
Stochastic Latent Actor-Critic: Deep Reinforcement Learning with a Latent Variable Model

Soft Actor Critic

Alex X. Lee^{1,2}

Anusha Nagabandi¹

Pieter Abbeel¹

Sergey Levine¹

¹University of California, Berkeley

²DeepMind

{alexlee_gk,nagaban2,pabbeel,svlevine}@cs.berkeley.edu

Reviewed by Susang Kim

Contents

1. Introduction & Motivation
 2. Preliminaries & Related Works
 3. Methods & Experiments
 4. Conclusion & Limitation
- Reference

1.Introduction - Motivation

Challenges in **high-dimensional observation**. (e.g., image, unstructured, a large number of pixels)

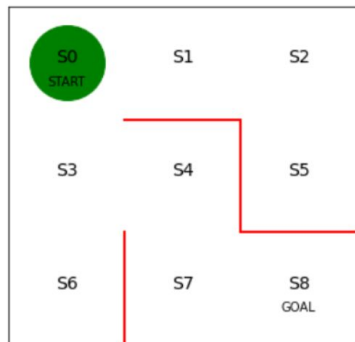
The standard approach relies **on sensors to obtain information that would be helpful for learning**.

It's hard to use reinforcement learning algorithms to solve tasks using **only low-level observations**, such as learning **robotic control using only unstructured raw image data**. (Robot Vision)

⇒ **tackle these two problems separately**, by explicitly learning latent representations

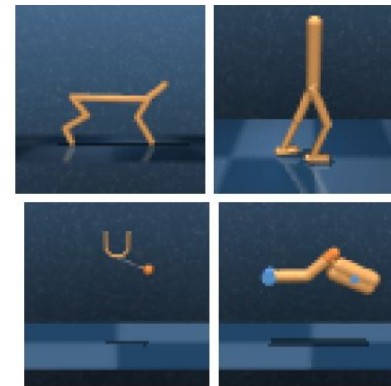
Learning from only image data is hard because the RL algorithm must learn **both a useful representation of the data and the task itself**. (+representation learning problem)

The agent (cheetah) didn't have any prior knowledge about movement.



```
state = np.array([[np.nan, 1, 1, np.nan], # s0
                  [np.nan, 1, np.nan, 1], # s1
                  [np.nan, np.nan, 1, 1], # s2
                  [1, 1, 1, np.nan], # s3
                  [np.nan, np.nan, 1, 1], # s4
                  [1, np.nan, np.nan, np.nan], # s5
                  [1, np.nan, np.nan, np.nan], # s6
                  [1, 1, np.nan, np.nan], # s7, s8
                  ])
```

Sensor numerical data



raw image data

1.Introduction - Stochastic Latent Actor-Critic

Standard model-free deep RL aims to unify these challenges of **representation learning and task learning** into a end-to-end training procedure. (policy or value function)

⇒ inefficient, practice be slow and sensitive to hyperparameters

SLAC: a sample-efficient and high-performing RL algorithm for learning policies for **complex continuous control tasks** directly from **high-dimensional image inputs**.

SLAC substantially outperforms both prior **model-free and model-based RL(agent in environment) algorithms on a range of image-based** continuous control benchmark tasks.

[SLAC]

- 1) High-dimensional observations as a **latent process, with a Gaussian prior and latent dynamics**.
- 2) **Partially Observed Markov Decision Process (POMDP)** : where the stochastic latent state enables the model to represent **uncertainty** about any of the state variables, given the past observations
- 3) Markovian **critic on latent state samples** and trains an **actor on a history of observations and actions**, resulting in our stochastic latent actor-critic (SLAC) algorithm.

The main contribution is a novel and principled approach that **integrates learning stochastic sequential models** and **RL into a single method** in the model's learned latent space.

1. Introduction - Contribution

Representation learning in RL : “Representation learning bottleneck”: a considerable portion of the learning period must be spent **acquiring good representations of the observation space**

⇒ Jointly modeling between consecutive latent states, we make it feasible for our proposed algorithm to perform bellman backups directly in the latent space of the learned model

Partial observability in RL : Recent work has proposed end-to-end RL methods that use recurrent neural networks to process histories of observations and actions, but without constructing a model of the POMDP. Other works, **learn latent-space dynamical system models and then use them to solve the POMDP** with model-based RL

⇒ Our approach does not use the model for prediction, and performs infinite horizon policy optimization. benefits from the good asymptotic performance of model-free RL, while at the same time leveraging the improved latent space representation for sample efficiency. our method learns a critic directly on latent state samples, which enables scaling to more complex tasks.

Sequential latent variable models : Various modeling choices to learn **stochastic sequential models. factorization of the generative and inference models**, their network architectures, and the objectives used in their training procedures.

⇒ Our approach is compatible with any of these sequential latent variable models, with the only requirement being that they provide a mechanism to sample latent states from the belief of the learned Markovian latent space.

2.Preliminaries - VAE: Variational Autoencoder

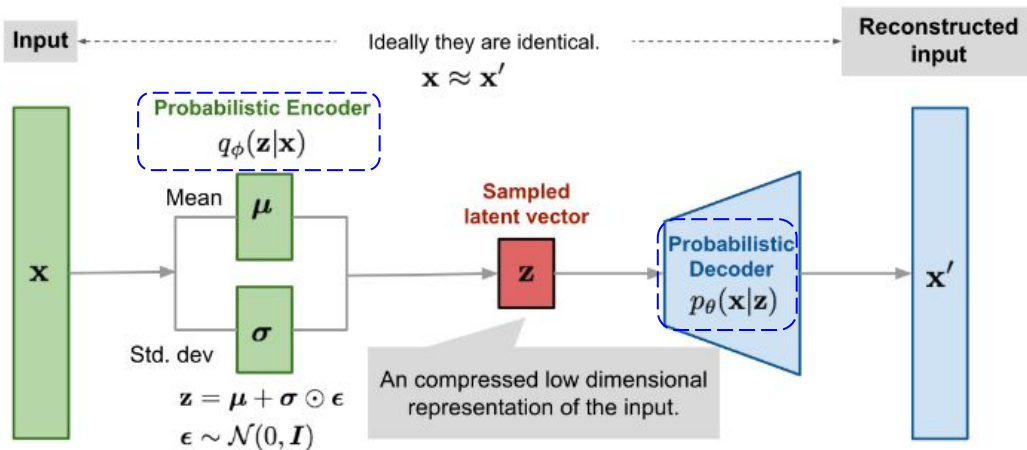
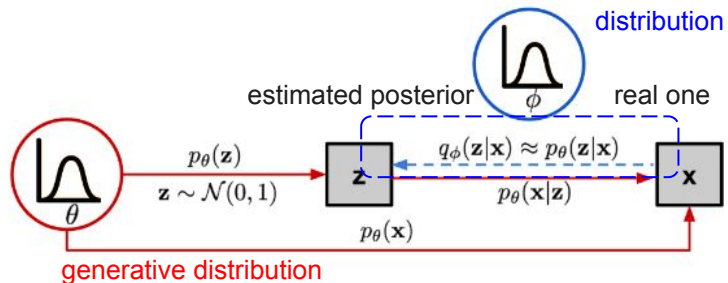


Fig. 9. Illustration of variational autoencoder model with the multivariate Gaussian assumption.



