text{ interactions} \times 15 \text{ iterations} = 15,000,000 \text{ interactions}$$

26 games, comparisons are Rainbow, PPO

---

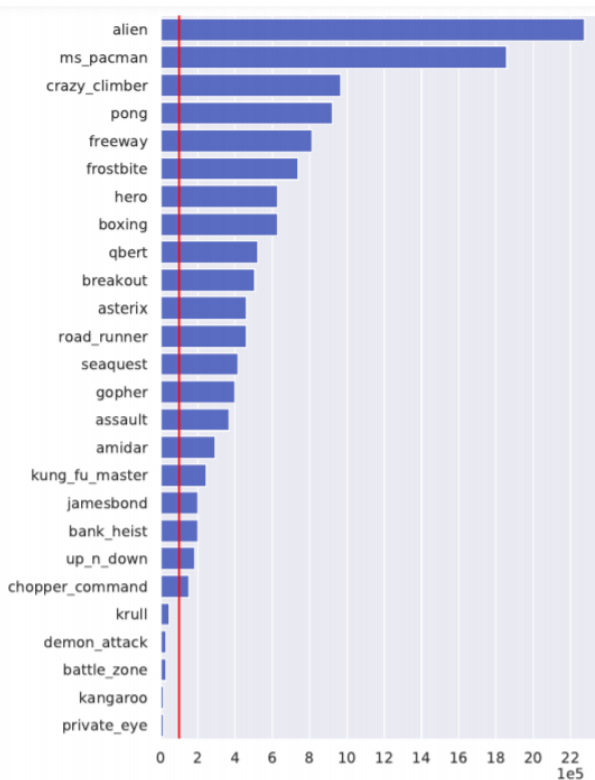
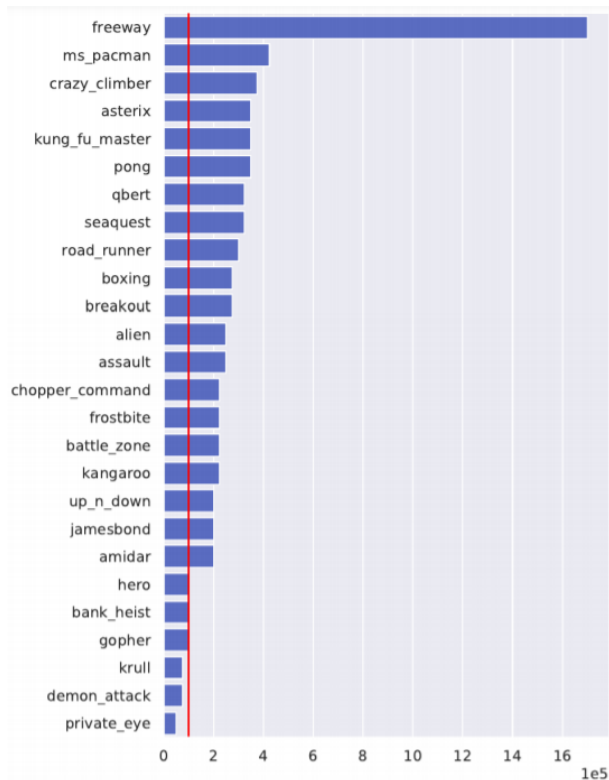
**Algorithm 1:** Pseudocode for SimPLe

---

```
Initialize policy  $\pi$ 
Initialize model parameters  $\theta$  of  $env'$ 
Initialize empty set  $\mathbf{D}$ 
while not done do
  ▷ collect observations from real env.
   $\mathbf{D} \leftarrow \mathbf{D} \cup \text{COLLECT}(env, \pi)$ 
  ▷ update model using collected data.
   $\theta \leftarrow \text{TRAIN\_SUPERVISED}(env', \mathbf{D})$ 
  ▷ update policy using world model.
   $\pi \leftarrow \text{TRAIN\_RL}(\pi, env')$ 
end while
```

---

# Experiment

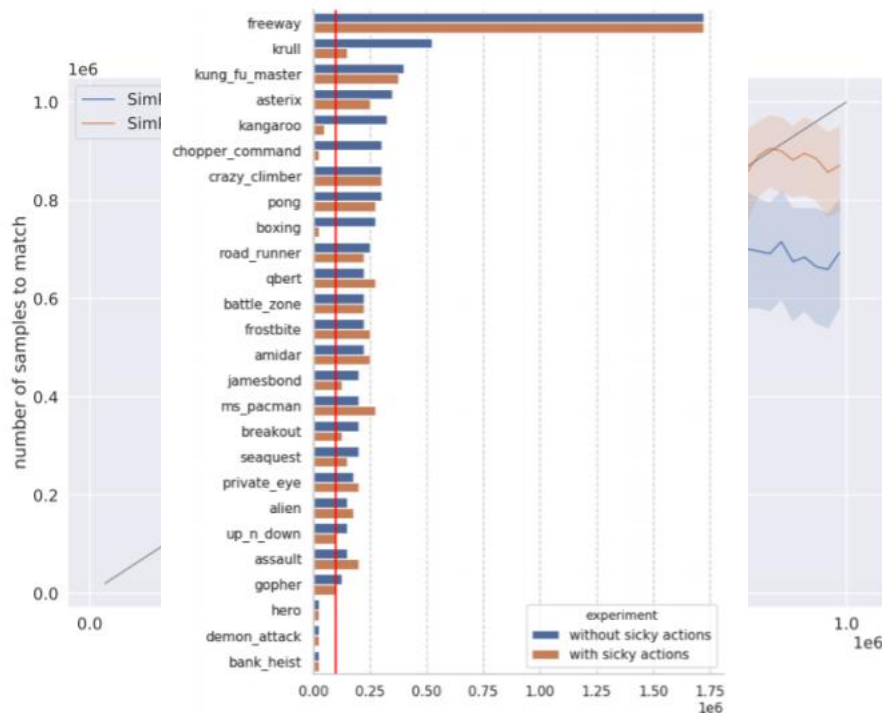


Model-based RL provides an effective approach to learning tasks!

For the 6 games, it exceeds the avg human score (future work)

In some cases, we observed high variance when mismatch the model and the real env

# Experiment



SimPLe excels in a low data regime, but its advantage disappears with a bigger amount of data

Also, It's useful initialization for model-free PPO training

SimPLe is also effective in stochastic environments.



# Conclusion



# Conclusion

This paper present SimPLe, a model-based RL approach that operates directly on raw pixel observations and learns effective policies to play games in the Atari Learning Environment

It can be applied in highly stochastic environments

The representation learned by model is likely be more meaningful than the raw pixel observations

## Limitations

1. Final scores are lower than SOTA model-free methods
2. The performance of our method generally varied between different runs on the same game  
→ capture uncertainty via Bayesian parameter posteriors or ensembles
3. Computational and time requirement of training are substantial  
→ developing lighter models & applied to other environments