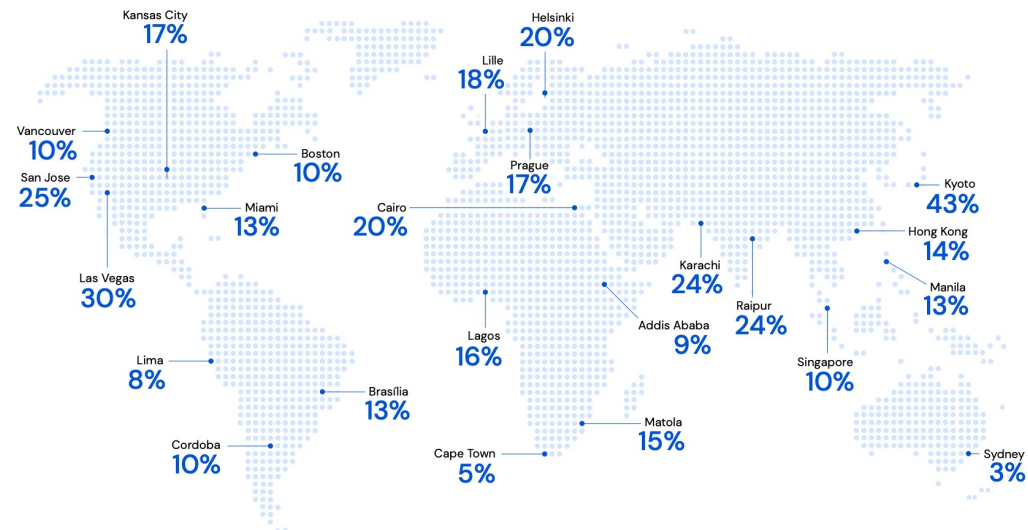




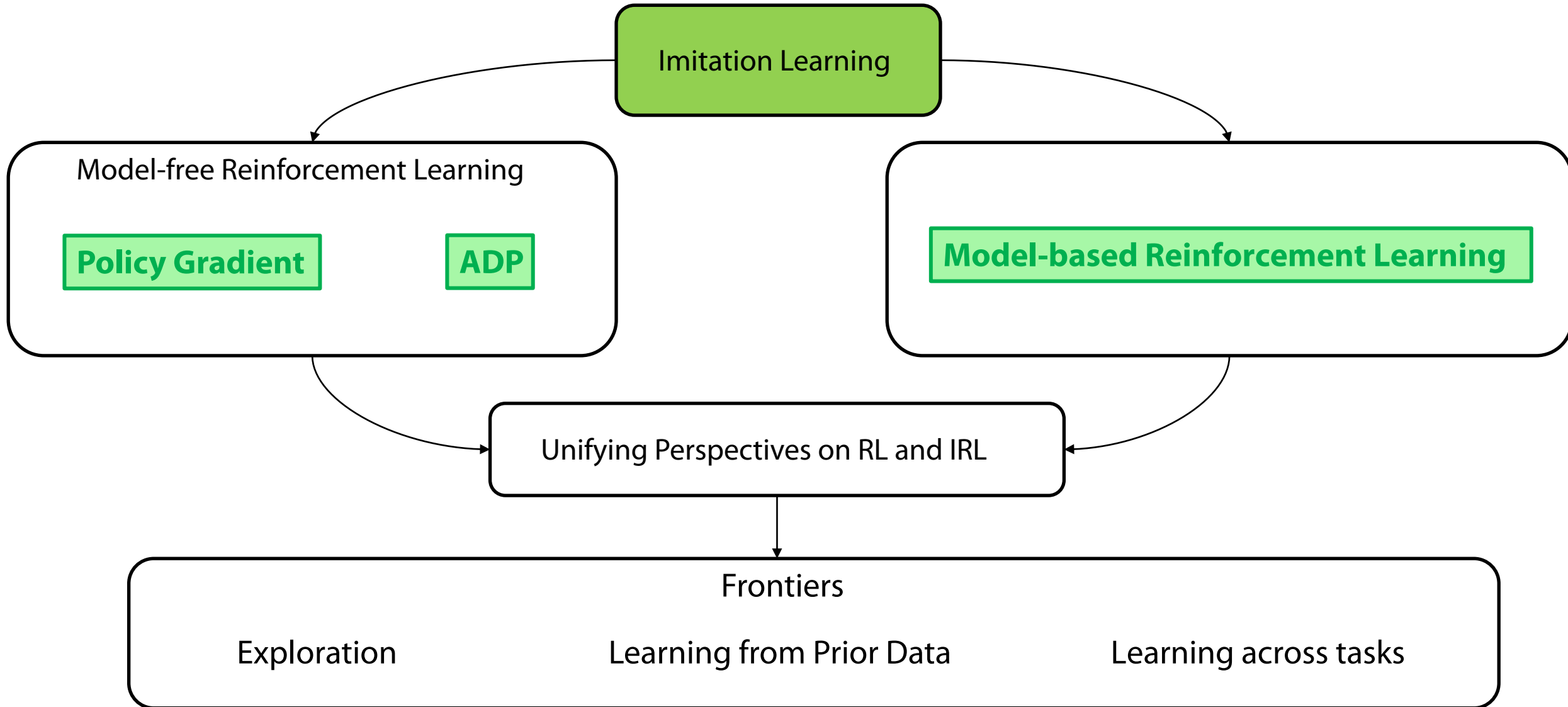
# Reinforcement Learning Spring 2026

Abhishek Gupta

TA: Mateo Guaman Castro



# Class Structure



# Ok, let's talk about "optimality"

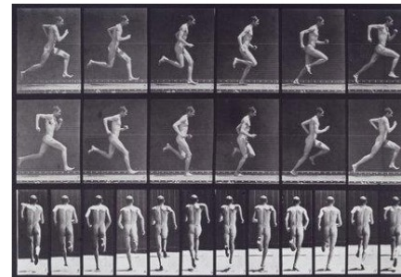
Optimal control problems aim to find the "max" reward policy

People are not perfectly rational, "noisily" rational

$$\arg \max_{a_0^j, a_1^j, \dots, a_T^j} \sum_{t=0}^T r(\hat{s}_t^j, a_t^j)$$
$$\hat{s}_{t+1}^j \sim \hat{p}_\theta(\cdot | \hat{s}_t^j, a_t^j)$$

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^T r(s_t, a_t) \right]$$

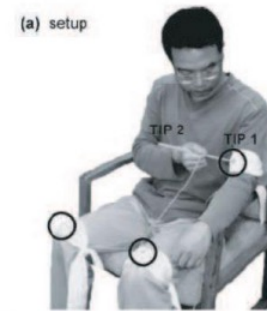
No notion of smooth suboptimality



Muybridge (c. 1870)



Mombaur et al. '09



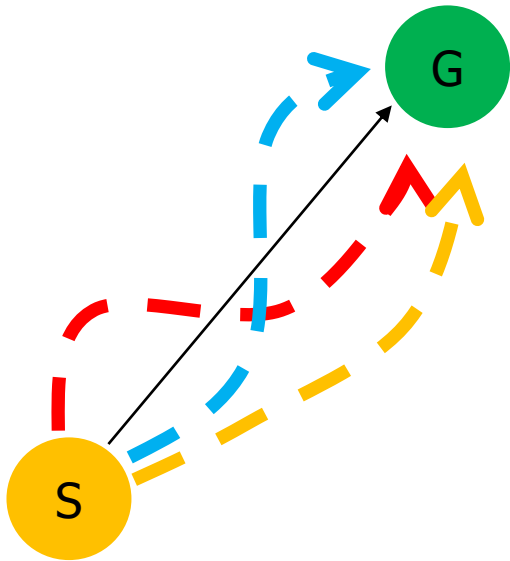
Li & Todorov '06



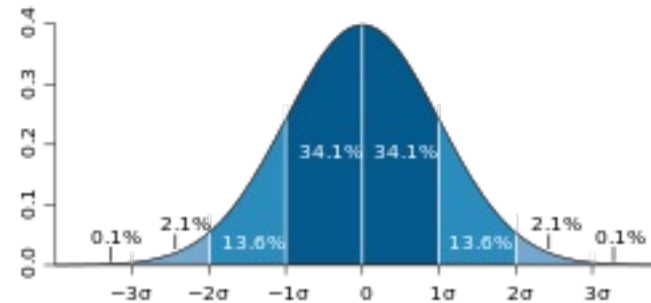
Ziebart '08

# Can we think about “soft optimality”?

So how can we properly model suboptimality?



