

Panasonic

Deep Generative Models for Robot Control

Masashi OKADA, Ph.D

Senior staff engineer, Panasonic

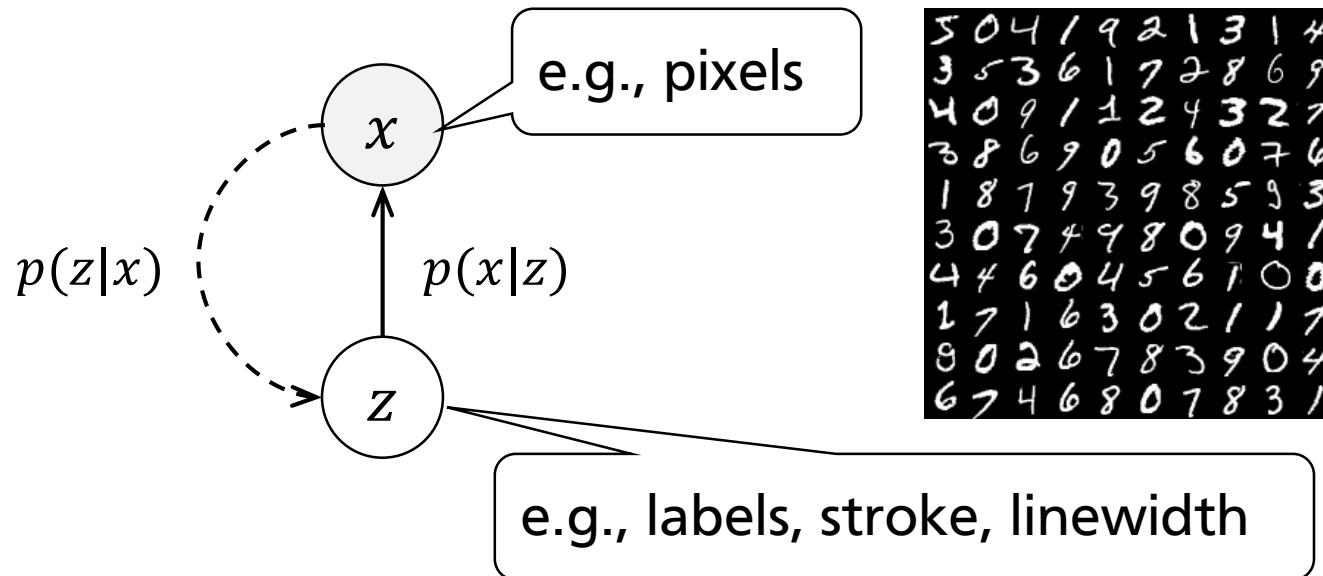
Lecture @ IROS2020 Tutorial sponsored by RSJ
Tutorial on Deep Probabilistic Generative Models for Robotics

Table of Contents

- A basic deep generative model:
variational autoencoder (VAE)
 - Theory
 - Applications
 - World model, PlaNet, Dreamer, etc...
- Control as probabilistic inference
 - Theory
 - Applications
 - SAC, SLAC, VI-MPC, PaETS, PlaNet-Bayes

Problem formulation of VAE

- Variables
 - Observation : x
 - Latent : z
- What we want:
 - $p(x, z)$
 - $= p(z|x)p_{\mathcal{D}}(x)$
 - We can infer the latent z from an observation x
 - $= p(x|z)p(z)$
 - We can generate samples x from latent z drawn from a prior $p(z)$



How to derive the models?

- Define parameterized deep generative models
 - $p(x, z) = p_\theta(x|z)p(z)$
 - θ : model parameter (i.e., weight of neural networks)
 - $p(z)$: prior (e.g., standard normal dist.)
- Find θ^* that maximizes the data log-likelihood
(MLE: Maximum Likelihood Estimation)
 - $\theta^* = \operatorname{argmax}_\theta \log p(x) = \operatorname{argmax}_\theta \log \int p_\theta(x|z)p(z)dz$
- However, this is intractable, because ...
 - No analytic solution
 - No efficient estimator

How to make it tractable

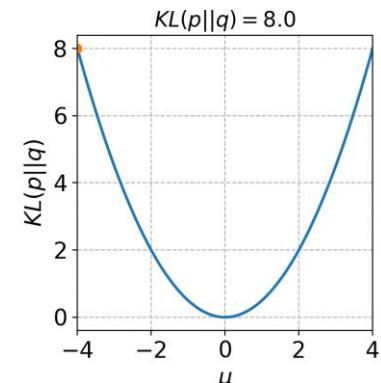
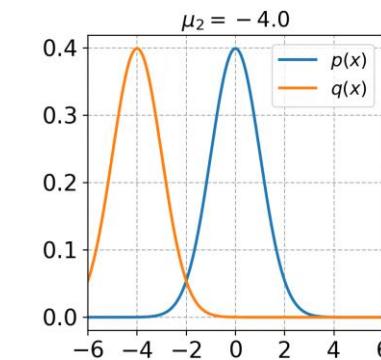
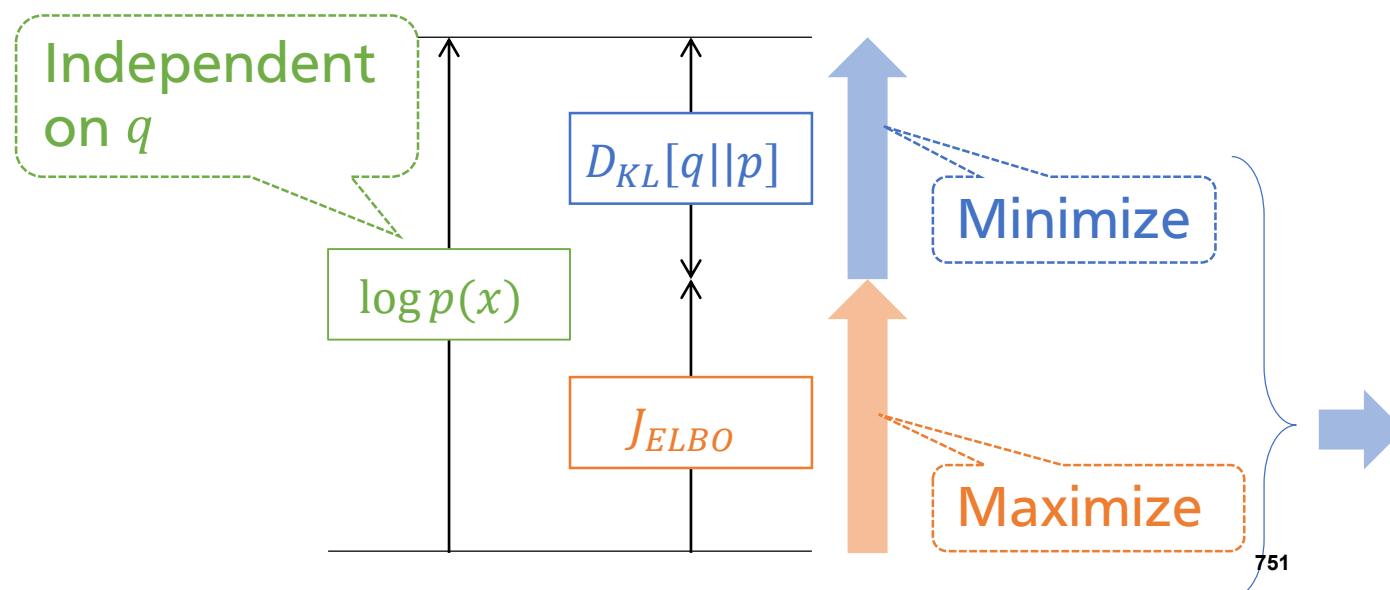
- $\log p(x) = \log \int p_\theta(x|z)p(z)dz$
- $= \log \int p_\theta(x|z)p(z) \frac{q_\phi(z|x)}{q_\phi(z|x)} dz$
- $\geq \int q_\phi(z|x) \log \frac{p_\theta(x|z)p(z)}{q_\phi(z|x)} dz$
- $\coloneqq J_{\text{ELBO}}$ (ELBO: Evidence Lower Bound)
 - J_{ELBO} can be efficiently optimized by using general deep learning techniques (discussed later)

$q_\phi(z|x)$: *variational distribution*
to approx. the posterior $p(z|x)$

Jensen's inequality
 $\log \mathbb{E}[x] \geq \mathbb{E}[\log x]$

What ELBO also optimizes?