$$D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x}))$$

$$= \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})} d\mathbf{z}$$

$$= \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})p_\theta(\mathbf{x})}{p_\theta(\mathbf{z}, \mathbf{x})} d\mathbf{z} \quad ; \text{ Because } p(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}, \mathbf{x})/p(\mathbf{x})$$

$$= \int q_\phi(\mathbf{z}|\mathbf{x}) (\log p_\theta(\mathbf{x}) + \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}, \mathbf{x})}) d\mathbf{z} \quad ; \text{ Because } \int q(\mathbf{z}|\mathbf{x}) d\mathbf{z} = 1$$

$$= \log p_\theta(\mathbf{x}) + \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}, \mathbf{x})} d\mathbf{z}$$

$$= \log p_\theta(\mathbf{x}) + \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{x}|\mathbf{z})p_\theta(\mathbf{z})} d\mathbf{z} \quad ; \text{ Because } p(\mathbf{z}, \mathbf{x}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$$

$$= \log p_\theta(\mathbf{x}) + \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} [\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z})} - \log p_\theta(\mathbf{x}|\mathbf{z})]$$

$$= \log p_\theta(\mathbf{x}) + D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})) - \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z})$$

$$\log p_\theta(\mathbf{x}) - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) = \mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z}) - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))$$

$$L_{\text{VAE}}(\theta, \phi) = -\log p_\theta(\mathbf{x}) + D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x}))$$

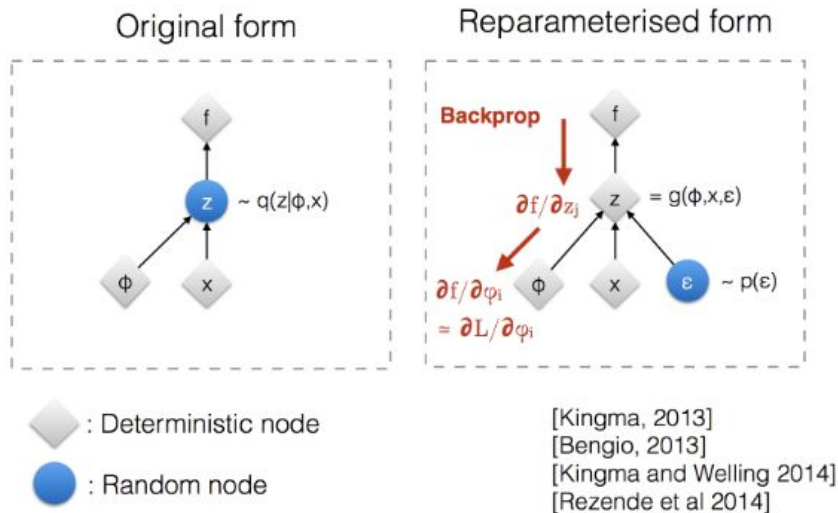
$$= -\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})} \log p_\theta(\mathbf{x}|\mathbf{z}) + D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))$$

$$\theta^*, \phi^* = \arg \min_{\theta, \phi} L_{\text{VAE}}$$

non-negative

$$-L_{\text{VAE}} = \log p_\theta(\mathbf{x}) - D_{\text{KL}}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) \leq \log p_\theta(\mathbf{x})$$

2.Preliminaries - Reparameterization Trick



The expectation term in the loss function invokes generating samples from $\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})$. Sampling is a stochastic process and therefore we cannot backpropagate the gradient. To make it trainable. (input -> model -> deterministic <-> stochastic)

$$\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{2(i)}\mathbf{I})$$

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \text{ where } \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$$

(SAC - backpropagation : noise vector)

$$\mathbf{u}_i^j = \mathbf{f}_{\theta}(\mathbf{x}_i, \boldsymbol{\eta}_j)$$

$$\mathbf{u}_i^j \sim \mathcal{N}(\boldsymbol{\mu}_{\theta}(\mathbf{x}_i), \boldsymbol{\sigma}_{\theta}^2(\mathbf{x}_i))$$

$$= \boldsymbol{\mu}_{\theta}(\mathbf{x}_i) + \boldsymbol{\sigma}_{\theta}(\mathbf{x}_i) \boldsymbol{\eta}_j, \boldsymbol{\eta}_j \sim \mathcal{N}(0, \mathbf{I})$$

$$L_{\pi}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}_i \sim \mathcal{D}} [\mathbb{E}_{\boldsymbol{\eta}_j \sim \mathcal{N}} [\alpha \log \pi_{\theta}(\mathbf{u}_i | \mathbf{x}_i) - Q_{\phi}(\mathbf{x}_i, \mathbf{u}_i)] | \mathbf{x}_i]$$

$$\nabla_{\boldsymbol{\theta}} L_{\pi}(\boldsymbol{\theta}) = \sum_i [\alpha \nabla_{\boldsymbol{\theta}} \log \pi_{\theta}(\mathbf{u}_i | \mathbf{x}_i) + \nabla_{\boldsymbol{\theta}} \mathbf{f}_{\theta}(\mathbf{x}_i, \boldsymbol{\eta}_i) (\alpha \nabla_{\mathbf{u}_i} \log \pi_{\theta}(\mathbf{u}_i | \mathbf{x}_i) - \nabla_{\mathbf{u}_i} Q_{\phi}(\mathbf{x}_i, \mathbf{u}_i))]$$

2.Preliminaries - ELBO(Evidence Lower Bound)

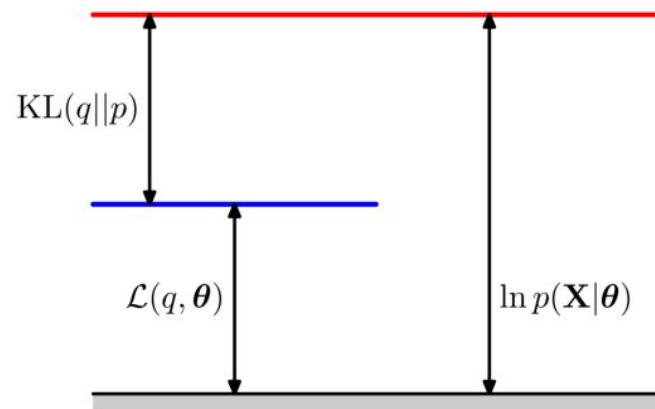
Variational Lower Bound

$$\log p(\mathbf{x}) \geq \underbrace{E_{\mathbf{z} \sim q} [\log p(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Error Decoder Cross Entropy Observation model}} - \underbrace{D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}))}_{\text{Regularization parameter Encoder KL divergence Variational approximate}}$$

Variational Inference
gaussian distribution
posterior

Reconstruction Error
Decoder
Cross Entropy
Observation model

Regularization parameter
Encoder
KL divergence
Variational approximate



Bishop – Pattern Recognition and Machine Learning

SLAC: RL objective for learning model, Actor and Critic.

$$\log p(\mathbf{x}_{1:\tau+1} | \mathbf{a}_{1:\tau}) \geq \mathbb{E}_{\mathbf{z}_{1:\tau+1} \sim q} \left[\sum_{t=0}^{\tau} \log p(\mathbf{x}_{t+1} | \mathbf{z}_{t+1}) - D_{\text{KL}}(q(\mathbf{z}_{t+1} | \mathbf{x}_{t+1}, \mathbf{z}_t, \mathbf{a}_t) \parallel p(\mathbf{z}_{t+1} | \mathbf{z}_t, \mathbf{a}_t)) \right]$$

$$q(\mathbf{z}_1 | \mathbf{x}_1, \mathbf{z}_0, \mathbf{a}_0) := q(\mathbf{z}_1 | \mathbf{x}_1)$$