Some mistakes are more important than others



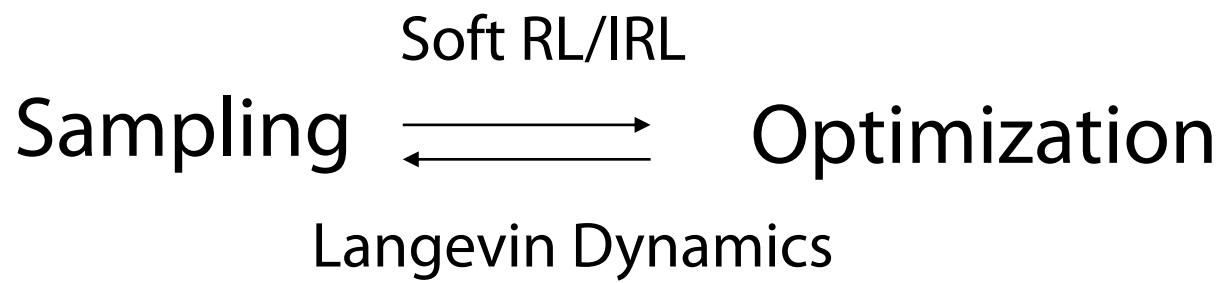
Let's use probability as a tool to represent “soft optimality”

- Going from deterministic to stochastic policies
- Better reward trajectories are “higher” likelihood
- Probabilistic measure of optimality, rather than an optimization one

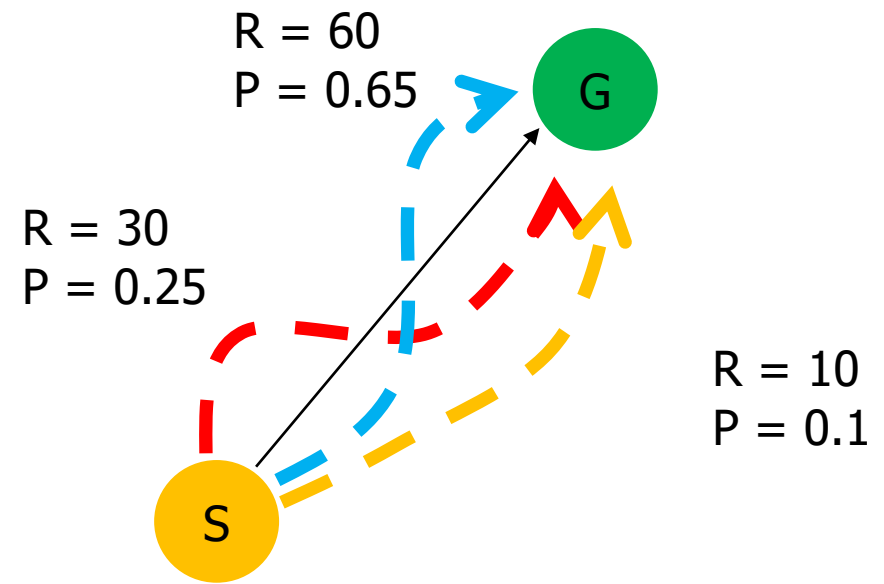
# Let's use probabilistic inference as a tool

$$\arg \max_{a_0^j, a_1^j, \dots, a_T^j} \sum_{t=0}^T r(\hat{s}_t^j, a_t^j)$$
$$\hat{s}_{t+1}^j \sim \hat{p}_\theta(\cdot | \hat{s}_t^j, a_t^j)$$

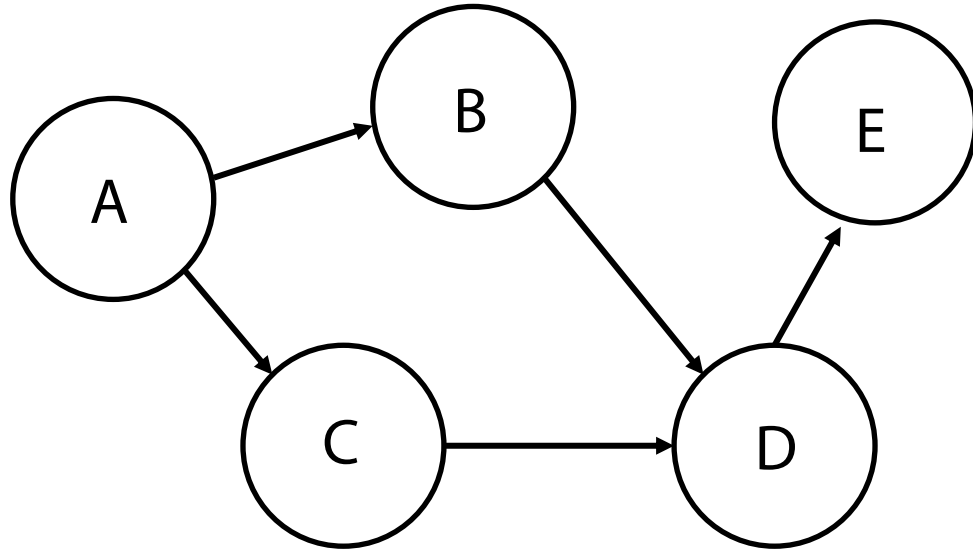
$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^T r(s_t, a_t) \right]$$



Rather than taking max wrt returns, sample proportional to returns



# Background: Probabilistic Graphical Models



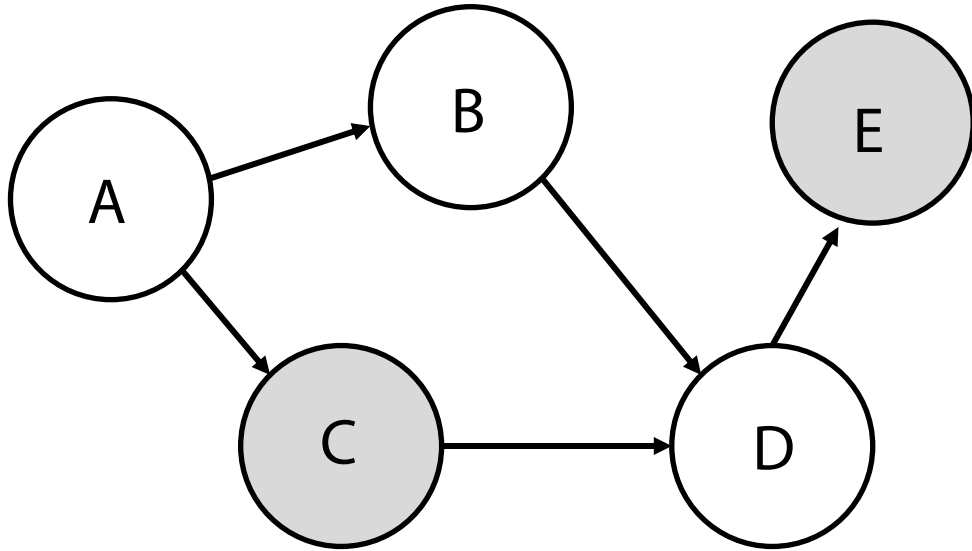
Convenient way to encode  
joint probability distribution

Encodes probabilities and conditional independences

$$P(A, B, \dots) = \prod_X P(X | \text{Parents}(X))$$

$$P(A, B, \dots) = P(A)P(B|A)P(C|A)P(D|B, C)P(E|D)$$

# Probabilistic Graphical Models



So what can you do with a probabilistic graphical model?

$$P(B|C, E)$$

Answer posterior inference queries

$$P(A, B|C, E)$$

What does this have to do with RL?

Isn't RL about maximizing expected reward?

Need to "eliminate" variables and use Bayes rule  
→ Easy in discrete space, challenging in continuous