- $\log p(x) - J_{\text{ELBO}} = D_{KL}[q_{\phi}(z|x)||p(z|x)]$
 - $D_{KL}[q||p] \left(:= \int q \log \frac{q}{p} dz \right)$: Kullback-Leibler divergence (KL-divergence)
 - A measure of how q, p are different



<https://qiita.com/ceptrree/items/9a473b5163d5655420e8>

ELBO max. is equivalent to variational inference $D_{KL}[q||p]$

How to compute J_{ELBO}

$$\begin{aligned} J_{\text{ELBO}} &:= \int q_{\phi}(z|x) \log \frac{p_{\theta}(x|z)p(z)}{q_{\phi}(z|x)} dz \\ &= \int q_{\phi}(z|x) \log p_{\theta}(x|z) dz - \int q_{\phi}(z|x) \log \frac{q_{\phi}(z|x)}{p(z)} dz \end{aligned}$$

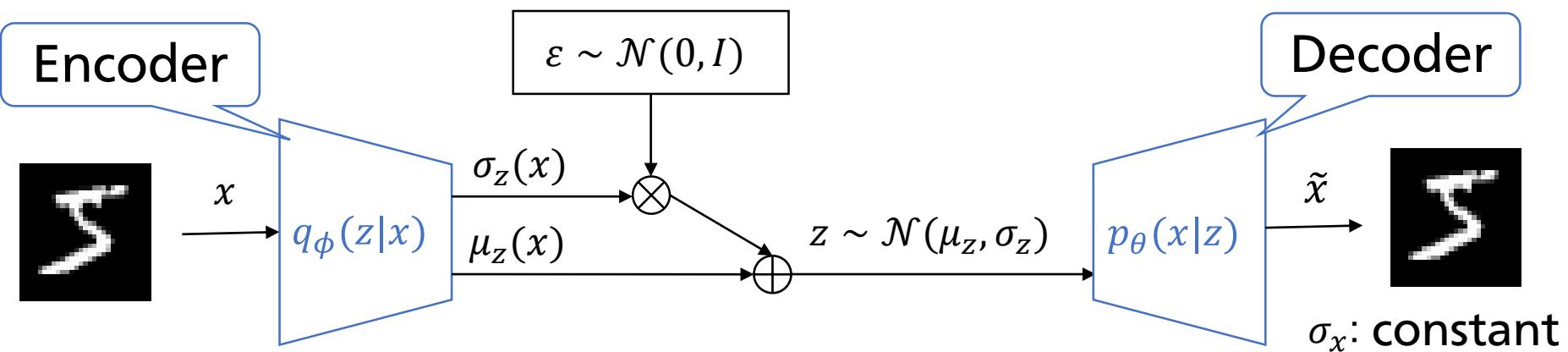
$$= \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{KL}[q_{\phi}(z|x) || p(z)]$$

Log-likelihood of
autoencoded x

Regularizer to encourage $q_{\phi}(z|x)$
to be close to the prior $p(z)$

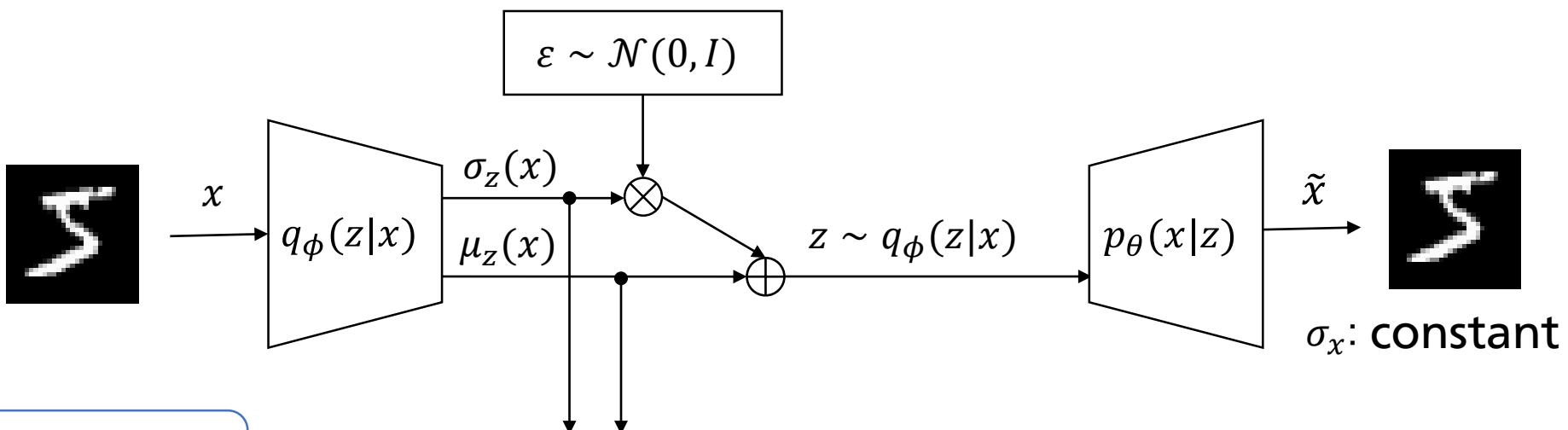
How to compute J_{ELBO}

- Computational Graph with Deep Gaussian Modeling of p_θ, q_ϕ
 - $J_{\text{ELBO}} = \mathbb{E}_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}[q_\phi(z|x) || p(z)]$



How to compute J_{ELBO}

- Computational Graph with Deep Gaussian Modeling of p_θ, q_ϕ
 - $J_{\text{ELBO}} = \mathbb{E}_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}[q_\phi(z|x) || p(z)]$



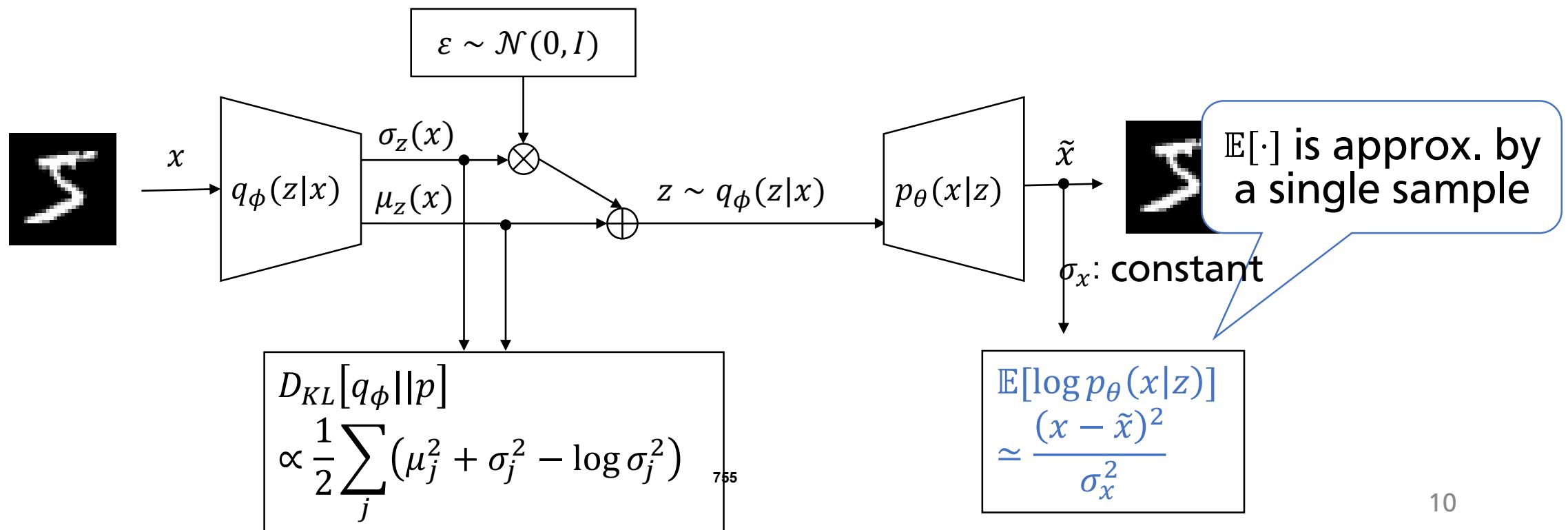
$p(z) := \mathcal{N}(0, I)$

$$D_{KL}[q_\phi || p] \propto \frac{1}{2} \sum_j (\mu_j^2 + \sigma_j^2 - \log \sigma_j^2)$$