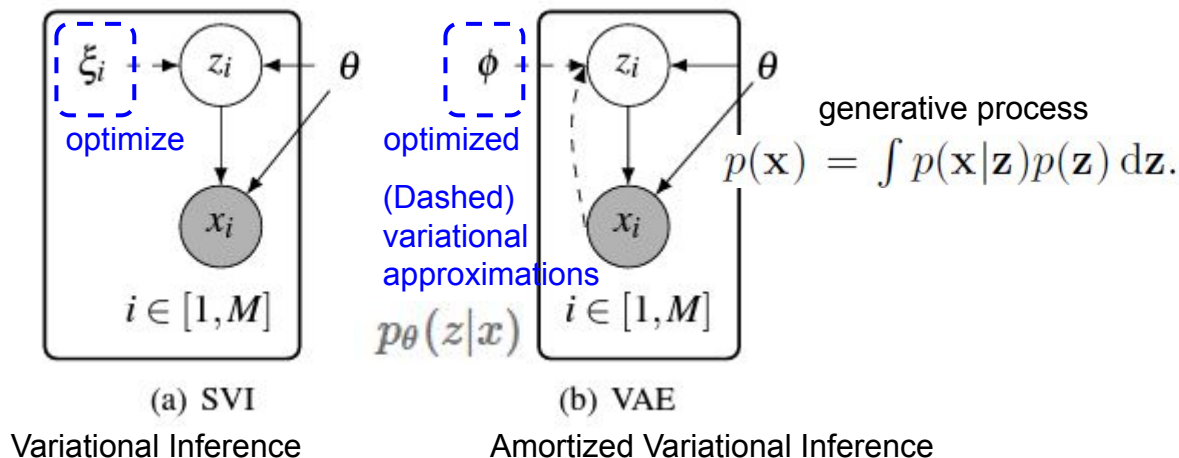
encoder

$$p(\mathbf{z}_1 | \mathbf{z}_0, \mathbf{a}_0) := p(\mathbf{z}_1)$$

decoder

2.Preliminaries - Amortized Variational Inference

Amortized variational inference don't need latent fitting.



Amortized inference usually refers to inference over local variables. Instead of approximating separate variables for each data point, as shown in Figure (a) amortized VI assumes that the local variational parameters can be predicted by a parameterized function of the data. Thus, once this function is estimated, the latent variables can be acquired by passing new data points through the function, as shown in Figure (b). Deep neural networks used in this context are also called inference networks. Amortized VI with inference networks thus combines probabilistic modeling with the representational power of deep learning.

2.Preliminaries - Partially Observed Markov Decision Process (POMDP) MLCF

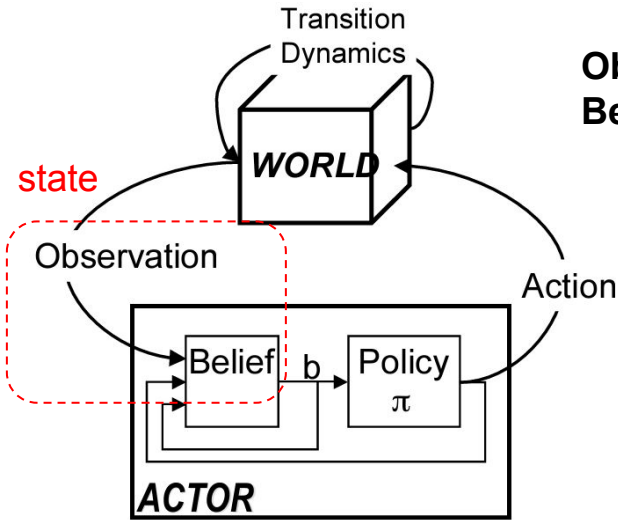
POMDP : Partially Observable Markov Decision Process.

Uncertainty about the action outcome and the world state due to partial information

MDP: Tractable to solve, Relatively easy to specify, Assumes perfect knowledge of state

POMDP : Treats all sources of uncertainty uniformly, Allows for information gathering actions, intractable to solve optimally (**Add 2 Parameter : Observations, Belief State**)

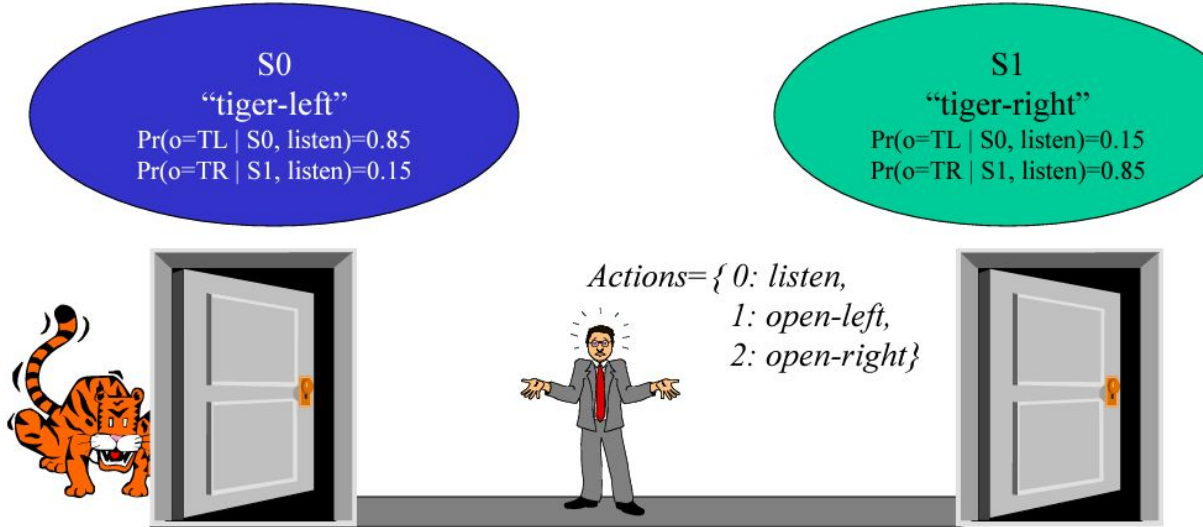
$$\langle S, A, P, R, \gamma \rangle \longrightarrow \langle S, A, P, R, \Omega, O, \gamma \rangle$$



Observation- indirect measurements (sensations of state and reward)
Beliefs - Understanding of state with **uncertainty**

Markov Models		Do we have control over the state transitions?	
		NO	YES
Are the states completely observable?	YES	Markov Chain	MDP Markov Decision Process
	NO	HMM Hidden Markov Model	POMDP Partially Observable Markov Decision Process

2.Preliminaries - POMDP : Continuous Space Belief MDP



Reward Function

- Penalty for wrong opening: -100
- Reward for correct opening: +10
- Cost for listening action: -1

Observations

- to hear the tiger on the left (TL)
- to hear the tiger on the right (TR)

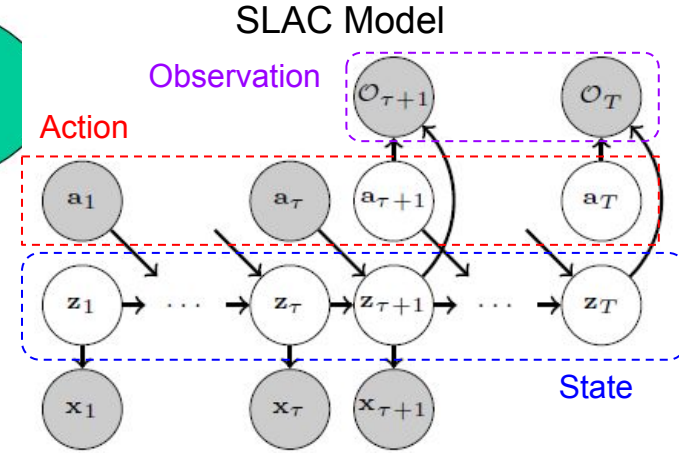
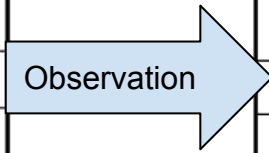