# Lecture outline

---

Control as Inference - Formulation



Variational Inference



Control as Inference to Derive Policy Gradient

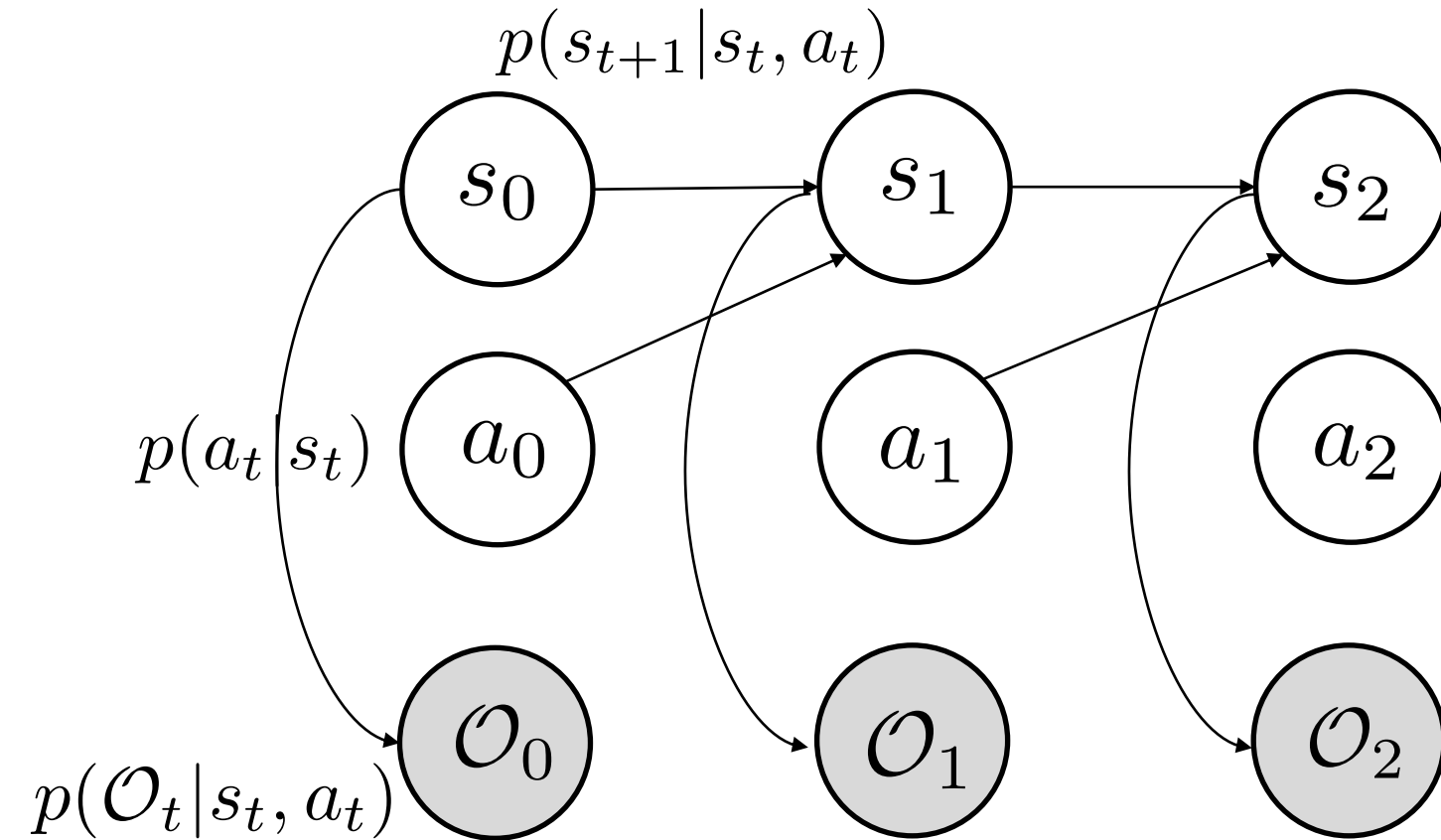


Control as Inference to Derive Q-learning



Control as Inference to Derive Model-Based RL

# Using Probabilistic Graphical Models for Decision Making



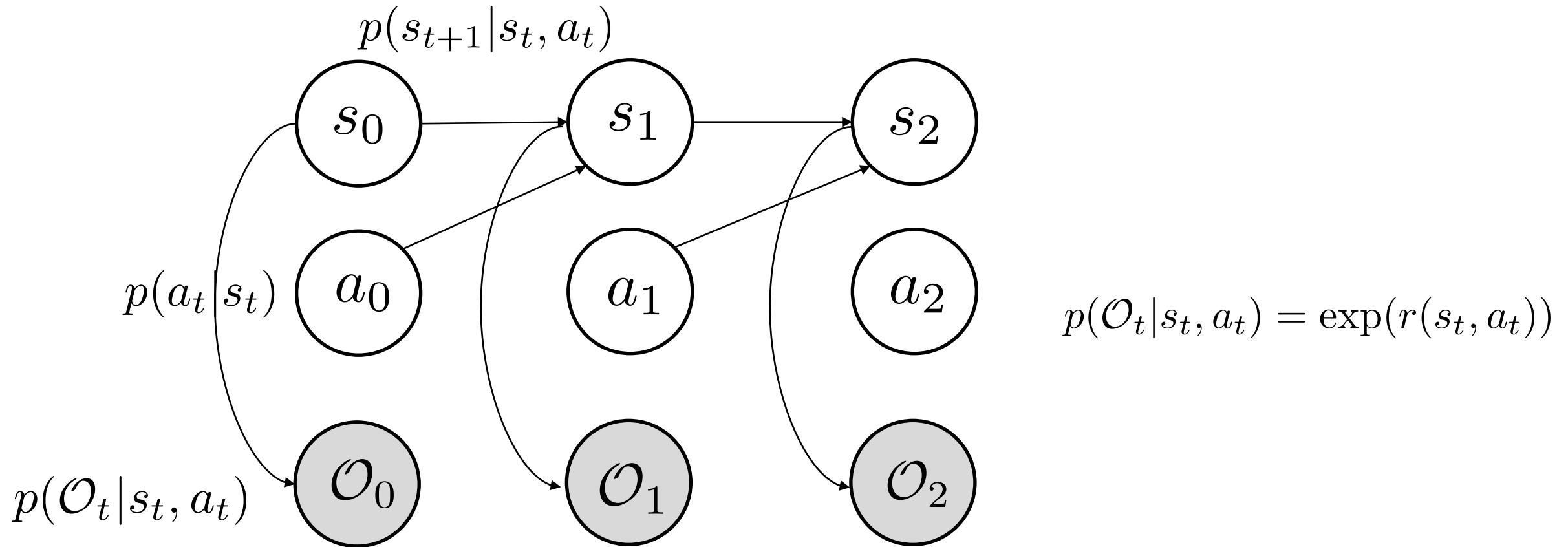
$$p(O_t|s_t, a_t) = \exp(r(s_t, a_t))$$

Rewards must be negative  
(subtract max reward WLOG)

Introduce binary “optimality” variables – optimal if  $O=1$ , suboptimal if  $O=0$

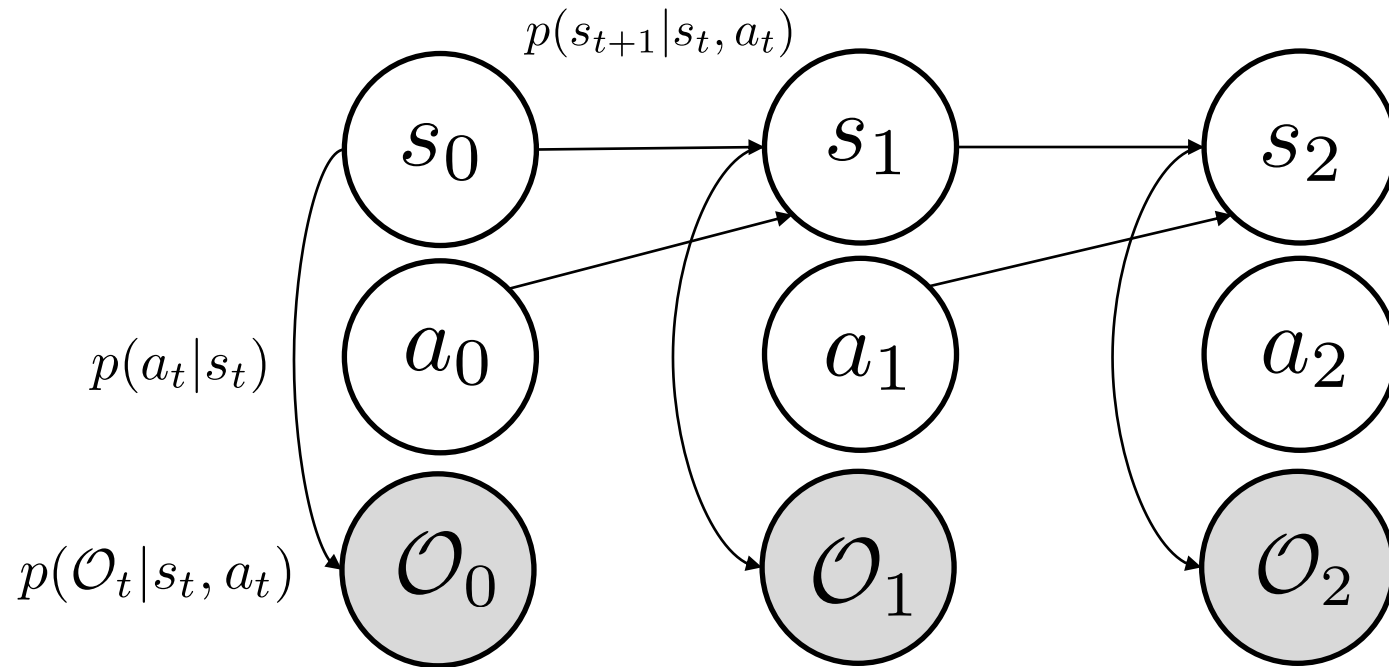
Agents are observed to be **optimal**

# Ok so how can we cast decision making as a PGM?



$$\begin{aligned}
 p(\tau | \mathcal{O}_{0:T} = 1) &\propto p(\tau) p(\mathcal{O}_{0:T} | \tau) = p(s_0) \prod_{t=0}^{T-1} p(s_{t+1} | s_t, a_t) p(a_t | s_t) p(\mathcal{O}_t | s_t, a_t) \\
 &= p(\tau) \exp\left(\sum_{t=0}^{T-1} r(s_t, a_t)\right) \quad \text{"Soft" optimality - higher return trajectories are higher likelihood}
 \end{aligned}$$

# Ok big whoop, what do we do this?



$$p(O_t|s_t, a_t) = \exp(r(s_t, a_t))$$

$$p(\tau|\mathcal{O}_{0:T} = 1) \propto p(\tau) \exp\left(\sum_{t=0}^T r(s_t, a_t)\right)$$

Use case 1:

Derive soft RL algorithms

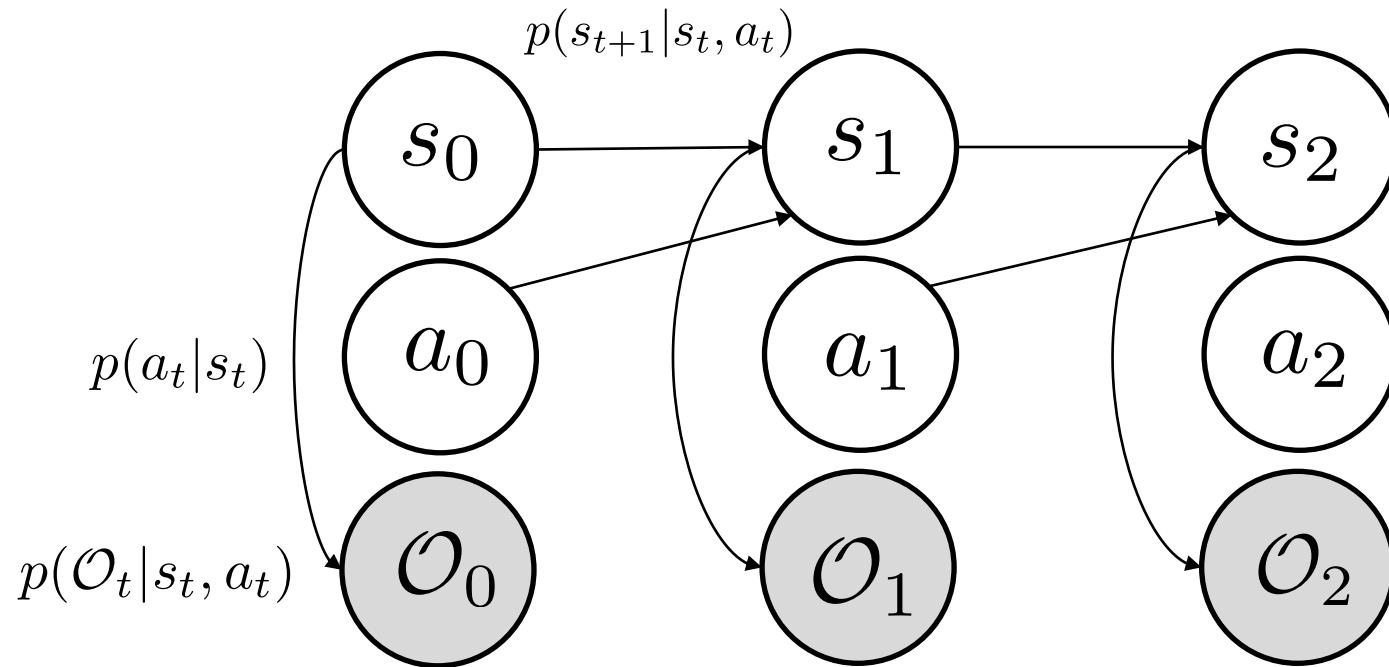
Use case 2:

Derive soft inverse RL algorithms

Use case 3:

Great algorithms for transfer

# So what are we doing inference over?



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

$$p(\tau|\mathcal{O}_{0:T} = 1) \propto p(\tau) \exp\left(\sum_{t=0}^T r(s_t, a_t)\right)$$

Use case 1:

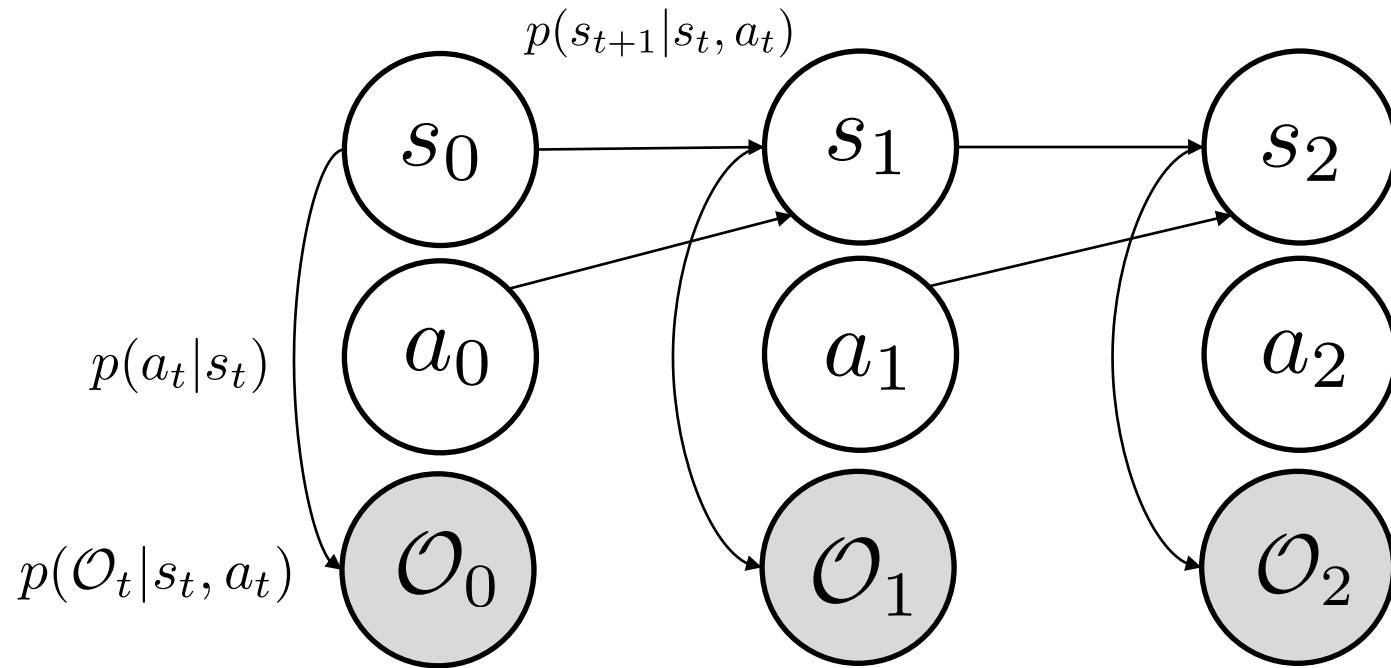
Derive soft RL algorithms

Insight: Computing optimal policy  $\rightarrow$  posterior inference

$$p(a_t|s_t, \mathcal{O}_{t:T} = 1)$$

“Given that you are acting optimally, what is the likelihood of a particular action at a state”

# Why isn't this trivial?



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

$$p(\tau | \mathcal{O}_{0:T} = 1) \propto p(\tau) \exp\left(\sum_{t=0}^T r(s_t, a_t)\right)$$

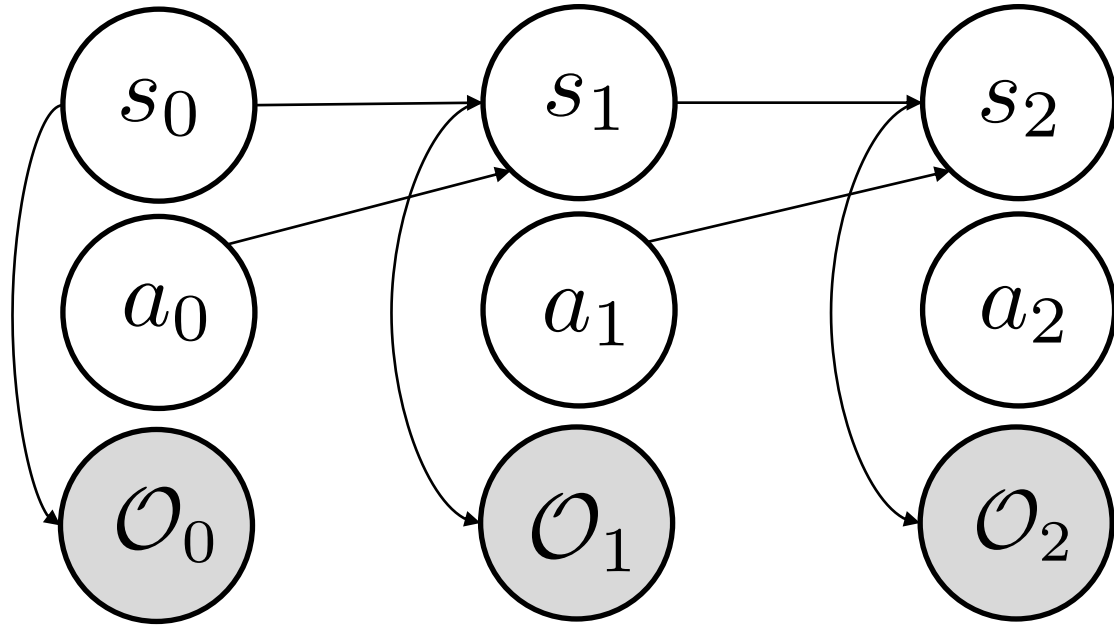
Optimal Policy  $\rightarrow$  Posterior Inference

$$p(a_t | s_t, \mathcal{O}_{t:T} = 1) = \frac{p(a_t, \mathcal{O}_{t:T} = 1 | s_t)}{p(\mathcal{O}_{t:T} = 1 | s_t)} = \frac{\int \int \cdots \int p(a_{t:T}, \mathcal{O}_{t:T} = 1, s_{t:T}) ds_{t+1:T} da_{t+1:T}}{\int \int \cdots \int p(a_{t:T}, \mathcal{O}_{t:T} = 1, s_{t:T}) ds_{t+1:T} da_{t:T}}$$