754

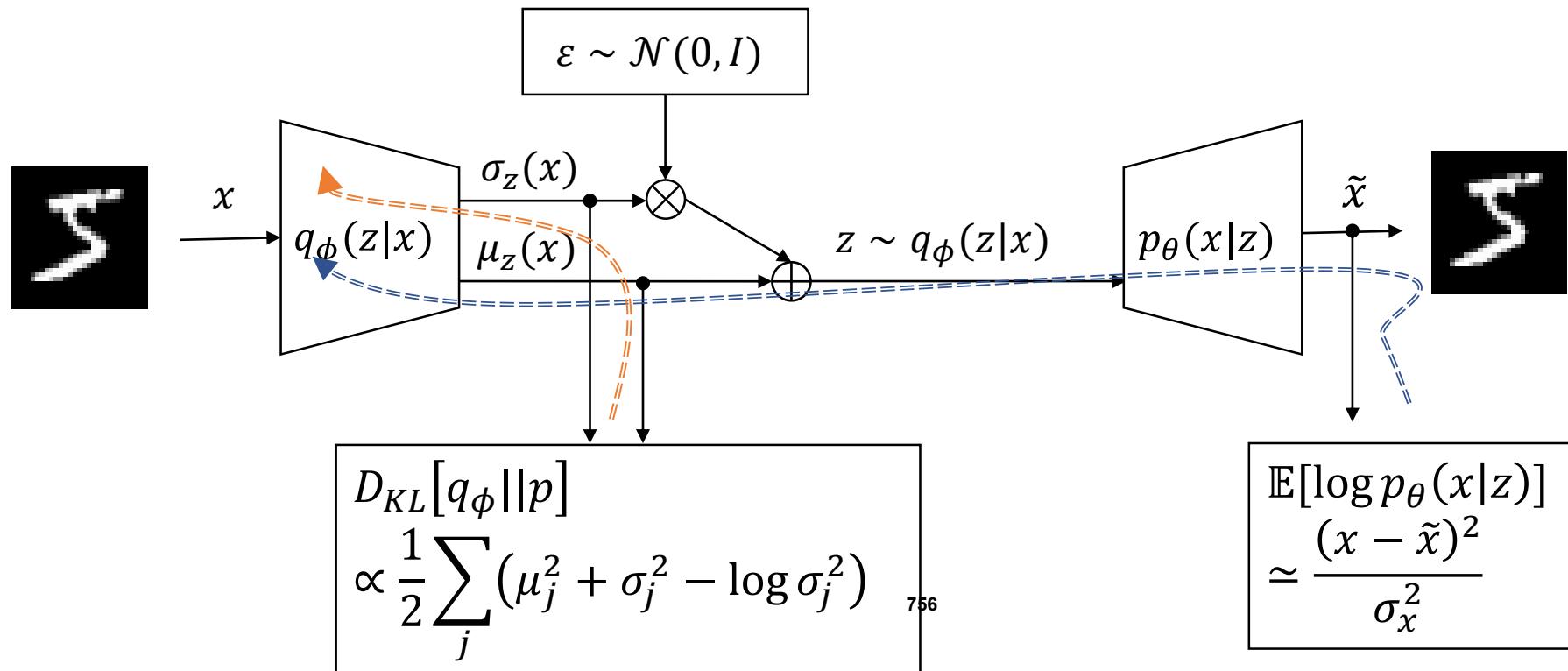
How to compute J_{ELBO}

- Computational Graph with Deep Gaussian Modeling of p_θ, q_ϕ
 - $J_{\text{ELBO}} = \mathbb{E}_{z \sim q_\phi(z|x)}[\log p_\theta(x|z)] - D_{KL}[q_\phi(z|x) || p(z)]$

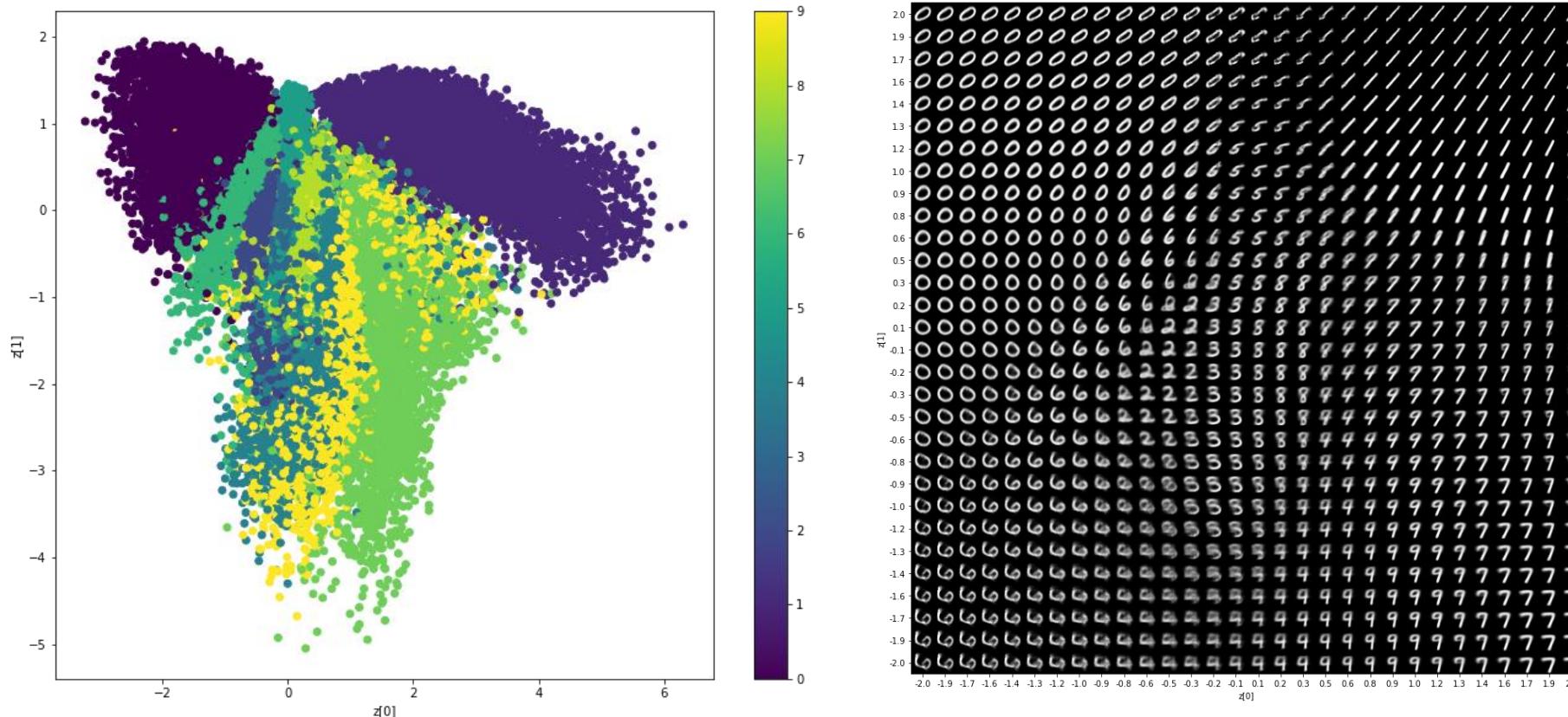


How to compute J_{ELBO}

All forward operations are differentiable
 $\Rightarrow \nabla_{\theta, \phi} J_{\text{ELBO}}$ can be efficiently computed by back-prop.



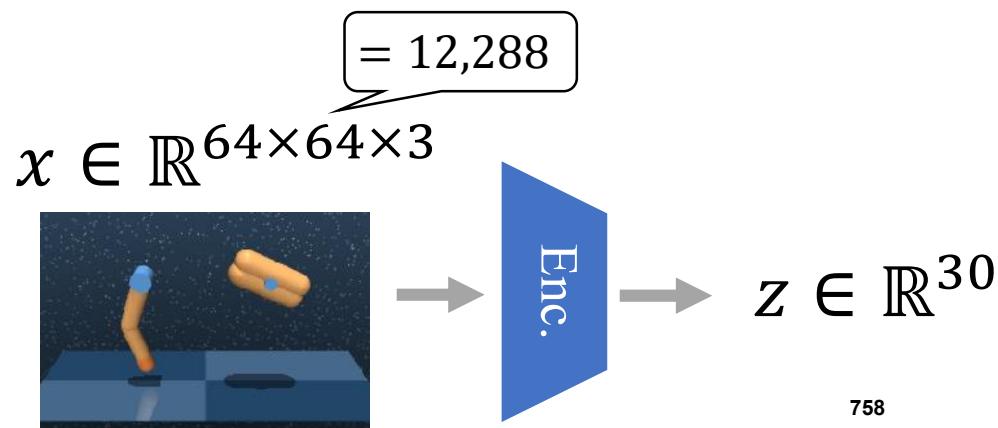
How MNIST data are encoded



<https://keras.io/examples/generative/vae/> ⁷⁵⁷

Application of VAE to robotics

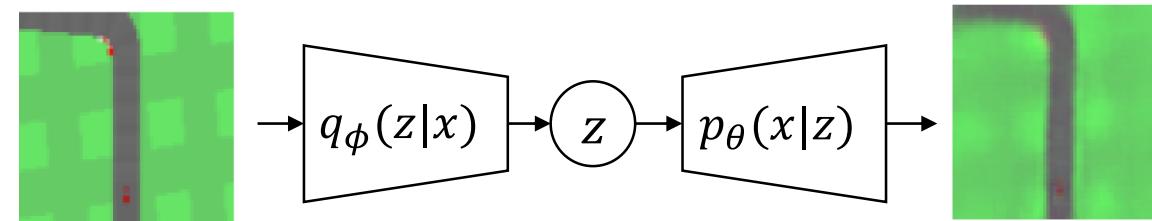
- Control from high-dimensional input (e.g., pixels)
- Compact latent representation makes planning and policy optimization much easier
 - Simple control schemes are also applicable; e.g.,
 - Linear programming [Water+, NeurIPS2015]
 - Linear controller [Ha+, NeurIPS2018]



758

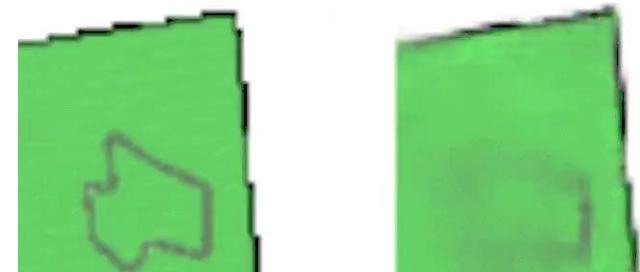
World Model [Ha+, NeurIPS2018]