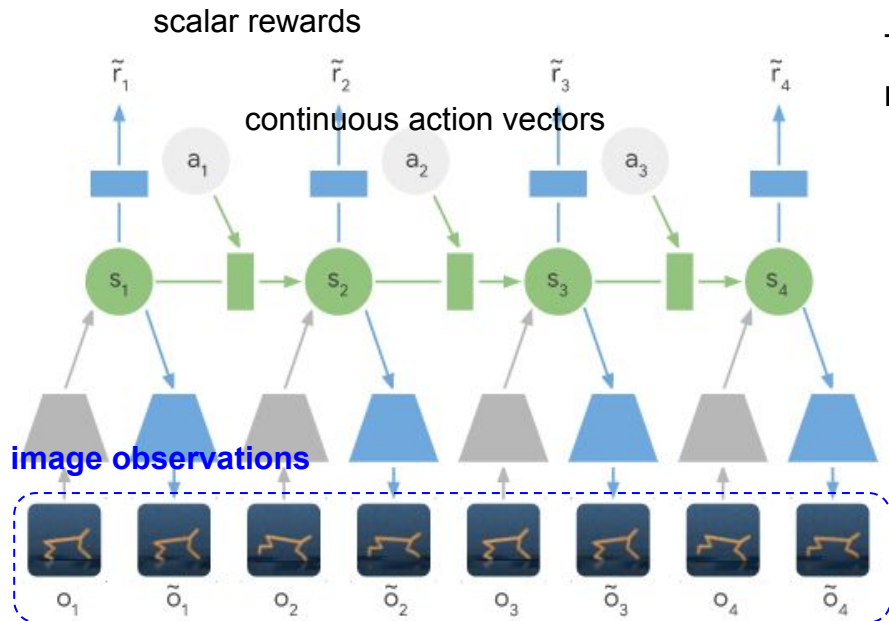
Figure 1: Graphical model of POMDP with optimality variables for $t \geq \tau + 1$.

Prob. (LEFT)	Tiger: left	Tiger: right
Tiger: left	0.5	0.5
Tiger: right	0.5	0.5

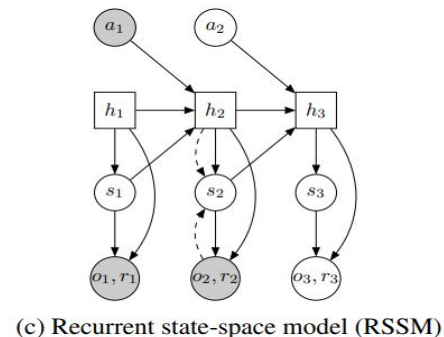
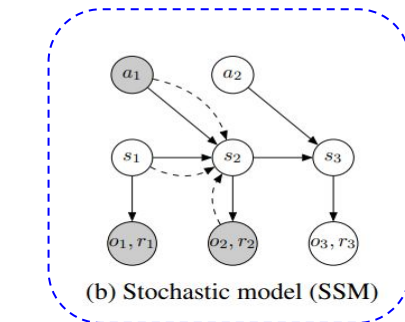
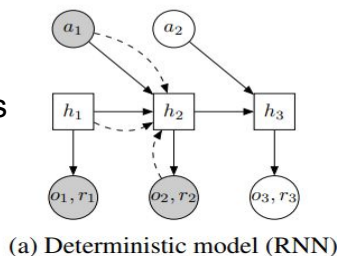


Prob. (LISTEN)	O: TL	O: TR
Tiger: left	0.85	0.15
Tiger: right	0.15	0.85

2.Preliminaries - Latent dynamics



Circles : stochastic variables
Squares : deterministic variables
Solid lines : generative process
Dashed lines : inference model.



To solve unknown environments via planning, we need to **model the environment dynamics from experience**

Sequences $\{o_t, a_t, r_t\}_{t=1}^T$

Transition function:

$$s_t \sim p(s_t | s_{t-1}, a_{t-1})$$

Observation function:

$$o_t \sim p(o_t | s_t) \quad (1)$$

Reward function:

$$r_t \sim p(r_t | s_t)$$

Policy:

$$a_t \sim p(a_t | o_{\leq t}, a_{< t}),$$

3.Methods - Maximum Entropy RL in Fully Observable MDPs

The probability distribution which best represents the current state of knowledge
 ⇒ minimize Worst-loss policy

Objective Function(maximize the expected entropy)

The entropy term **encourages exploration**

$$\mathcal{H}(\pi_\phi(\cdot|\mathbf{s}_t)). \quad \text{Entropy Bonus}$$

Maximum entropy objective function

$$\sum_{t=1}^T \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \rho_\pi} [r(\mathbf{s}_t, \mathbf{a}_t)] + \alpha \mathcal{H}(\pi_\phi(\cdot|\mathbf{s}_t))$$

RL

temperature parameter that trades off between
 maximizing for the reward and for the policy entropy

Minimize the soft Bellman residual

discount factor

$$J_Q(\theta) = \frac{1}{2} \left(Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \left(r_t + \gamma \mathbb{E}_{\mathbf{a}_{t+1} \sim \pi_\phi} [Q_\theta(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}) - \alpha \log \pi_\phi(\mathbf{a}_{t+1}|\mathbf{s}_{t+1})] \right) \right)^2$$

delayed parameters

The policy parameters ϕ are optimized to update
 the policy towards the exponential of the soft
 Q-function, resulting in the policy loss

$$J_\pi(\phi) = \mathbb{E}_{\mathbf{a}_t \sim \pi_\phi} [\alpha \log(\pi_\phi(\mathbf{a}_t|\mathbf{s}_t)) - Q_\theta(\mathbf{s}_t, \mathbf{a}_t)]$$

SLAC builds on this maximum entropy RL, by integrating **representation learning and partial observability**

3.Methods - POMDP & Factorize variational distribution

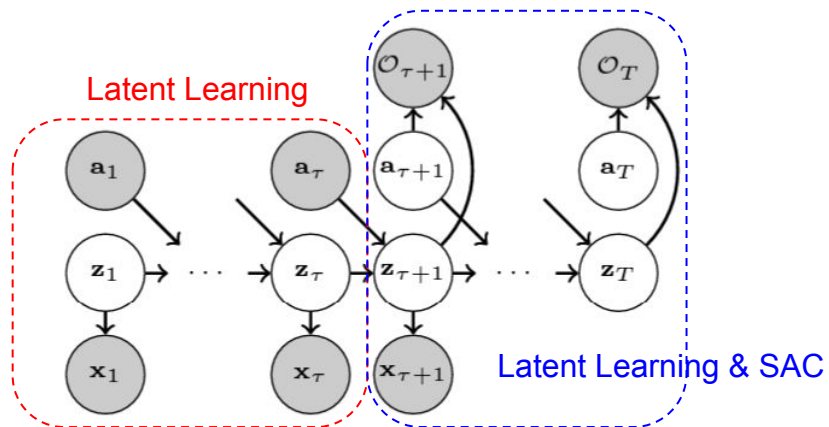


Figure 1: Graphical model of POMDP with optimality variables for $t \geq \tau + 1$.

binary random variable 0 is not optimal

$$p(O_t = 1 | s_t, \mathbf{a}_t) = \exp(r(s_t, \mathbf{a}_t))$$

maximize the marginal likelihood

$$p(\mathbf{x}_{1:\tau+1}, O_{\tau+1:T} | \mathbf{a}_{1:\tau})$$

factorize variational distribution

$$q(\mathbf{z}_{1:T}, \mathbf{a}_{\tau+1:T} | \mathbf{x}_{1:\tau+1}, \mathbf{a}_{1:\tau}) = \prod_{t=0}^{\tau} \underbrace{q(\mathbf{z}_{t+1} | \mathbf{x}_{t+1}, \mathbf{z}_t, \mathbf{a}_t)}_{\text{recognition terms}} \prod_{t=\tau+1}^{T-1} \underbrace{p(\mathbf{z}_{t+1} | \mathbf{z}_t, \mathbf{a}_t)}_{\text{dynamics terms}} \prod_{t=\tau+1}^T \underbrace{\pi(\mathbf{a}_t | \mathbf{x}_{1:t}, \mathbf{a}_{1:t-1})}_{\text{policy terms}}$$