“Given that you are acting optimally, what is the likelihood of a particular action at a state”

Difficult/intractable to compute  
 $\rightarrow$  Most RL algorithms are approximations to this

# What makes this so cool?



Policy Gradient

Approximate DP

Model-Based RL

Variational Inference lower bound  
solved with Gradient Ascent

Variational Inference lower bound  
solved with dynamic programming

Posterior Inference Approximated with  
Monte-Carlo Samples

Optimal Policy  $\rightarrow$  Posterior Inference

$$\begin{aligned} & p(a_t | s_t, \mathcal{O}_{t:T} = 1) \\ &= \frac{p(a_t, \mathcal{O}_{t:T} = 1 | s_t)}{p(\mathcal{O}_{t:T} = 1 | s_t)} \\ &= \frac{\int \int \cdots \int p(a_{t:T}, \mathcal{O}_{t:T} = 1, s_{t:T}) ds_{t+1:T} da_{t+1:T}}{\int \int \cdots \int p(a_{t:T}, \mathcal{O}_{t:T} = 1, s_{t:T}) ds_{t+1:T} da_{t:T}} \end{aligned}$$

Can derive old algorithms + new classes of algorithms from the same framework!

# Lecture outline

---

Control as Inference - Formulation



Variational Inference



Control as Inference to Derive Policy Gradient

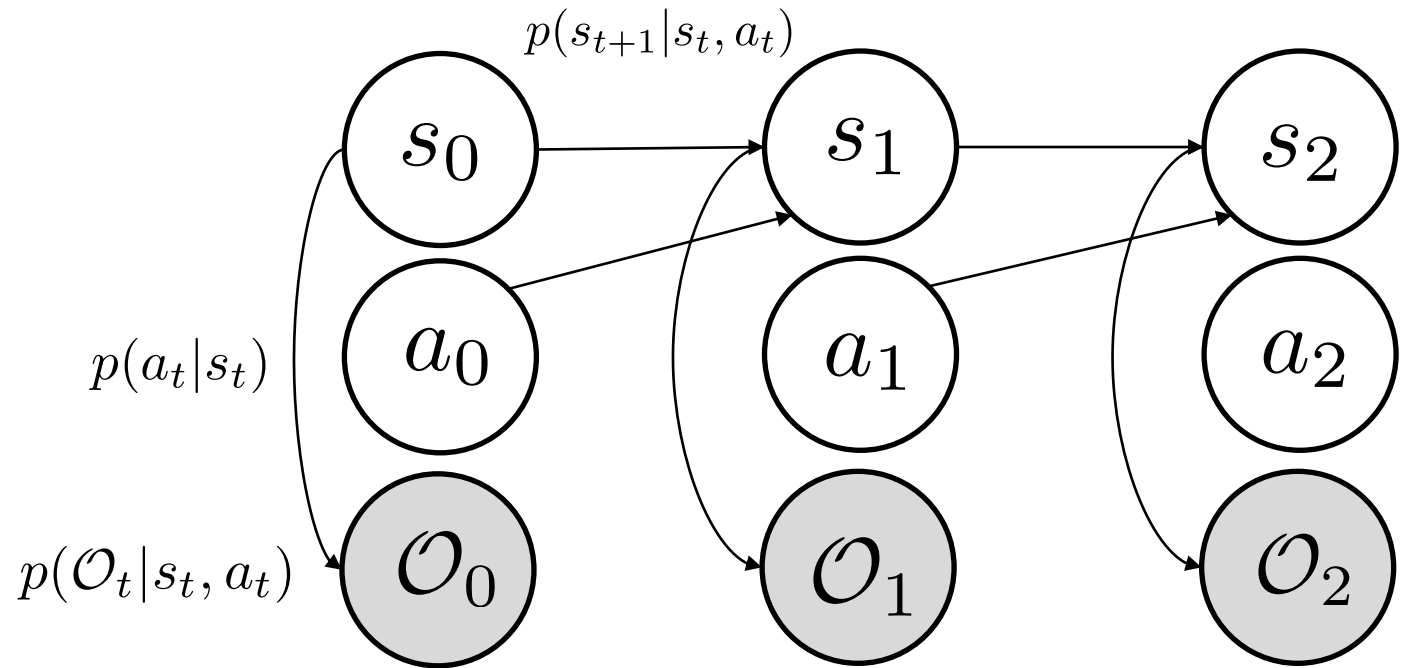


Control as Inference to Derive Q-learning



Control as Inference to Derive Model-Based RL

# Why isn't this trivial?



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

$$p(\tau | \mathcal{O}_{0:T} = 1) \propto p(\tau) \exp\left(\sum_{t=0}^T r(s_t, a_t)\right)$$

Optimal Policy  $\rightarrow$  Posterior Inference

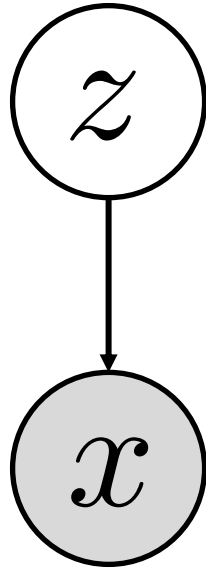
$$p(a_t | s_t, \mathcal{O}_{t:T} = 1) = \frac{p(a_t, \mathcal{O}_{t:T} = 1 | s_t)}{p(\mathcal{O}_{t:T} = 1 | s_t)} = \frac{\int \int \cdots \int p(a_{t:T}, \mathcal{O}_{t:T} = 1, s_{t:T}) ds_{t+1:T} da_{t+1:T}}{\int \int \cdots \int p(a_{t:T}, \mathcal{O}_{t:T} = 1, s_{t:T}) ds_{t+1:T} da_{t:T}}$$

“Given that you are acting optimally, what is the likelihood of a particular action at a state”

Difficult/intractable to compute  
 $\rightarrow$  Most RL algorithms are approximations to this

# Let's take the simplest possible example

Let us assume  $p(x|z)$  is known, as is  $p(z)$



Standard latent-variable model

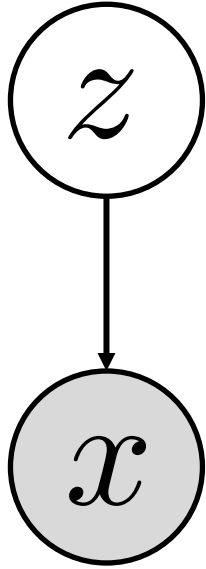
Goal: Infer posterior  $p(z|x)$

$$p(z|x) = \frac{p(x, z)}{p(x)} = \frac{p(x|z)p(z)}{p(x)}$$

$$= \frac{p(x|z)p(z)}{\int p(x|z)p(z)dz}$$

Challenging to compute efficiently with samples

# So how can we solve this posterior inference problem?



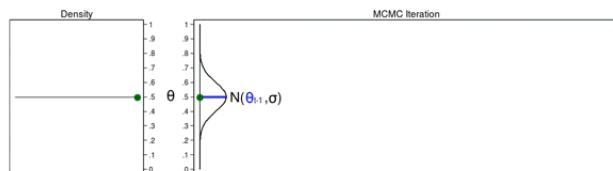
Let us assume  $p(x|z)$  is known, as is  $p(z)$

Goal: Infer posterior  $p(z|x)$

$$p(z|x) = \frac{p(x|z)p(z)}{\int p(x|z)p(z)dz}$$

Challenging to compute efficiently with samples

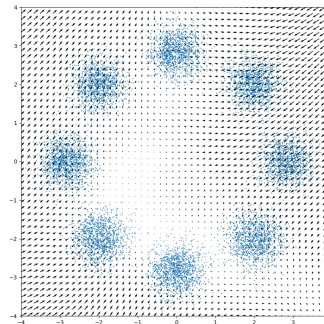
## MCMC



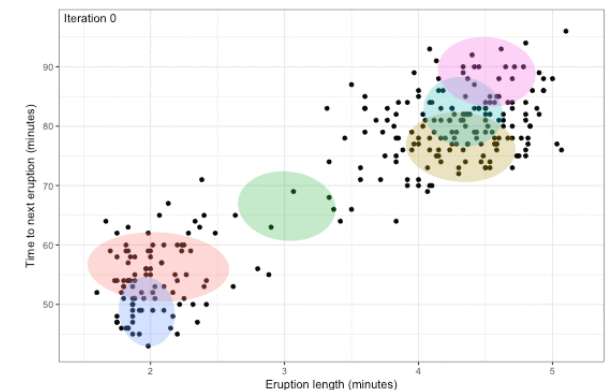
Draw  $\theta_t \sim \text{Normal}(\theta_{t-1}, \sigma)$