- Model-based reinforcement learning utilizing VAE
- Procedure
 - 1. Train a VAE



- 2. Train a latent dynamics model
 $p(z_{t+1}|z_t, a_t)$

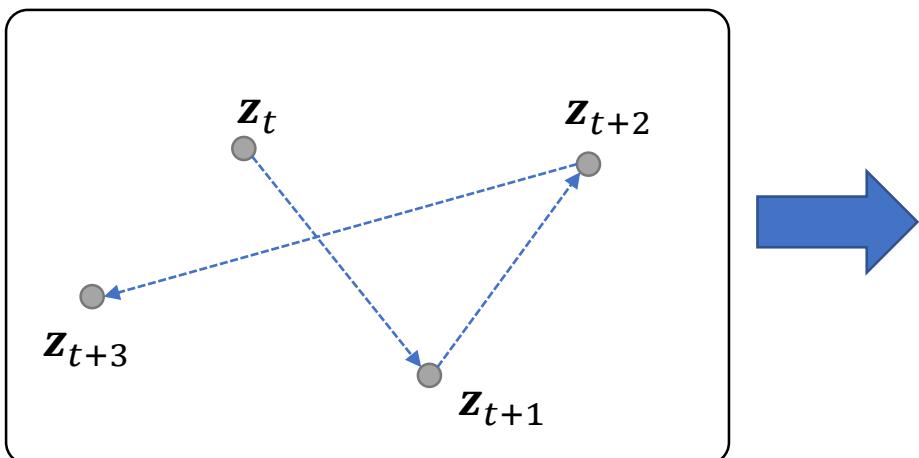
- 3. Train a policy utilizing *imagined* trajectories
(learning inside of a dream)



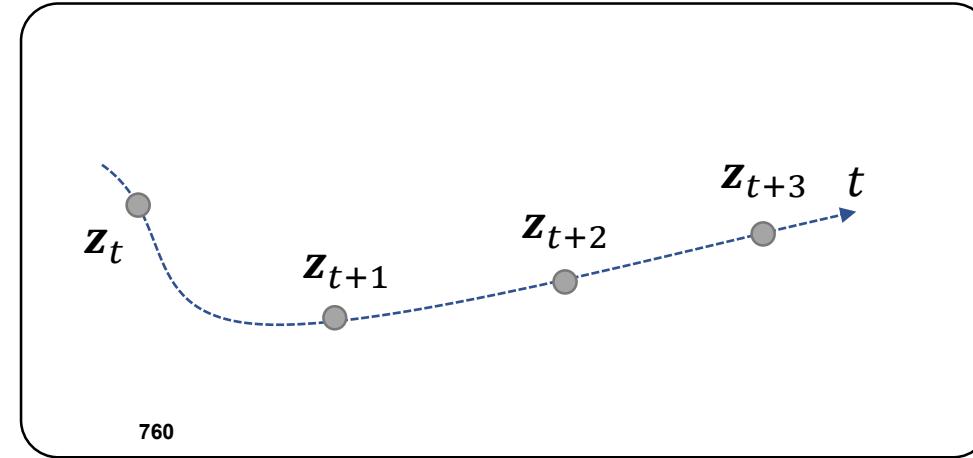
Is (vanilla-)VAE enough for control?

- World Model trains VAE and dynamics **independently**
- VAE does not consider time correlation
 - Consecutive latent might be distributed unsmoothly
 - How to incorporate time correlation for more smooth space?

Latent space

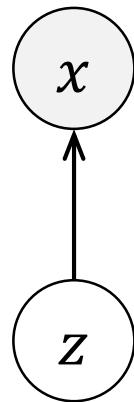


Latent space



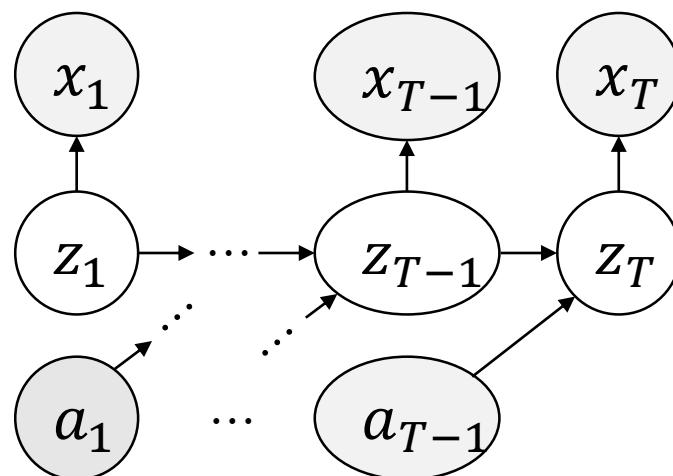
Time-series VAE

PGM of vanilla-VAE



$$J_{\text{ELBO}} := \mathbb{E}_{q(z|x)}[\log p(x|z)] - D_{KL}[q(z|x)||p(z)]$$

PGM of time-series VAE



a.k.a. POMDP
(Partially observable
Markov Decision Process)

$$J_{\text{ELBO}} := \sum_t \mathbb{E}_{q(z_t|x_{\leq t}, a_{<t})}[\log p(x_t|z_t)] - \sum_t \mathbb{E}_{q(z_{t-1}|\cdot)}^{761}[D_{KL}[q(z_t|x_{\leq t}, a_{<t})||p(z_t|z_{t-1}, a_{t-1})]]$$

Time-series VAE

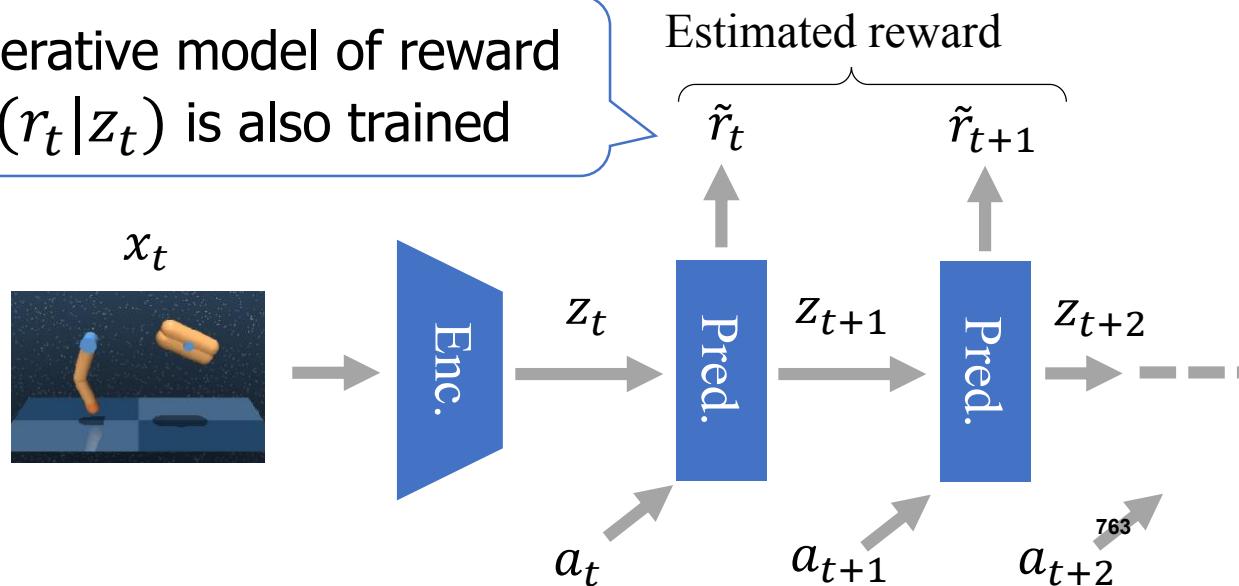
- $J_{\text{ELBO}} := \sum_t \mathbb{E}_{q(z_t|x_{\leq t}, a_{<t})} [\log p(x_t|z_t)] - \sum_t \mathbb{E}_{q(z_t|x_{\leq t}, a_{<t})} [D_{KL}[q(z_t|x_{\leq t}, a_{<t})||p(z_t|z_{t-1}, a_{t-1})]]$
 - $q(z_t|x_{\leq t}, a_{<t})$: encoder
 - $p(x_t|z_t)$: decoder
 - $p(z_t|z_{t-1}, a_{t-1})$: dynamics model
- ⇒ We can jointly train VAE and dynamics
- ⇒ Smooth latent space is constructed so that consecutive latent is easily predictable

Deep Planning Network (PlaNet)

[Hafner+, ICML2018]

- Planning (trajectory optimization) in latent space utilizing a learned time-series VAE
- Optimize $a_{\geq t}$ to maximize the predicted return
 - Opt. method: Cross-entropy Method (Monte Carlo based)

Generative model of reward
 $p(r_t|z_t)$ is also trained



Dreamer [Hafner+, ICLR2019]

- Policy optimization in latent space utilizing learned time-series VAE
- Train parametrized policy π_θ (and value function) by backpropagating the estimated return