3.Methods - Derivation of the ELBO

$$\log p(\mathbf{x}_{1:\tau+1}, \mathcal{O}_{\tau+1:T} | \mathbf{a}_{1:\tau})$$

$$= \log \int \int p(\mathbf{x}_{1:\tau+1}, \mathcal{O}_{\tau+1:T}, \mathbf{z}_{1:T}, \mathbf{a}_{\tau+1:T} | \mathbf{a}_{1:\tau}) d\mathbf{z}_{1:T} d\mathbf{a}_{\tau+1:T} \quad (13)$$

$$\geq \mathbb{E}_{(\mathbf{z}_{1:T}, \mathbf{a}_{\tau+1:T}) \sim q} \left[\log p(\mathbf{x}_{1:\tau+1}, \mathcal{O}_{\tau+1:T}, \mathbf{z}_{1:T}, \mathbf{a}_{\tau+1:T} | \mathbf{a}_{1:\tau}) - \log q(\mathbf{z}_{1:T}, \mathbf{a}_{\tau+1:T} | \mathbf{x}_{1:\tau+1}, \mathbf{a}_{1:\tau}) \right] \quad (14)$$

$$= \mathbb{E}_{(\mathbf{z}_{1:T}, \mathbf{a}_{\tau+1:T}) \sim q} \left[\underbrace{\sum_{t=0}^{\tau} \left(\log p(\mathbf{x}_{t+1} | \mathbf{z}_{t+1}) - D_{\text{KL}}(q(\mathbf{z}_{t+1} | \mathbf{x}_{t+1}, \mathbf{z}_t, \mathbf{a}_t) \| p(\mathbf{z}_{t+1} | \mathbf{z}_t, \mathbf{a}_t)) \right)}_{\text{model objective terms}} \right. \\ \left. + \underbrace{\sum_{t=\tau+1}^T \left(r(\mathbf{z}_t, \mathbf{a}_t) + \log p(\mathbf{a}_t) - \log \pi(\mathbf{a}_t | \mathbf{x}_{1:t}, \mathbf{a}_{1:t-1}) \right)}_{\text{policy objective terms}} \right], \quad (15)$$

SLAC which maximizes the **ELBO** using function approximators to model the prior and posterior distributions. (Latent variable model & Actor and critic)

3.Methods - Latent Variable Factorization

Solid arrows as part of the **generative model** and all dashed arrows as part of **approximate posterior**.

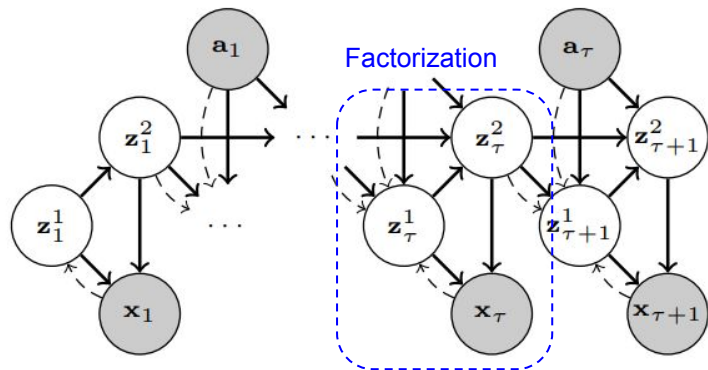


Figure 7: Diagram of our full model. Solid arrows show the generative model, dashed arrows show the inference model. Rewards are not shown for clarity.

The generative model consists of the probability distributions:

$$\begin{aligned} \mathbf{z}_1^1 &\sim p(\mathbf{z}_1^1) && \text{multivariate standard normal distribution } \mathcal{N}(0, \mathbf{I}). \\ \mathbf{z}_1^2 &\sim p_\psi(\mathbf{z}_1^2 | \mathbf{z}_1^1) \\ \mathbf{z}_{t+1}^1 &\sim p_\psi(\mathbf{z}_{t+1}^1 | \mathbf{z}_t^2, \mathbf{a}_t) \\ \mathbf{z}_{t+1}^2 &\sim p_\psi(\mathbf{z}_{t+1}^2 | \mathbf{z}_{t+1}^1, \mathbf{z}_t^2, \mathbf{a}_t) \\ \mathbf{x}_t &\sim p_\psi(\mathbf{x}_t | \mathbf{z}_t^1, \mathbf{z}_t^2) \\ r_t &\sim p_\psi(r_t | \mathbf{z}_t^1, \mathbf{z}_t^2, \mathbf{a}_t, \mathbf{z}_{t+1}^1, \mathbf{z}_{t+1}^2). \end{aligned}$$

two fully connected layers, each with 256 hidden units, and a Gaussian output layer

The variational distribution q , also referred to as the inference model or the posterior, is following factorization

$$\begin{aligned} \mathbf{z}_1^1 &\sim q_\psi(\mathbf{z}_1^1 | \mathbf{x}_1) \\ \mathbf{z}_1^2 &\sim p_\psi(\mathbf{z}_1^2 | \mathbf{z}_1^1) \\ \mathbf{z}_{t+1}^1 &\sim q_\psi(\mathbf{z}_{t+1}^1 | \mathbf{x}_{t+1}, \mathbf{z}_t^2, \mathbf{a}_t) \\ \mathbf{z}_{t+1}^2 &\sim p_\psi(\mathbf{z}_{t+1}^2 | \mathbf{z}_{t+1}^1, \mathbf{z}_t^2, \mathbf{a}_t). \end{aligned}$$

5 convolutional layers (32 5x5, 64 3x3, 128 3x3, 256 3x3, and 256 4x4 filters, respectively, stride 2 each, except for the last layer), 2 fully connected layers (256 units each), and a Gaussian output layer. The parameters of the convolution layers are shared among both distributions.

3.Methods - Compared with Soft Actor Critic(SAC)

SAC learns two Q-networks, a V-network, and a policy network. Two Q-networks are used to mitigate overestimation bias. A V-network is used to stabilize training. Taking gradients through the expectations is done using the reparameterization trick
 Off-Policy(DDPG: ICLR 2016)+Soft Bellman(Soft Q-Learning: ICML 2017) + Stable Actor-Critic(TD3:ICML 2018)

Algorithm 1 Soft Actor-Critic

Initialize parameter vectors $\psi, \bar{\psi}, \theta, \phi$.
for each iteration **do**
 for each environment step **do**
 $\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)$
 $\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$
 $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}_{t+1})\}$
 end for
 for each gradient step **do**
 $\psi \leftarrow \psi - \lambda_V \nabla_\psi J_V(\psi)$
 $\theta_i \leftarrow \theta_i - \lambda_Q \nabla_{\theta_i} J_Q(\theta_i)$ for $i \in \{1, 2\}$
 $\phi \leftarrow \phi - \lambda_\pi \nabla_\phi J_\pi(\phi)$
 $\bar{\psi} \leftarrow \tau \bar{\psi} + (1 - \tau) \psi$
 end for
end for

Algorithm 1 Stochastic Latent Actor-Critic (SLAC)