Normal(0.500,  $\sigma$ ) = 0.497

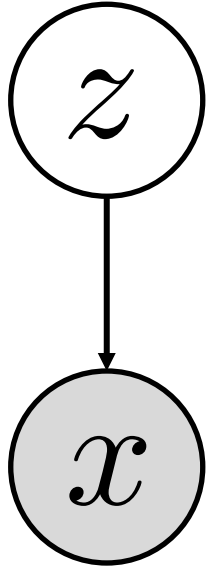
## EBMs and Score Matching



## Variational Inference

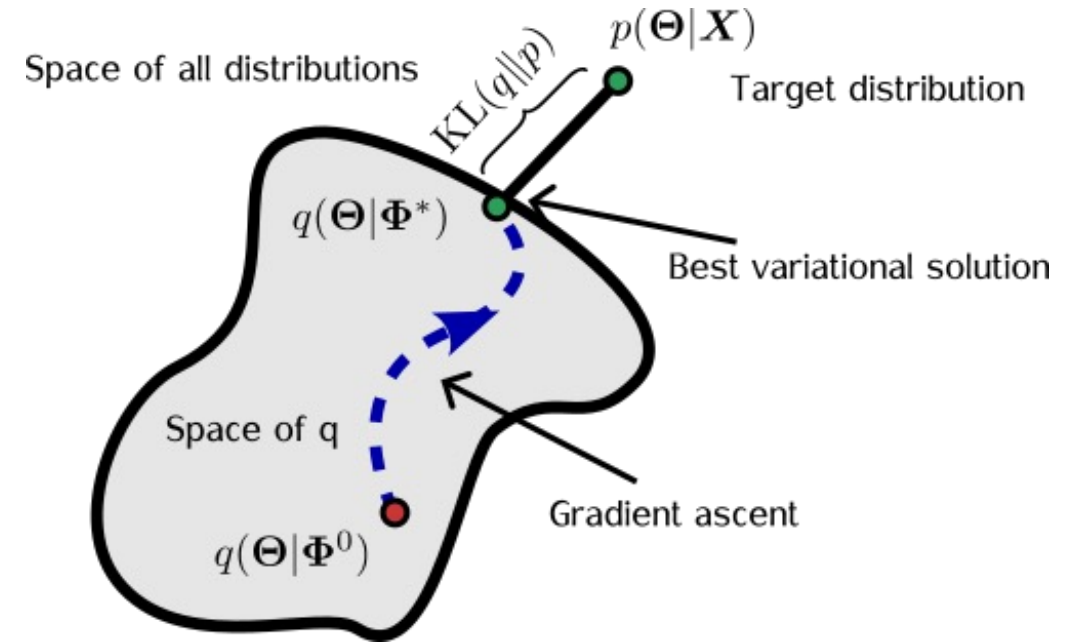


# What is the key idea behind variational inference?



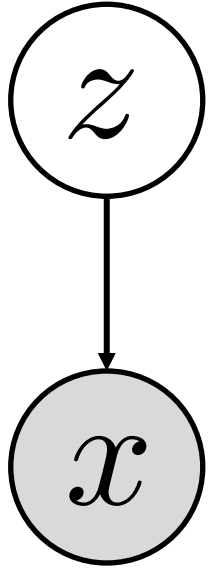
$$p(z|x) = \frac{p(x|z)p(z)}{\int p(x|z)p(z)dz}$$

Intractable!



Approximate challenging posterior with closest possible "tractable" posterior

# Let's derive the Evidence Lower Bound



$$p(z|x) = \frac{p(x|z)p(z)}{\int p(x|z)p(z)dz}$$

Intractable!

Introduce a "tractable" approximation  $q(z|x)$   
e.g. Gaussian

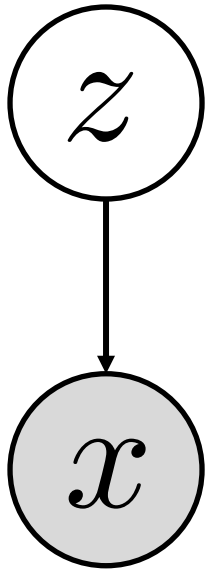
Can choose **whatever** variational family you want  
→ it's an approximation! 🙄

$$\phi^* \leftarrow \arg \min_{\phi} D_{KL}(q_{\phi}(z|x) || p(z|x)) \quad \text{Unknown}$$

Known

How can we tractably approximate this objective?

# Let's derive the Evidence Lower Bound



$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

Intractable!

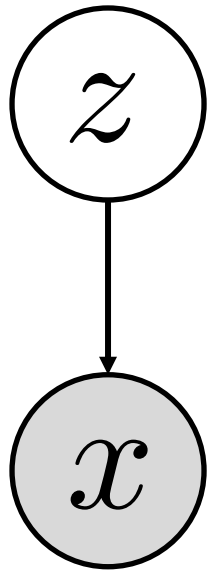
Unknown

$$\phi^* \leftarrow \arg \min_{\phi} D_{KL}(q_{\phi}(z|x) || p(z|x))$$

Known

$$\begin{aligned} D_{KL}(q_{\phi}(z|x) || p(z|x)) &= \int q(z|x) \log \frac{q(z|x)}{p(z|x)} dz = \int q(z|x) \log \frac{q(z|x)p(x)}{p(x|z)p(z)} dz \\ &= \int q(z|x) \log \frac{q(z|x)}{p(z)} dz - \int q(z|x) \log p(x|z) dz + \log p(x) \\ &= D_{KL}(q(z|x) || p(z)) - \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] + \log p(x) \end{aligned}$$

# Let's derive the Evidence Lower Bound



$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

Intractable!

$$\phi^* \leftarrow \arg \min_{\phi} D_{KL}(q_{\phi}(z|x) || p(z|x))$$

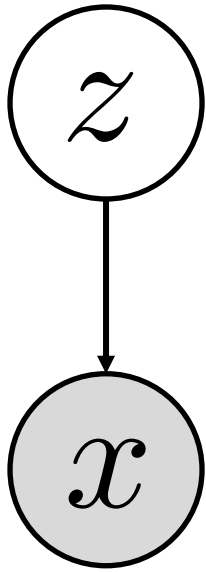
Unknown  
Known

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

View 1: Find best posterior

View 2: Maximize marginal likelihood

# Evidence Lower Bound: Best Posterior



$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

Intractable!

View 1: Find best posterior

$$D_{KL}(q_\phi(z|x) || p(z|x))$$

$$= D_{KL}(q(z|x) || p(z)) - \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] + \log p(x)$$

Likelihood/prior known – posterior hard to compute

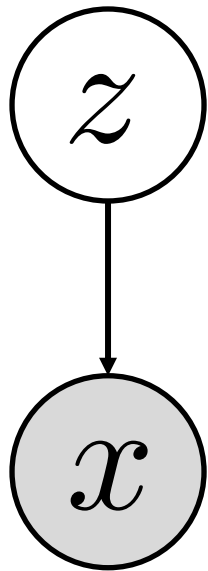
Maximum likelihood

Stay close to the prior

$$\max_q \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z)) \right]$$

Learn a tractable posterior  $q(z|x)$  with known likelihood and sampling

# Evidence Lower Bound: Max Marginal Likelihood



$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

Intractable!

View 2: Maximize marginal likelihood

$$\begin{aligned} & D_{KL}(q_\phi(z|x) || p(z|x)) \\ &= D_{KL}(q(z|x) || p(z)) - \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] + \log p(x) \end{aligned}$$