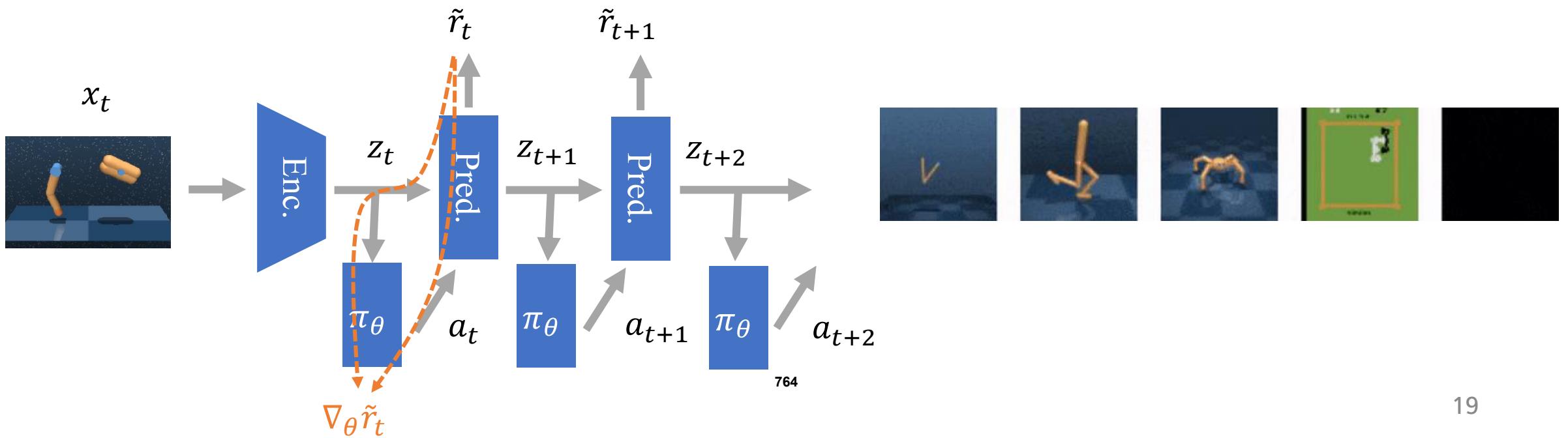


Table of Contents

- A basic deep generative model:
variational autoencoder (VAE)
 - Theory
 - Applications
 - World model, PlaNet, Dreamer, etc...
- Control as probabilistic inference
 - Theory
 - Applications
 - SAC, SLAC, VI-MPC, PaETS, PlaNet-Bayes

The screenshot shows a search results page from Google Scholar. The search term is "Reinforcement learning and control as probabilistic inference: Tutorial and review". There are approximately 63,000 results. The top result is a paper by S. Levine on arXiv.org, which provides a mathematical formalization of intelligent decision making. The second result is a paper by C. Daniel et al. in Machine Learning, Springer, which discusses probabilistic inference for determining options in reinforcement learning. The third result is a paper by M. Toussaint in KI, researchgate.net, which explores probabilistic inference as a model of planned behavior.

Scholar About 63,000 results (0.22 sec) YEAR

Reinforcement learning and control as probabilistic inference: Tutorial and review [PDF] arxiv.org

S. Levine - arXiv preprint arXiv:1805.00909, 2018 - arxiv.org

The framework of **reinforcement learning** or optimal **control** provides a mathematical formalization of intelligent decision making that is powerful and broadly applicable. While the general form of the **reinforcement learning** problem enables effective reasoning about ...

☆ 99 Cited by 144 Related articles All 3 versions

[HTML] Probabilistic inference for determining options in reinforcement learning [HTML] springer.com Full View

C. Daniel, H. Van Hoof, J. Peters, G. Neumann - Machine Learning, 2016 - Springer

... **Machine Learning** ... of this region, the imitation **learning** solution will not be able to successfully solve the task and a **reinforcement learning** solution is ... In the proposed method, the activation policy will **learn** to initiate sub-policies according to their responsibility of state-space ...

☆ 99 Cited by 70 Related articles All 13 versions

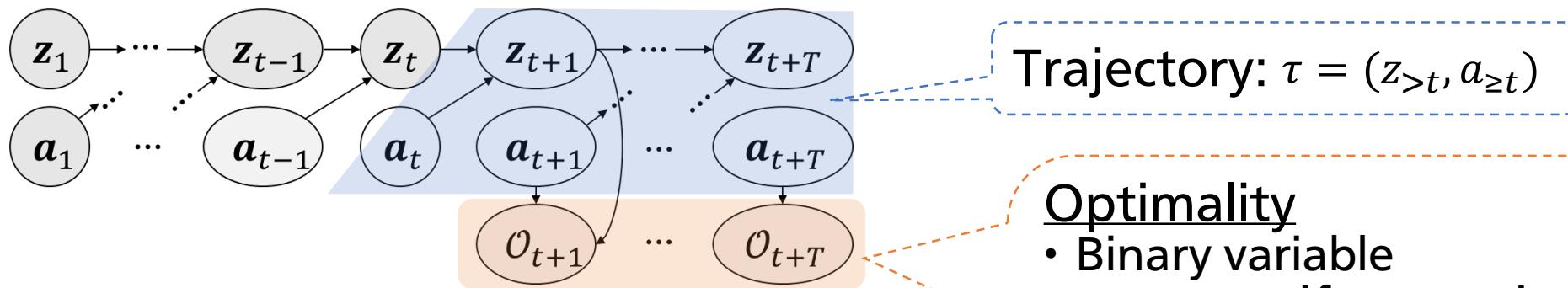
[PDF] Probabilistic inference as a model of planned behavior [PDF] researchgate.net

M. Toussaint - KI, 2009 - researchgate.net

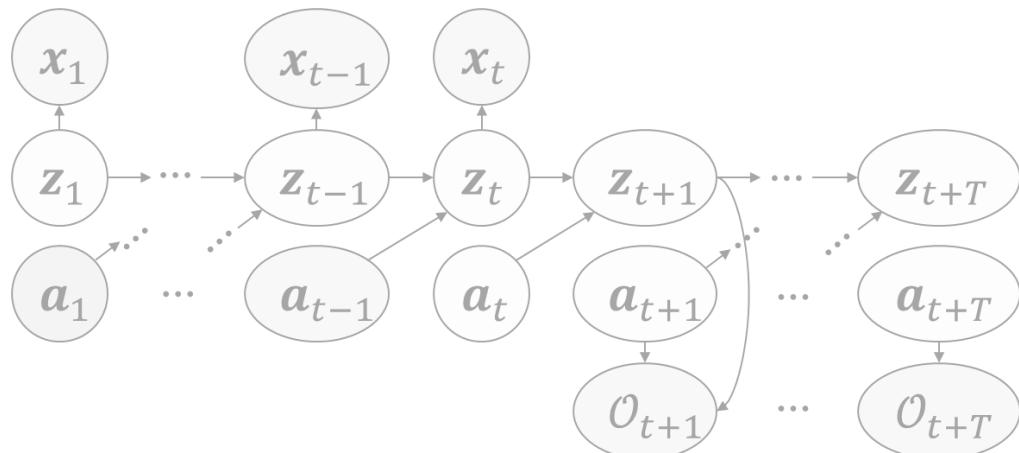
... Recently there is increasing efforts to marry classical AI representations with the **probabilistic** frame-work and, in fact, it is possible to express and **learn** compact models of the ... Using expectation maximiza-tion for **reinforcement learning** ... on **Machine Learning** (ICML 2009), 2009

Problem Formulation

Fully Observable Markov Decision Process (MDP)



Partially Observable MDP (discussed later)



Trajectory: $\tau = (z_{>t}, a_{\geq t})$

Optimality

- Binary variable
 - $\mathcal{O}_t = 1$ if $r(z_t, a_t)$ is optimal
 - $r(z_t, a_t)$: reward function
- $p(\mathcal{O}_t = 1 | z_t, a_t) := f(r(z_t, a_t))$
- $f(\cdot)$: increasing function (discussed later)

What we want:
 $p(\tau | \mathcal{O}_{>t} = 1, z_t)$

will be omitted for readability

How to infer $p(\tau | \mathcal{O}_{>t} = 1)$

- Analytical inference $p(\tau | \mathcal{O}_{>t} = 1)$ is intractable
⇒ Variational inference $D_{KL}[q(\tau) || p(\tau | \mathcal{O}_{>t} = 1)]$ by ELBO maximization
- How to define a variational distribution $q(\tau)$
 - $q(\tau) := q(a_{\geq t}) \prod_t p(z_{t+1} | z_t, a_t)$
 - To derive trajectory optimization
 - $q(\tau) := \prod_t p(z_{t+1} | z_t, a_t) \pi(a_t | s_t)$
 - To derive policy optimization
 - Note: $p(z_{t+1} | z_t, a_t)$ governs the state transition

Trajectory optimization as inference

- Derivation of ELBO

- $\log p(\mathcal{O}_{>t} = \mathbf{1}) = \log \int p(\tau, \mathcal{O}_{>t} = \mathbf{1}) d\tau$

- $= \log \int p(\tau, \mathcal{O}_{>t} = \mathbf{1}) \cdot \frac{q(\tau)}{q(\tau)} d\tau$

- $= \log \mathbb{E}_{q(\tau)} \left[\frac{p(\tau, \mathcal{O}_{>t} = \mathbf{1})}{q(\tau)} \right]$

- $= \log \mathbb{E}_{q(\tau)} \left[\frac{\prod p(z_{t+1}|z_t, a_t) \cdot \prod p(\mathcal{O}_t = 1|z_t, a_t)}{q(a_{\geq t}) \prod p(z_{t+1}|z_t, a_t)} \right]$

- $\geq \mathbb{E}_{q(\tau)} [\log \prod p(\mathcal{O}_t = 1|z_t, a_t) - \log q(a_{\geq t})] := J_{\text{ELBO}}$