Require: Environment E and initial parameters
 $\psi, \phi, \theta_1, \theta_2$ for the model, actor, and critics.
 $\mathbf{x}_1 \sim E_{\text{reset}}()$
 $\mathcal{D} \leftarrow (\mathbf{x}_1)$
for each iteration **do**
 for each environment step **do**
 $\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t | \mathbf{x}_{1:t}, \mathbf{a}_{1:t-1})$
 $r_t, \mathbf{x}_{t+1} \sim E_{\text{step}}(\mathbf{a}_t)$
 $\mathcal{D} \leftarrow \mathcal{D} \cup (\mathbf{a}_t, r_t, \mathbf{x}_{t+1})$
 for each gradient step **do**
 $\mathbf{x}_{1:\tau+1}, \mathbf{a}_{1:\tau}, r_\tau \sim \mathcal{D}$
 $\mathbf{z}_{1:\tau+1} \sim q_\psi(\mathbf{z}_{1:\tau+1} | \mathbf{x}_{1:\tau+1}, \mathbf{a}_{1:\tau})$
 $\psi \leftarrow \psi - \lambda_M \nabla_\psi J_M(\psi)$
 $\theta_i \leftarrow \theta_i - \lambda_Q \nabla_{\theta_i} J_Q(\theta_i)$ for $i \in \{1, 2\}$
 $\phi \leftarrow \phi - \lambda_\pi \nabla_\phi J_\pi(\phi)$
 $\bar{\theta}_i \leftarrow \nu \bar{\theta}_i + (1 - \nu) \theta_i$ for $i \in \{1, 2\}$
 end for
end for

3.Methods - Algorithm

Algorithm 1 Stochastic Latent Actor-Critic (SLAC)

Require: Environment E and initial parameters

$\psi, \phi, \theta_1, \theta_2$ for the model, actor, and critics.

$\mathbf{x}_1 \sim E_{\text{reset}}()$ initial observation from the environment

$\mathcal{D} \leftarrow (\mathbf{x}_1)$ initialize replay buffer with initial obs.

for each iteration **do**

for each environment step **do**

$\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t | \mathbf{x}_{1:t}, \mathbf{a}_{1:t-1})$ sample action from the policy

$r_t, \mathbf{x}_{t+1} \sim E_{\text{step}}(\mathbf{a}_t)$ sample transition from the environment

$\mathcal{D} \leftarrow \mathcal{D} \cup (\mathbf{a}_t, r_t, \mathbf{x}_{t+1})$ Store the transition in the replay buffer

for each gradient step **do**

$\mathbf{x}_{1:\tau+1}, \mathbf{a}_{1:\tau}, r_\tau \sim \mathcal{D}$

$\mathbf{z}_{1:\tau+1} \sim q_\psi(\mathbf{z}_{1:\tau+1} | \mathbf{x}_{1:\tau+1}, \mathbf{a}_{1:\tau})$

$\psi \leftarrow \psi - \lambda_M \nabla_\psi J_M(\psi)$

$\theta_i \leftarrow \theta_i - \lambda_Q \nabla_{\theta_i} J_Q(\theta_i)$ for $i \in \{1, 2\}$

$\phi \leftarrow \phi - \lambda_\pi \nabla_\phi J_\pi(\phi)$

$\theta_i \leftarrow \nu \theta_i + (1 - \nu) \theta_i$ for $i \in \{1, 2\}$

Putting it all Together

Both the latent variable model and agent are trained together. SLAC improves the SAC by learning the representation spaces with a latent variable model which is more stable and efficient for complex continuous control tasks. It can improve both the exploration and robustness of the learned model.

$$V_\theta(\mathbf{z}_{\tau+1}) = \mathbb{E}_{\mathbf{a}_{\tau+1} \sim \pi_\phi} [Q_\theta(\mathbf{z}_{\tau+1}, \mathbf{a}_{\tau+1}) - \alpha \log \pi_\phi(\mathbf{a}_{\tau+1} | \mathbf{x}_{1:\tau+1}, \mathbf{a}_{1:\tau})],$$

Model

$$J_M(\psi) = \mathbb{E}_{\mathbf{z}_{1:\tau+1} \sim q_\psi} \left[\sum_{t=0}^{\tau} -\log p_\psi(\mathbf{x}_{t+1} | \mathbf{z}_{t+1}) + D_{\text{KL}}(q_\psi(\mathbf{z}_{t+1} | \mathbf{x}_{t+1}, \mathbf{z}_t, \mathbf{a}_t) \| p_\psi(\mathbf{z}_{t+1} | \mathbf{z}_t, \mathbf{a}_t)) \right]$$

Critic

$$J_Q(\theta) = \mathbb{E}_{\mathbf{z}_{1:\tau+1} \sim q_\psi} \left[\frac{1}{2} (Q_\theta(\mathbf{z}_\tau, \mathbf{a}_\tau) - (r_\tau + \gamma V_\theta(\mathbf{z}_{\tau+1})))^2 \right]$$

Actor

$$J_\pi(\phi) = \mathbb{E}_{\mathbf{z}_{1:\tau+1} \sim q_\psi} \left[\mathbb{E}_{\mathbf{a}_{\tau+1} \sim \pi_\phi} [\alpha \log \pi_\phi(\mathbf{a}_{\tau+1} | \mathbf{x}_{1:\tau+1}, \mathbf{a}_{1:\tau}) - Q_\theta(\mathbf{z}_{\tau+1}, \mathbf{a}_{\tau+1})] \right]$$

3.Methods - Network Architectures

Latent variables $\mathbf{z}_t^1 \in \mathbb{R}^{32}$ and $\mathbf{z}_t^2 \in \mathbb{R}^{256}$.

Image observations $\mathbf{x}_t \in [0, 1]^{64 \times 64 \times 3}$.

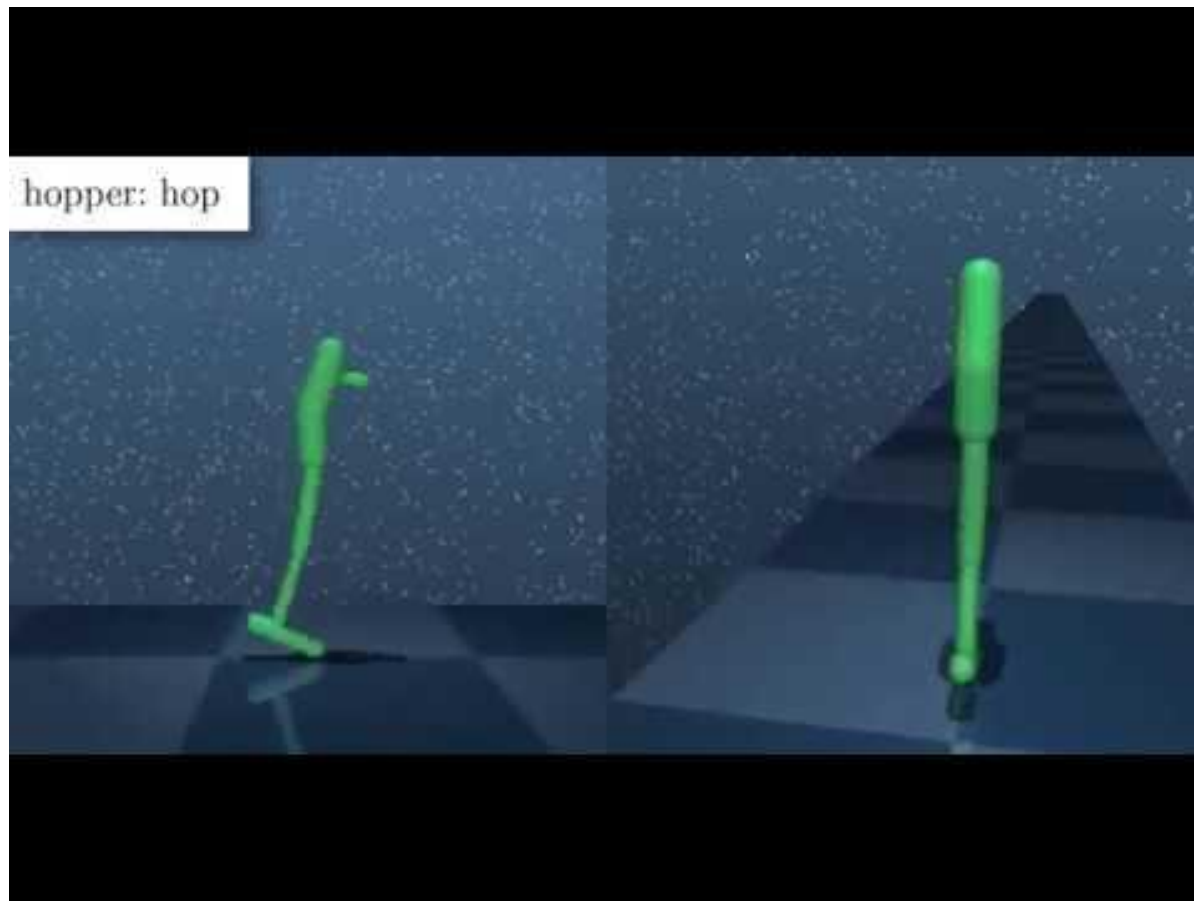
leaky ReLU non-linearities.