Likelihood unknown and posterior hard to compute

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

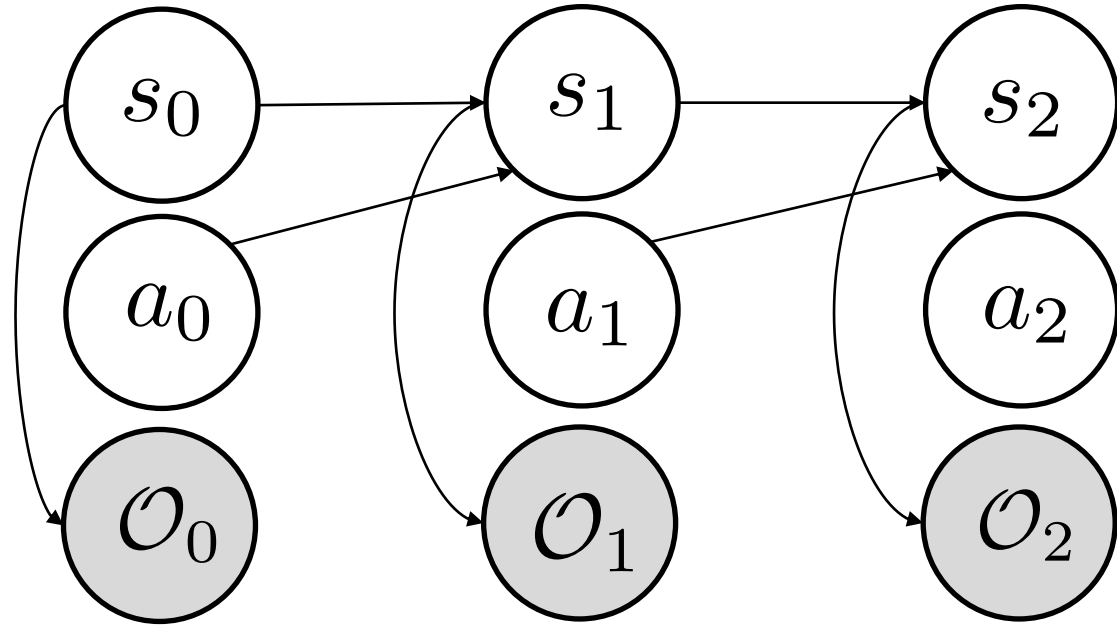
$$D_{KL}(p || q) \geq 0$$

$$\log p(x) \geq \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z))$$

Evidence **lower** bound – maximize to maximize likelihood

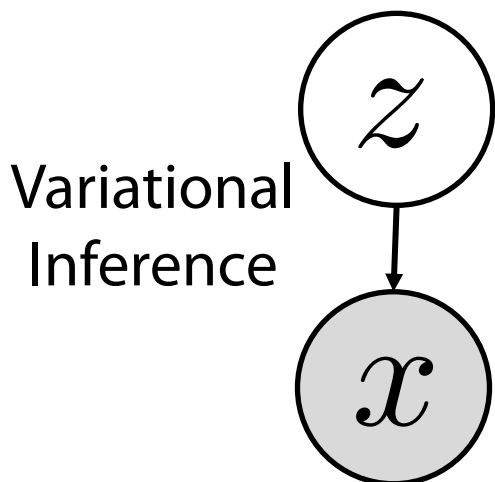
Learned

# Lets revisit our original inference problem in control



Optimal Policy  $\rightarrow$  Posterior Inference

$$\begin{aligned}
 & p(a_t | s_t, \mathcal{O}_{t:T} = 1) \\
 &= \frac{p(a_t, \mathcal{O}_{t:T} = 1 | s_t)}{p(\mathcal{O}_{t:T} = 1 | s_t)} \\
 &= \frac{\int \int \cdots \int p(a_{t:T}, \mathcal{O}_{t:T} = 1, s_{t:T}) ds_{t+1:T} da_{t+1:T}}{\int \int \cdots \int p(a_{t:T}, \mathcal{O}_{t:T} = 1, s_{t:T}) ds_{t+1:T} da_{t:T}}
 \end{aligned}$$



Approximate  $p(a_t | s_t, \mathcal{O}_{t:T} = 1)$  by  $q(a_t | s_t, \mathcal{O}_{t:T} = 1)$

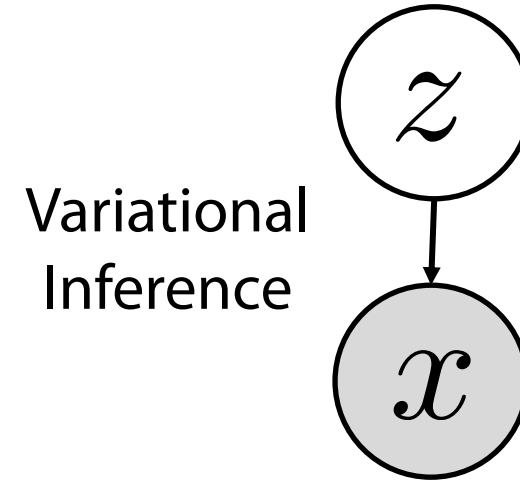
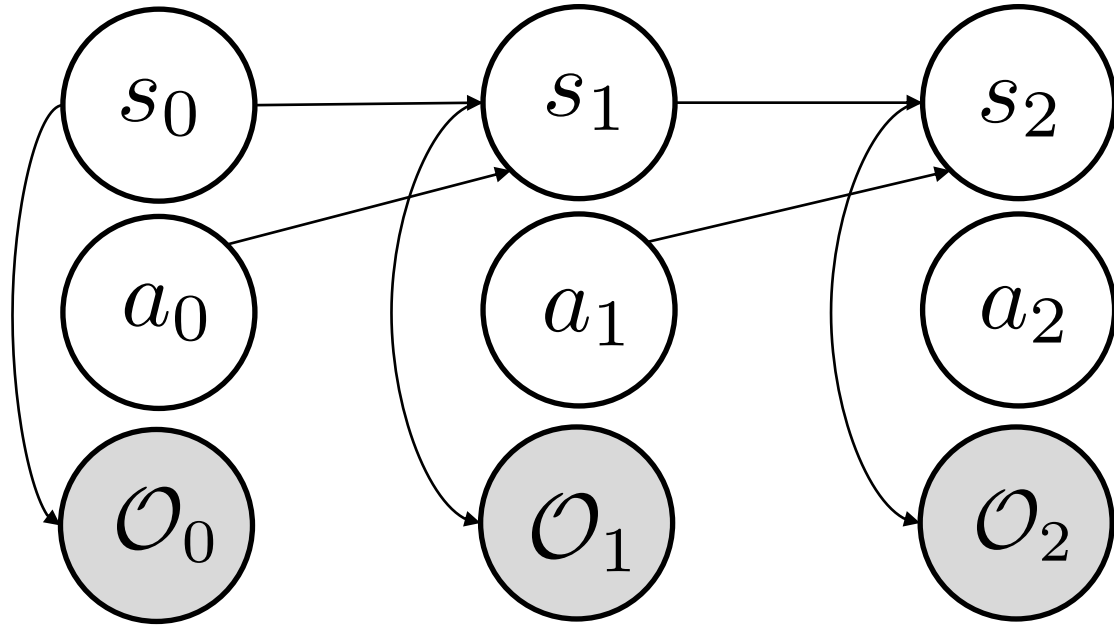
$$\max_q \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z)) \right]$$

$x$   
 $\updownarrow$   
 $(\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_T)$

Tractable techniques for  
 posterior policy computation

$z$   
 $\updownarrow$   
 $(s_0, a_0, s_1, a_1, \dots, s_T, a_T)$

# Lets revisit our original inference problem in control



$$\max_q \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z)) \right]$$

$x$



$(\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_T)$

$z$



$(s_0, a_0, s_1, a_1, \dots, s_T, a_T)$

Next –

derive ELBO and work out how to compute  
Policy gradient/Actor-Critic

# Lecture outline

---

Control as Inference - Formulation



Variational Inference



Control as Inference to Derive Policy Gradient

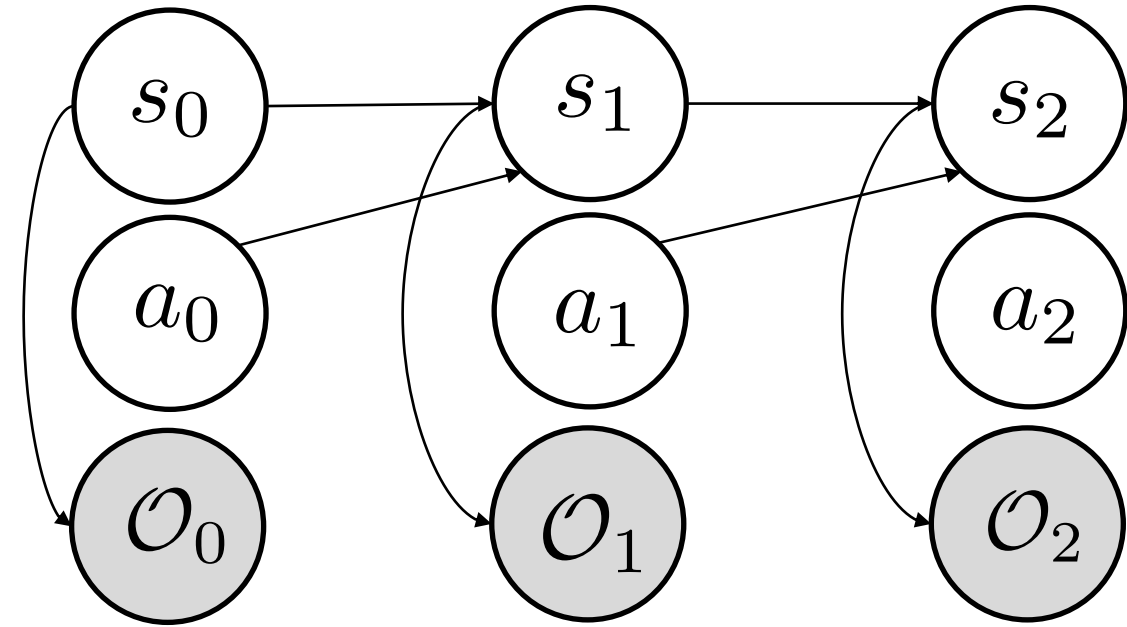


Control as Inference to Derive Q-learning



Control as Inference to Derive Model-Based RL

# Lets revisit our original inference problem in control



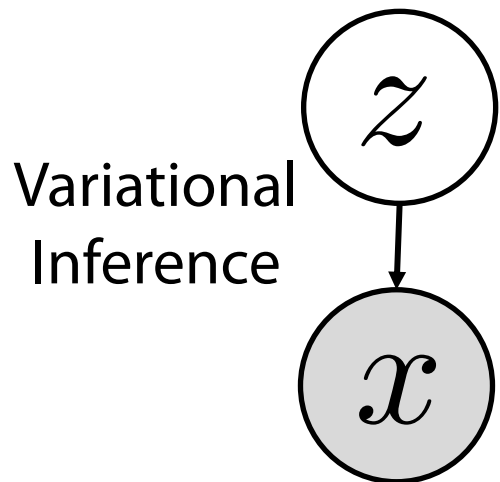
Optimal Policy  $\rightarrow$  Posterior Inference