$q(\tau)$: variational distribution

$p(a_{\geq t})$: uninformative prior

Jensen's inequality

$\sum r(z_t, a_t)$

if $p(\mathcal{O}_t|z_t, a_t) := \exp(r(z_t, a_t))$

Entropy of $q(a_{\geq t})$

Max. J_{ELBO}
=

Max. Return + Entropy

Variational Inference MPC

[Okada+, CoRL2019]

- Max. $J_{ELBO} := \mathbb{E}_{q(\tau)}[\log p(\mathcal{O}_{>t} = 1 | \tau) - \log q(a_{\geq t})]$

↓ Mirror Descent

$$:= \prod p(\mathcal{O}_t | z_t, a_t)$$

- VI-MPC (variational inference model predictive control)

$$q(a) \leftarrow \frac{q(a) \cdot \mathbb{E}_\tau[p(\mathcal{O}_{>t} = 1 | \tau)] \cdot q(a)^{-\kappa}}{\mathbb{E}_{q(a)}[\mathbb{E}_\tau[p(\mathcal{O}_{>t} = 1 | \tau)] \cdot q(a)^{-\kappa}]}$$

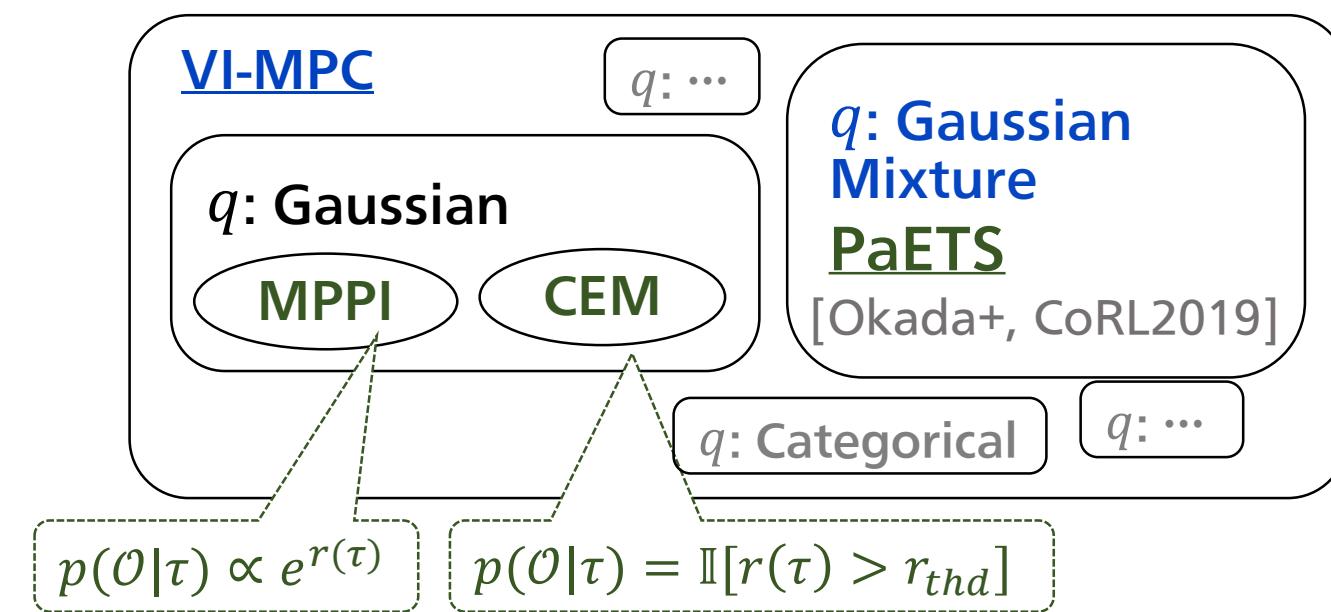
- This generalizes various MPC methods

$p(\mathcal{O}_{>t} \tau)$	MPPI [18]	CEM [16]	Prop-CEM [20]	CMA-ES [19]
$f(r(\tau))$	$\propto e^{r(\tau)}$	$\mathbb{1}[r(\tau) > r_{thd}]$	$\frac{r(\tau) - r_{min}}{r_{max} - r_{min}}$	$\propto \log g(r(\tau)) \cdot \mathbb{1}[r(\tau) > r_{thd}]$

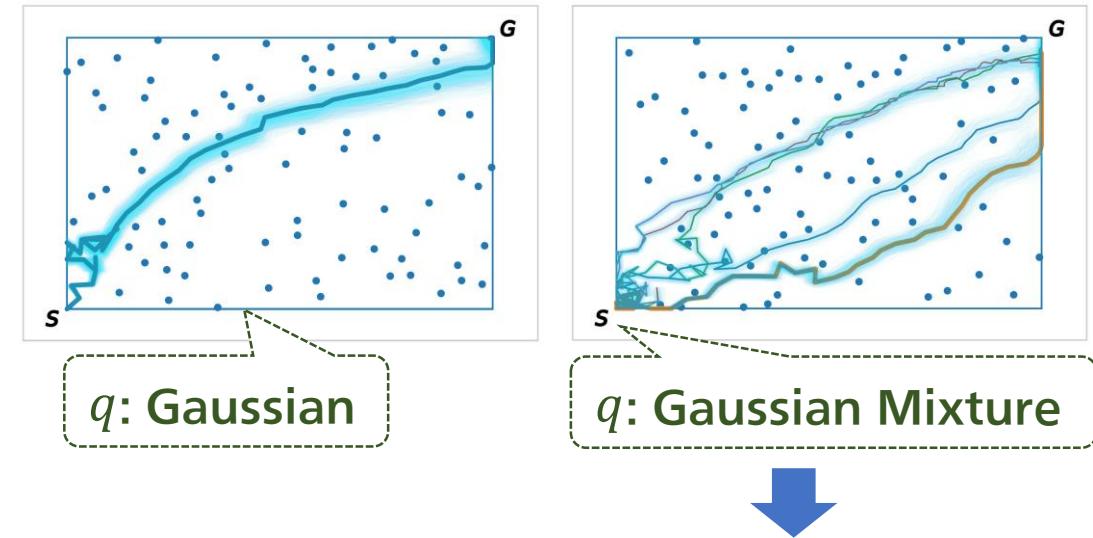
Model Predictive Path Integral Cross Entropy Method ⁷⁶⁹

VI-MPC & PaETS [Okada+, CoRL2019] (Probabilistic Action Ensemble with Traj. Sampling)

- Different definitions of $p(\mathcal{O}|\tau)$ & $q(a)$ derive various MPC methods



e.g., Toy navigation task



Encourage active exploration
in reinforcement learning₂₅

Policy optimization as inference

- Derivation of ELBO

- $\log p(\mathcal{O}_{>t} = 1) \geq$

- $\mathbb{E}_{\tau \sim q(\tau)} [\sum r(z_t, a_t) - \sum \log \pi(a_t | z_t)] := J_{\text{ELBO}}$

$$p(\mathcal{O}_t | z_t, a_t) := \exp(r(z_t, a_t))$$

Entropy of π

$$q(\tau) := \prod p(z_{t+1} | z_t, a_t) \pi(a_t | z_t)$$

- Solve ELBO by dynamic programming

- $J_{\text{ELBO}} = \mathbb{E}_{a_t \sim \pi(a_t | z_t)} [r(z_t, a_t) - \log \pi(a_t | z_t) + \mathbb{E}_{z_{t+1} \sim p(z_{t+1} | z_t, a_t)} [V(z_{t+1})]]$

- $= \mathbb{E}_{a_t \sim \pi(a_t | z_t)} [Q(z_t, a_t) - \log \pi(a_t, z_t)]$

- where

- $V(z_{t+1}) := \mathbb{E}_{\pi(a_{t+1} | z_{t+1})} [r(z_{t+1}, a_{t+1}) - \log \pi(a_{t+1}, z_{t+1}) + \mathbb{E}_{p(z_{t+2} | \cdot)} [V(z_{t+2})]]$

- $Q(z_t, a_t) := r(z_{t+1}, a_{t+1}) + \mathbb{E}_{p(z_{t+1} | z_t, a_t)} [V(z_{t+1})]$

Policy optimization as inference

- Solve ELBO by dynamic program
 - $J_{\text{ELBO}} = \mathbb{E}_{a_t \sim \pi(a_t|z_t)} [Q(z_t, a_t) - \log \pi(a_t, z_t)]$
 - $= \mathbb{E}_{a_t \sim \pi(a_t|z_t)} [\log e^{Q(z_t, a_t)} - \log \pi(a_t, z_t)]$
 - $= \mathbb{E}_{a_t \sim \pi(a_t|z_t)} [\log e^{Q(z_t, a_t)} - \log \pi(a_t, z_t) + \log Z - \log Z]$
 - $Z (= \int e^Q da)$: normalizer
 - $\propto -D_{KL} [\pi(a_t|z_t) || \frac{e^{Q(z_t, a_t)}}{Z}]$
- Optimal policy
 - $\pi^*(a_t|z_t) = \frac{e^{Q(z_t, a_t)}}{Z}$
 - Boltzmann distribution

Soft Actor Critic (SAC)

[Haarnoja+, ICML2018]

- Soft bellman equation

- $V^*(z_t) = \log \int \exp Q^*(z_t, a_t) da_t (= \log Z)$
- $Q^*(z_t, a_t) = r(z_t, a_t) + \mathbb{E}_{p(z_{t+1}|z_t, a_t)}[V^*(z_{t+1})]$
- $\pi^*(a_t|z_t) = \exp[Q^*(z_t, a_t) - V^*(z_t)] = \text{argmin } D_{KL}[\pi || \exp[Q^* - V^*]]$