Critic network(Q) consisting of 2 fully connected layers (256 units each) and a linear output layer.

Actor network(P) consists of 5 convolutional layers, 2 fully connected layers (256 units each), a Gaussian layer, and a tanh bijector, which constrains the actions to be in the bounded action space of $[-1, 1]$.

The convolutional layers are shared with the ones from the latent variable model, but the parameters of these layers are only updated by the model objective and not by the actor objective.

3.Methods - Experimental Evaluation



The DeepMind Control Suite is a set of **continuous control tasks with a standardised structure** and interpretable rewards, intended to serve as performance benchmarks for reinforcement learning agent

Pixel observations by default, control suite environments return **low-dimensional feature observations**. The pixel.Wrapper adds or replaces these with images

3.Methods - Experimental Evaluation

Image observations for continuous control benchmark tasks

DeepMind
Control Suite



Cheetah run

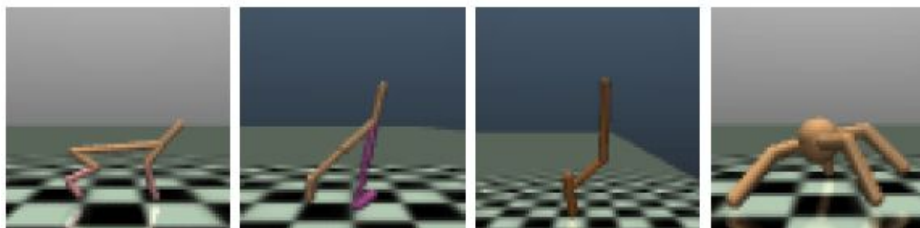
Walker walk

Ball-in-cup catch

Finger spin

64 x 64 pixels

OpenAI
Gym's



Cheetah

Walker

Hopper

Ant

9-DoF
3-fingered
DClaw robot



Push a door

Close a drawer

Reach out / pick up an object

SLAC algorithm
achieves above action

3.Methods - Related Works

Comparative Evaluation on Continuous Control Benchmark Tasks

SAC : Off-policy actor-critic algorithm, model-free learning. We include experiments showing the performance of SAC based on true state (upper bound on performance) from raw images.

D4PG : Off-policy actor-critic algorithm, learning directly from images

MPO : Off-policy actor-critic algorithm that performs an expectation maximization form of policy iteration, learning directly from raw images.

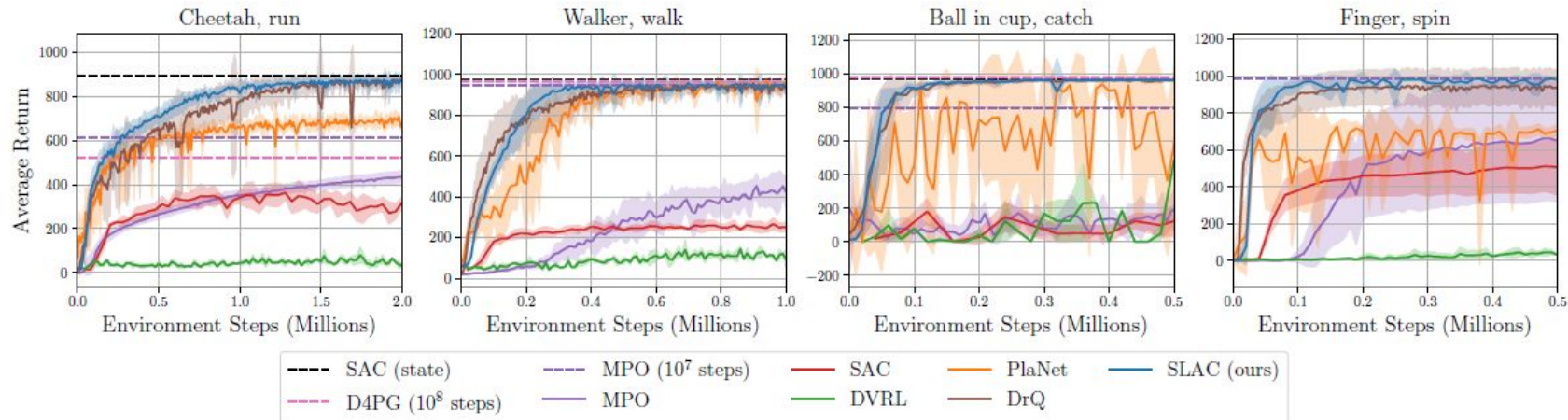
DVRL : On-policy model-free RL algorithm, Mixed deterministic/stochastic latent-variable POMDP model as opposed to our method, which trains the critic with the latent state and the actor with a history of actions and observations

PlaNet : Model-based RL method for learning from images, which uses a partially stochastic sequential latent variable model, but without explicit policy learning.

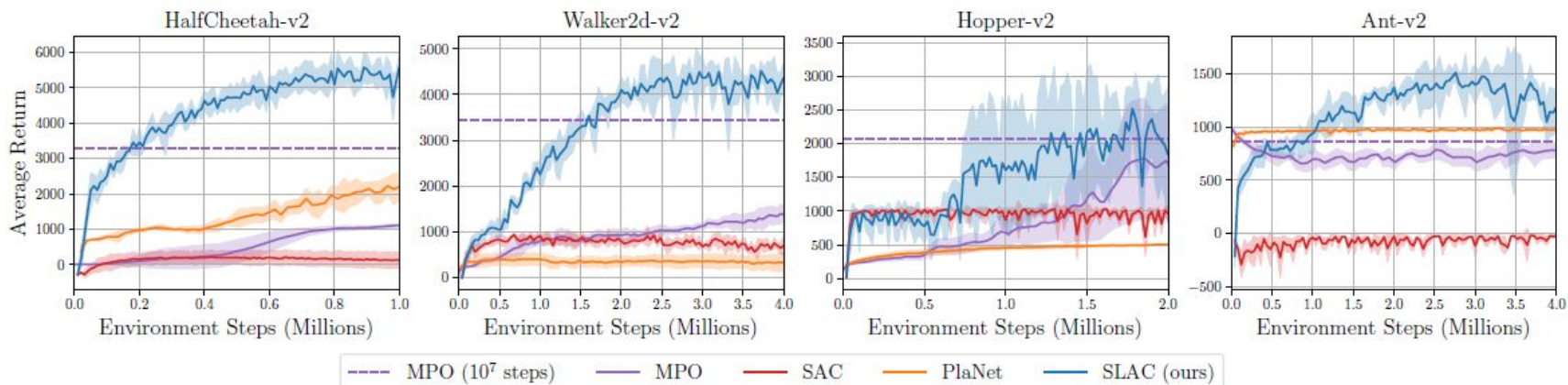
DrQ : the same as the SAC algorithm, but combined with data augmentation on the image inputs.