- Concept of SAC

- Define deep parameterized models $V_\psi, Q_\theta, \pi_\phi$
- Train $V_\psi, Q_\theta, \pi_\phi$ to be V^*, Q^*, π^* by optimizing bellman errors

- $J_V = (V_\psi - \log \int \exp Q_\theta da)^2 \geq (V_\psi - \mathbb{E}[Q_\theta - \log \pi_\phi])^2$

- $J_Q = (Q_\theta - (r + \mathbb{E}[V_\psi]))^2$

- $J_\pi = -D_{KL}[\pi_\phi || \exp[Q_\theta - V_\psi]]$

Lower bound

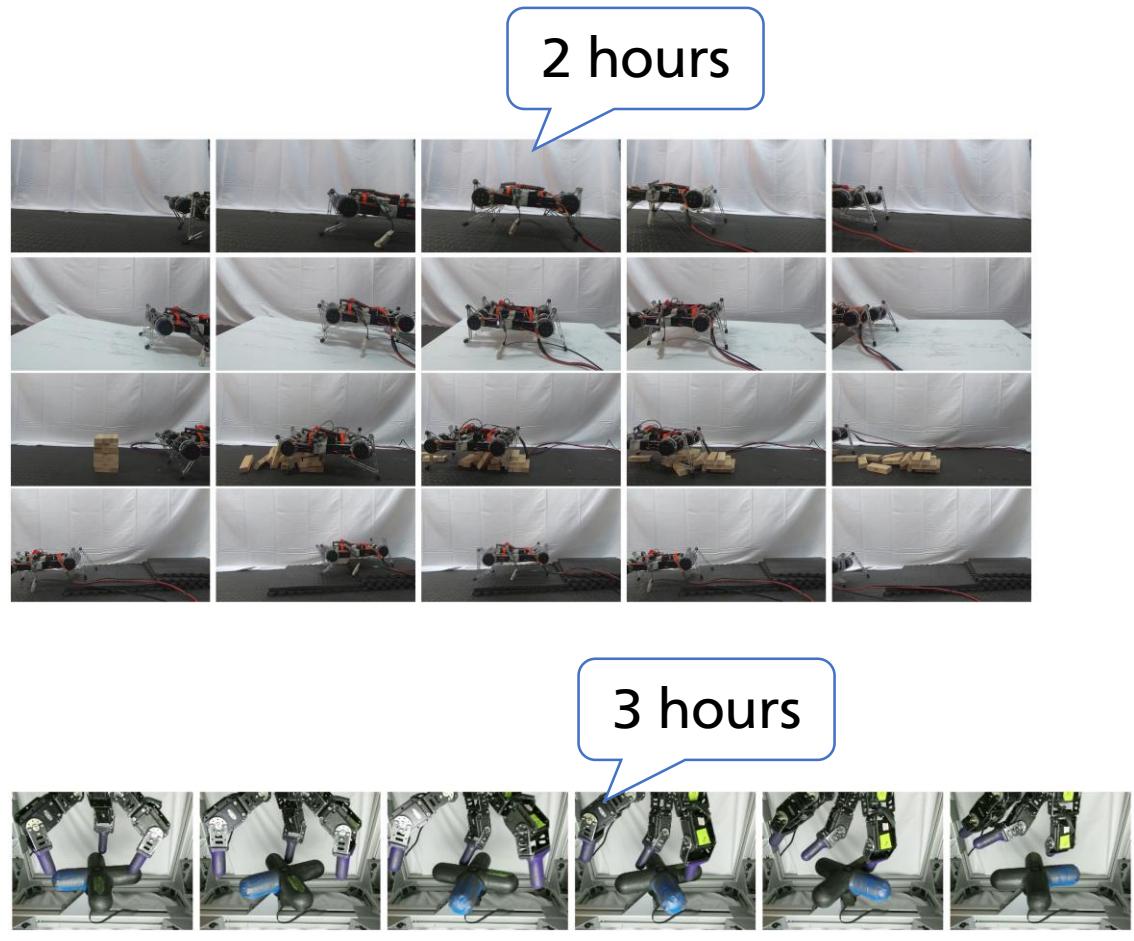
Soft Actor Critic (SAC)

[Haarnoja+, ICML2018]

Screenshot @ Oct. 5th, 2020

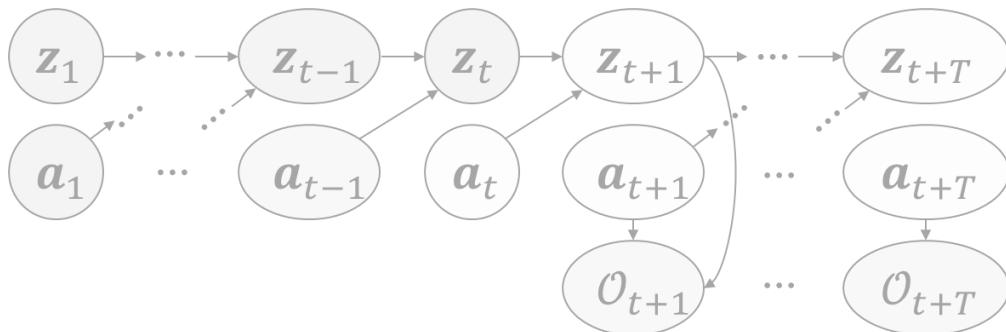
The screenshot shows a search results page for "soft actor critic" on Google Scholar. The search bar at the top contains the query. Below it, a "Scholar" section indicates about 88,000 results found in 0.05 seconds. A dropdown menu for "YEAR" is visible. The results are listed in a grid:

- Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor** [PDF] arxiv.org
T Haarnoja, A Zhou, P Abbeel, S Levine - arXiv preprint arXiv:1801.01290, 2018 - arxiv.org
Model-free deep reinforcement learning (RL) algorithms have been demonstrated on a range of challenging decision making and control tasks. However, these methods typically suffer from two major challenges: very high sample complexity and brittle convergence ...
☆ 99 Cited by 867 Related articles All 6 versions
- Soft actor-critic algorithms and applications** [PDF] arxiv.org
T Haarnoja, A Zhou, K Hartikainen, G Tucker... - arXiv preprint arXiv ..., 2018 - arxiv.org
Model-free deep reinforcement learning (RL) algorithms have been successfully applied to a range of challenging sequential decision making and control tasks. However, these methods typically suffer from two major challenges: high sample complexity and brittleness to ...
☆ 99 Cited by 227 Related articles All 2 versions
- Improving exploration in soft-actor-critic with normalizing flows policies** [PDF] arxiv.org
PN Ward, A Smofsky, AJ Bose - arXiv preprint arXiv:1906.02771, 2019 - arxiv.org
Deep Reinforcement Learning (DRL) algorithms for continuous action spaces are known to be brittle toward hyperparameters as well as cut {being} sample inefficient. **Soft Actor Critic** (SAC) proposes an off-policy deep **actor critic** algorithm within the maximum entropy RL ...
☆ 99 Cited by 9 Related articles All 2 versions
- Soft actor-critic for discrete action settings** [PDF] arxiv.org

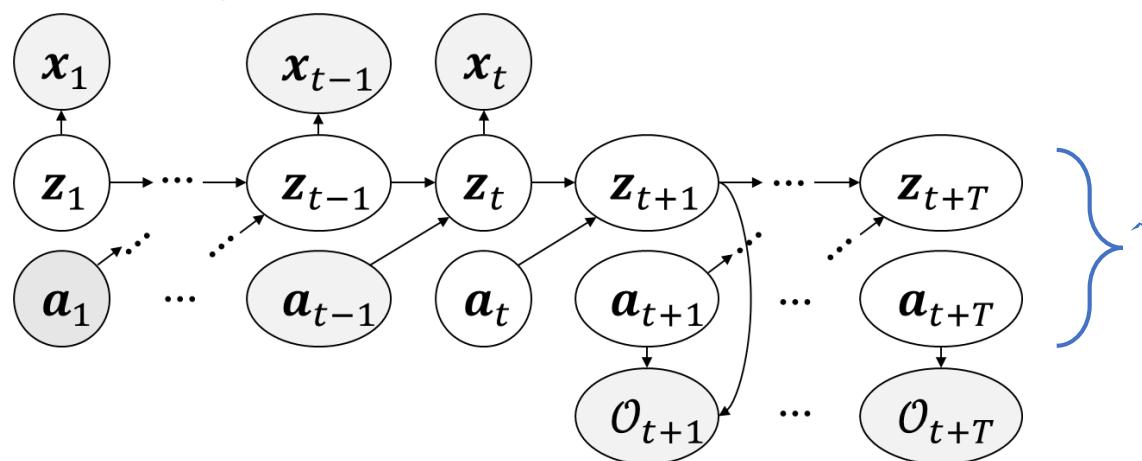


Problem Formulation

Fully Observable Markov Decision Process (MDP)



Partially Observable MDP



Trajectory: $\tau = (z_{\leq t+T}, a_{\geq t})$

What we want:
 $p(\tau | \mathcal{O}_{>t} = 1, x_{\leq t}, a_{<t})$

Trajectory optimization as inference

- Derivation of ELBO

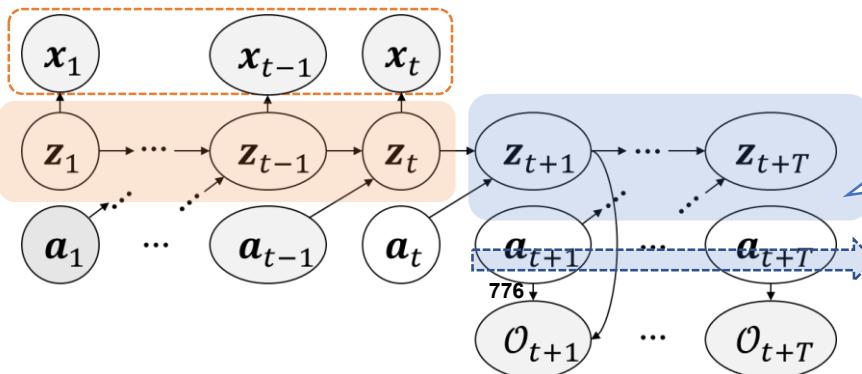
- $\log p(\mathcal{O}_{>t} = \mathbf{1}, x_{\leq t} | a_{<t})$
- $\geq \mathbb{E}_{q(\tau)}[\log p(\mathcal{O}_{>t} | \tau) - \log q(a_{\geq t}) + \log p(x_{\leq t} | z_{\leq t})]$

Observation likelihood

- Optimizing ELBO by latent VI-MPC [Okada+, IROS2020]

- $q(a_{>t}) \leftarrow \frac{q(a_{>t}) \cdot \mathbb{E}_{z_{\leq t+T} \sim \mathbb{P}}[p(\mathcal{O}_{>t} | \tau)] \cdot q^{-\kappa}}{\mathbb{E}_{q(a_{>t})}[\mathbb{E}_{z_{\leq t+T} \sim \mathbb{P}}[p(\mathcal{O}_{>t} | \tau)] \cdot q^{-\kappa}]}$
- $\mathbb{P} := (\prod_{t'=1}^t q(z_{t'} | x_{\leq t'}, a_{<t'})) (\prod_{t'=t}^{t+T} p(z_{t'+1} | z_{t'}, a_{t'}))$

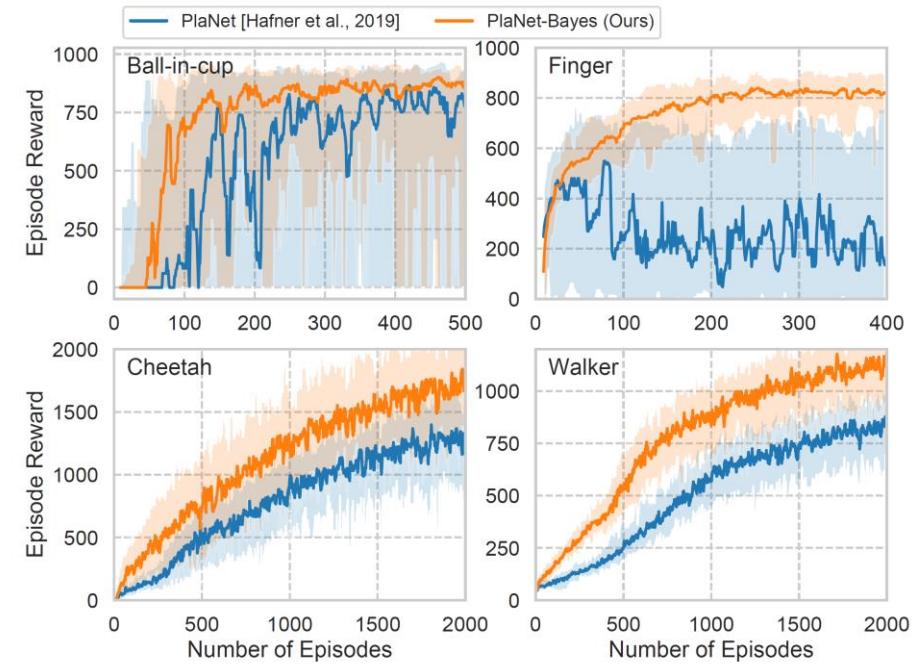
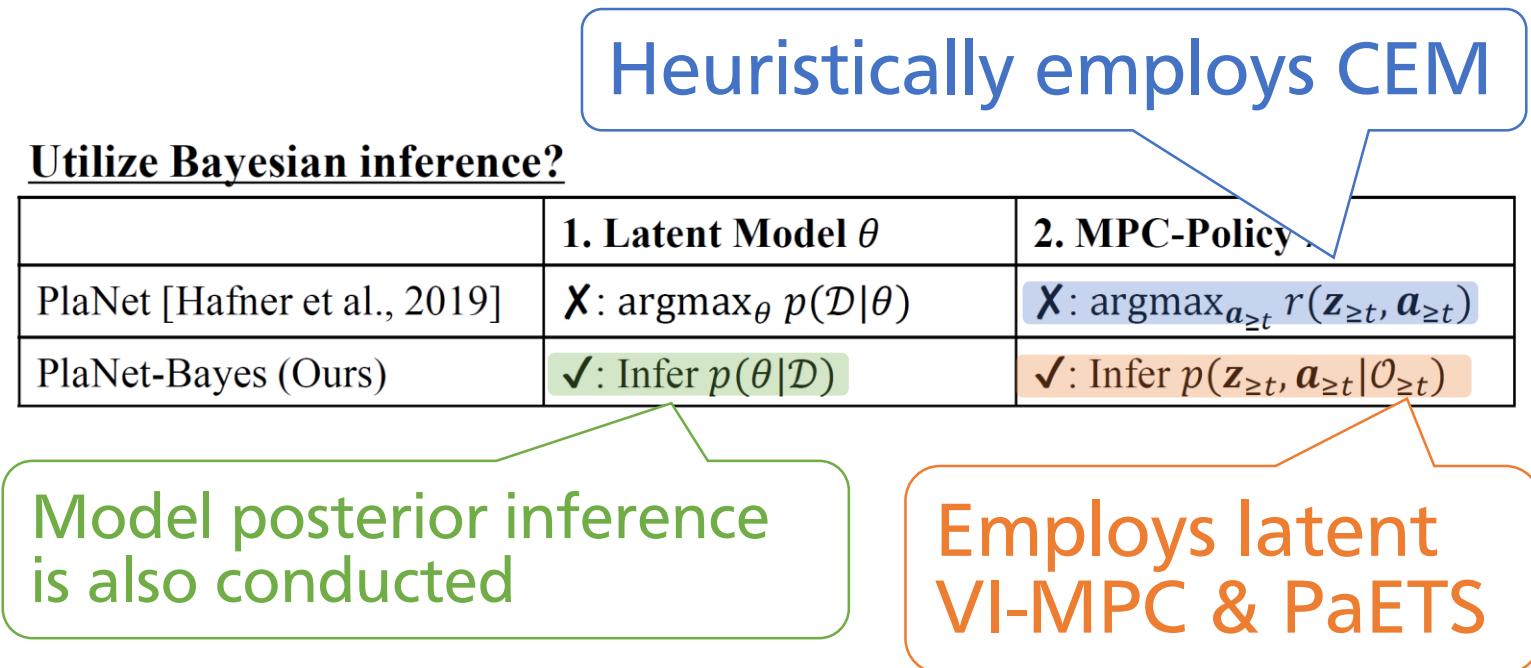
Inference by encoder
 $q(z_{t'} | \cdot)$



Prediction by latent dynamics $p(z_{t'+1} | \cdot)$

PlaNet of the Bayesians (PlaNet-Bayes) [Okada+, IROS2020]

- Extension of PlaNet based on (Bayesian) inference



Policy optimization as inference

SLAC (stochastic latent actor critic) [Lee+, arXiv2019]

- Derivation of ELBO
 - $\log p(\mathcal{O}_{>t} = \mathbf{1}, x_{\leq t} | a_{<t})$
 - $\geq \sum_{t' \leq t} \mathbb{E}_{q(z_{t'} | x_{\leq t'}, a_{<t'})} [\log p(x_{t'} | z_{t'})] - \sum_{t' > t} \mathbb{E}_{q(z_{t'-1} | \cdot)} \left[D_{KL}[q(z_{t'} | x_{\leq t'}, a_{<t'}) || p(z_{t'} | z_{t'-1}, a_{t'-1})] \right] + \sum_{t' > t} \mathbb{E}_{\pi(a_{t'} | x_{\leq t'}, a_{<t'})} [r(z_{t'}, a_{t'}) - \log \pi(a_{t'} | \cdot)]$
 - $\coloneqq \mathcal{J}_{VAE} + \mathcal{J}_{SAC}$
- Concept of SLAC
 - Define parameterized models $V_\psi, Q_\theta, \pi_\phi, q_\phi, p$
 - Jointly optimize \mathcal{J}_{VAE} and \mathcal{J}_{SAC} by similar training procedures of time-series VAE and SAC

Takeaways

- Various popular algorithms are based on variational inference (or ELBO maximization)
- Derive your novel algorithm by
 - Formulate your graphical model (or generative model)
 - Define variational distribution q
 - Derive J_{ELBO}
- And implement your algorithm by
 - TensorFlow, Keras, Pytorch,
 - [Pixyz, SERKET](#)