3.Methods - Experiments

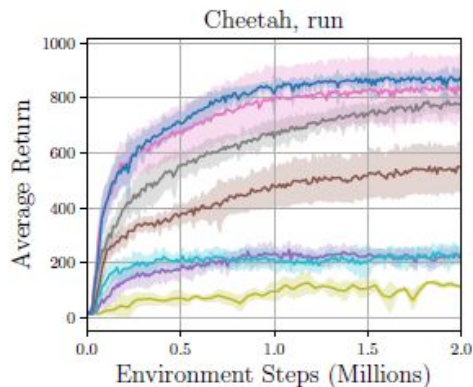
DeepMind Control Suite from images



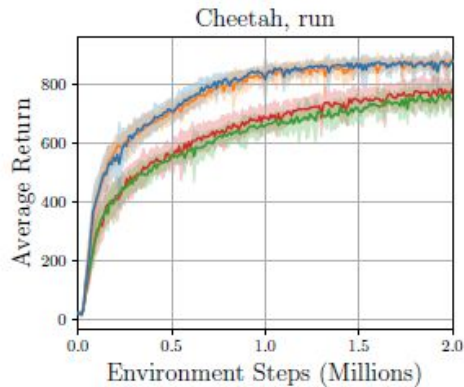
OpenAI Gym benchmark tasks from images



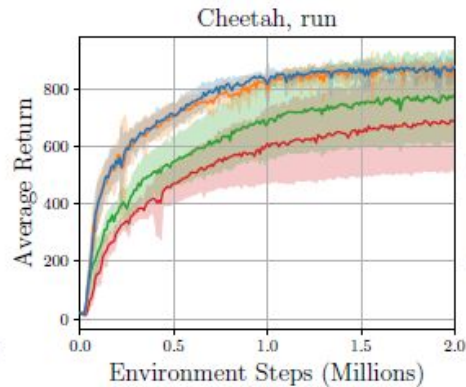
3.Methods - Ablation Experiments



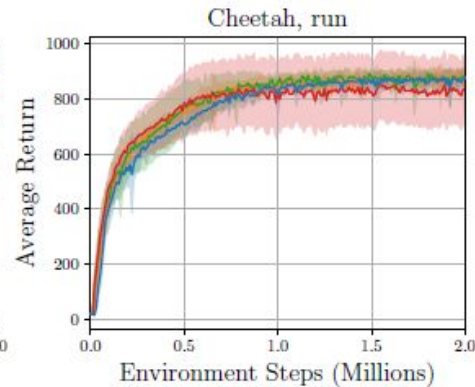
(a) Latent variable model



(b) Actor and critic inputs



(c) Model pretraining



(d) Train steps per iteration

(a) **Latent variable model** : 1) fully stochastic model to a standard non-sequential VAE model, 2) two-variable factorization but without any temporal dependencies, 3) sequential filtering model that uses temporal dependencies but without the two-variable factorization, 4) without the two-variable factorization, the partially stochastic model used by PlaNet

(b) **Actor and critic inputs** : actor and critic inputs as either the observation-action history or the latent sample

(c) **Model pre-training** : the agent first collects a small amount of data by executing random actions, and then the model is pre-trained with that data

(d) **Training updates per iteration** : the effect of the number of training updates per iteration

4. Conclusion

[Contribution]

Combined [off-policy model-free RL with representation learning](#) via a [sequential stochastic](#) state model.

SLAC algorithm for learning from [high-dimensional image inputs](#). (learn directly from raw image observations)

Previous work have used mixed deterministic-stochastic models, but SLAC's model is purely stochastic

SLAC's [fully stochastic model outperforms other latent variable models](#).

- Achieved improved sample efficiency and final task performance
(Four DeepMind Control Suite tasks and four OpenAI Gym tasks)

Simulation on robotic manipulation tasks (9-DoF 3-fingered DCIaw robot on four tasks)

[Limitation]

In real-world setups such as robotics how to choice reward function?

It takes lots of time to train latent space. (Change environment and complexity)

For fairness, performance evaluations for other models seems necessary

(not just SLAC RL framework, compare on different latent variable models)

Performance on other image-based continuous control tasks

States benefits of using two layers of latent variables (Insufficient explanation on why it brings good balance)

Paper that are [not properly backed up by evidence](#), it is not sufficiently clear to what degree the shown performance improvement is due to the stochastic nature of the model used.

Official Blind Review #4 (Rating: 3: Weak Reject)

- There are quite a few recent works about representation learning and yet the paper makes no references to previous works. I find it surprising that [none of the past related works are mentioned in the paper](#). It does not explicitly require references to a POMDP. I do not understand why the [critic evaluation in the latent space is even a good approach?](#) It is more of an experimental design and engineering approach that combines previous known techniques.

Official Blind Review #2 (Rating: 8: Accept)

- Operates on a learned latent state, rather than an observed one, and therefore aims to jointly learn to represent high dimensional inputs and execute continuous control based on this representation. ([Novel formulation, SOTA results, well written](#)) the role of making the primary [latent variable stochastic](#), is investigated, a deeper investigation of what makes the model more effective than existing techniques would be insightful.

Official Blind Review #3 (Rating: 3: Weak Reject)

- The real-benefit of the method may not be consistent. (My biggest concern) Full stochasticity contributed to the practical gain but in the experiment Figure 6, The [simple filtering does not perform well](#). It seems that the benefit of the method is rather from such [particular latent space design](#) rather than [the stochastic vs deterministic](#). the [experiments may be unfair](#), because, another [partially stochastic](#) method can easily utilize such design and further improve the performance.

Simple Filtering (without factoring model)

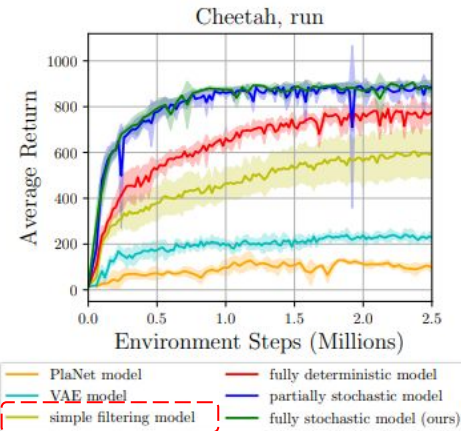


Figure 6: Comparison of different design choices for the latent variable model.

References

ICLR 2020 Conference homepage : <https://openreview.net/forum?id=HJxDugSFDB>

VAE : <https://lilianweng.github.io/posts/2018-08-12-vae/>

Zhang, Cheng, et al. "Advances in variational inference." *TPAMI 2018* <https://arxiv.org/pdf/1711.05597.pdf>

POMDP : <http://www.pomdp.org/faq.html>

https://www.techfak.uni-bielefeld.de/~skopp/Lehre/STdKI_SS10/POMDP_tutorial.pdf

CS 330: Deep Multi-Task and Meta Learning : <http://cs330.stanford.edu/fall2019/index.html>

Hafner, Danijar, et al. "Learning latent dynamics for planning from pixels." ICML 2019 <https://planetrl.github.io/>

Thanks

Any Questions?

You can send mail to
Susang Kim(healess1@gmail.com)