Approximate  $p(a_t | s_t, \mathcal{O}_{t:T} = 1)$  by  $q(a_t | s_t, \mathcal{O}_{t:T} = 1)$

$$\max_q \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z)) \right]$$

$x$   
 $\updownarrow$   
 $(\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_T)$

$z$   
 $\updownarrow$   
 $(s_0, a_0, s_1, a_1, \dots, s_T, a_T)$



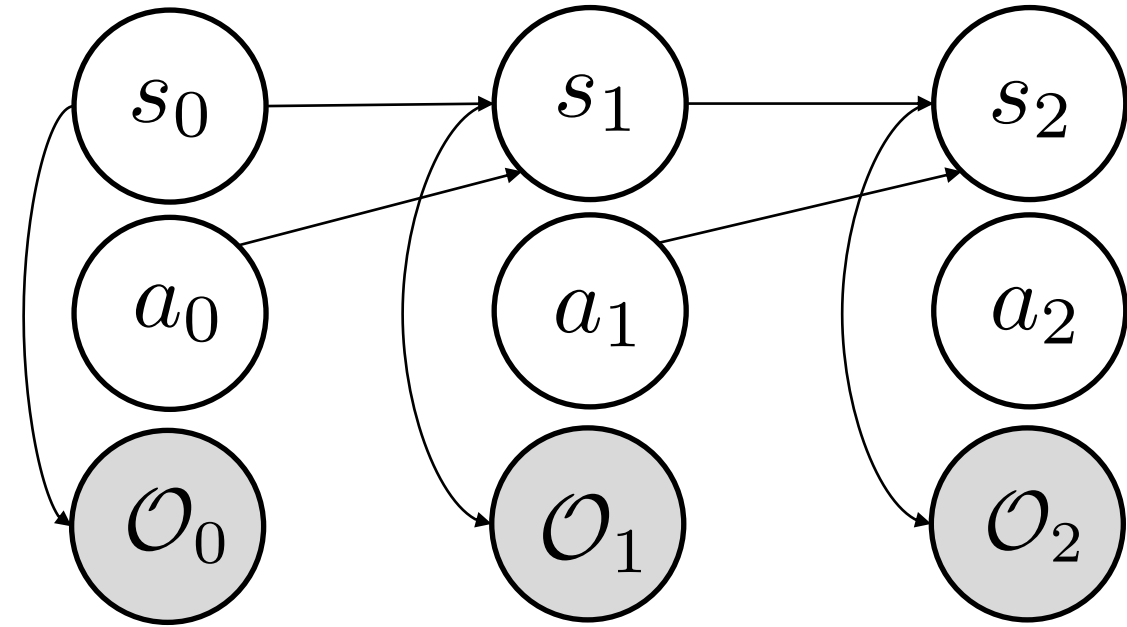
So what do we need to compute this?

$\rightarrow$  Choice for  $q(z|x)$  --  $q(s_0, a_0, \dots, s_T, a_T | \mathcal{O}_0, \dots, \mathcal{O}_T)$

Key desiderata:

- 1) Can sample,
- 2) Compute KL Divergence

# Choice of Variational Family for Approximate Inference



$$\max_q \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z)) \right]$$

$x$   $z$   
 $\updownarrow$   $\updownarrow$   
 $(\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_T)$   $(s_0, a_0, s_1, a_1, \dots, s_T, a_T)$

So what do we need to compute this?

→ Choice for  $q(s_0, a_0, \dots, s_T, a_T | \mathcal{O}_0, \dots, \mathcal{O}_T)$

Use temporal structure of the **true** dynamics in  $q$

$$q(s_0, a_0, \dots, s_T, a_T | \mathcal{O}_0, \dots, \mathcal{O}_T) = p(s_0) p(s_1 | s_0, a_0) q(a_0 | s_0) p(s_2 | s_1, a_1) q(a_1 | s_1) \dots q(a_T | s_T)$$

Does it satisfy:

- 1) Can sample,
- 2) Compute KL Divergence

$$= p(s_0) \prod_{t=0}^T p(s_{t+1} | s_t, a_t) q(a_t | s_t)$$

True dynamics and initial state

Approximate policy for

$$p(a_t | s_t, \mathcal{O}_{t:T} = 1)$$



# Computing Evidence Lower Bound

$$x (\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_T) \quad z (s_0, a_0, s_1, a_1, \dots, s_T, a_T)$$

$$\max_q \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z)) \right]$$

$$q(s_0, a_0, \dots, s_T, a_T | \mathcal{O}_0, \dots, \mathcal{O}_T) = p(s_0) \prod_{t=0}^T p(s_{t+1} | s_t, a_t) q(a_t | s_t)$$

1)  $\mathbb{E}_{x \sim p(x), z \sim q(z|x)} [\log p(x, z) - \log q(z|x)]$

2)  $\mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \log p(s_0, \dots, s_T, a_0, \dots, a_T, \mathcal{O}_0, \dots, \mathcal{O}_T) - \log q(s_0, a_0, \dots, s_T, a_T | \mathcal{O}_0, \dots, \mathcal{O}_T) \right]$

3)  $\mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \left[ \log p(s_0) + \sum_t [\log p(a_t | s_t) + \log p(s_{t+1} | s_t, a_t) + \log p(\mathcal{O}_t | s_t, a_t)] \right] - \left[ \log p(s_0) + \sum_t [\log q(a_t | s_t) + \log p(s_{t+1} | s_t, a_t)] \right] \right]$

# Computing Evidence Lower Bound

$$\max_q \mathbb{E}_{x \sim p(x)} [\mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z))]$$

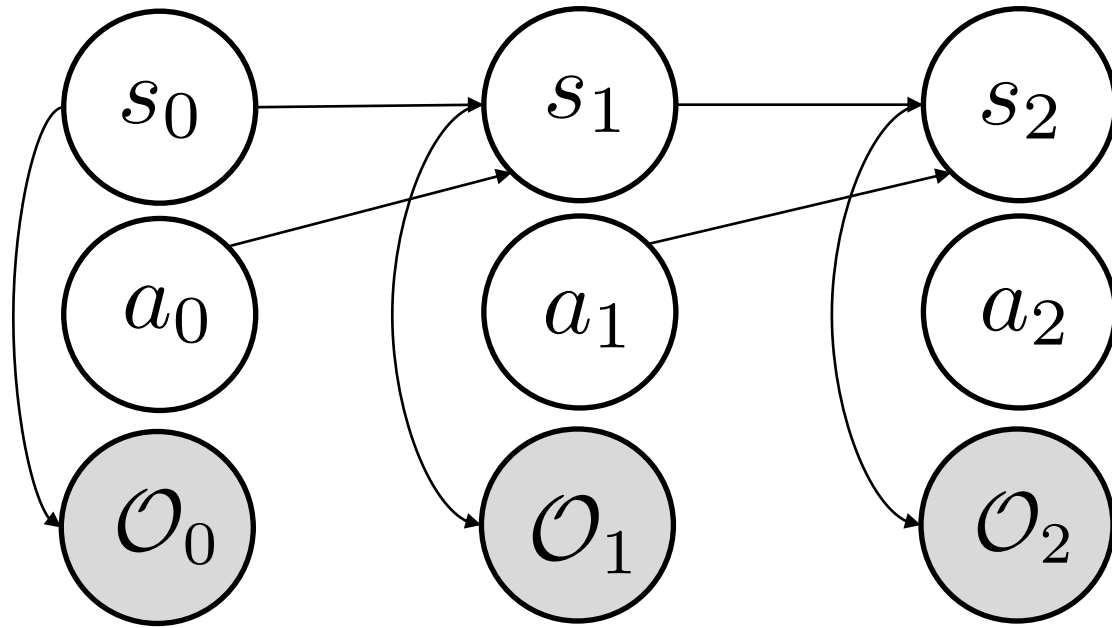
Set to uniform

$$3) \mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \left[ \cancel{\log p(s_0)} + \sum_t \left[ \cancel{\log p(a_t|s_t)} + \log p(s_{t+1}|s_t, a_t)} + \log p(\mathcal{O}_t|s_t, a_t) \right] \right] - \left[ \cancel{\log p(s_0)} + \sum_t \left[ \log q(a_t|s_t)} + \cancel{\log p(s_{t+1}|s_t, a_t)} \right] \right] \right]$$

$$4) \mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_t \log p(\mathcal{O}_t|s_t, a_t) - \log q(a_t|s_t) \right] \quad p(\mathcal{O}_t|s_t, a_t) = \exp(r(s_t, a_t))$$

$$5) \mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_t r(s_t, a_t) + \mathcal{H}(q(\cdot|s_t)) \right] \quad \begin{array}{l} \text{Maximum entropy RL} \\ \text{Gradient ascent = PG!} \end{array}$$

# Ok so what did we show?

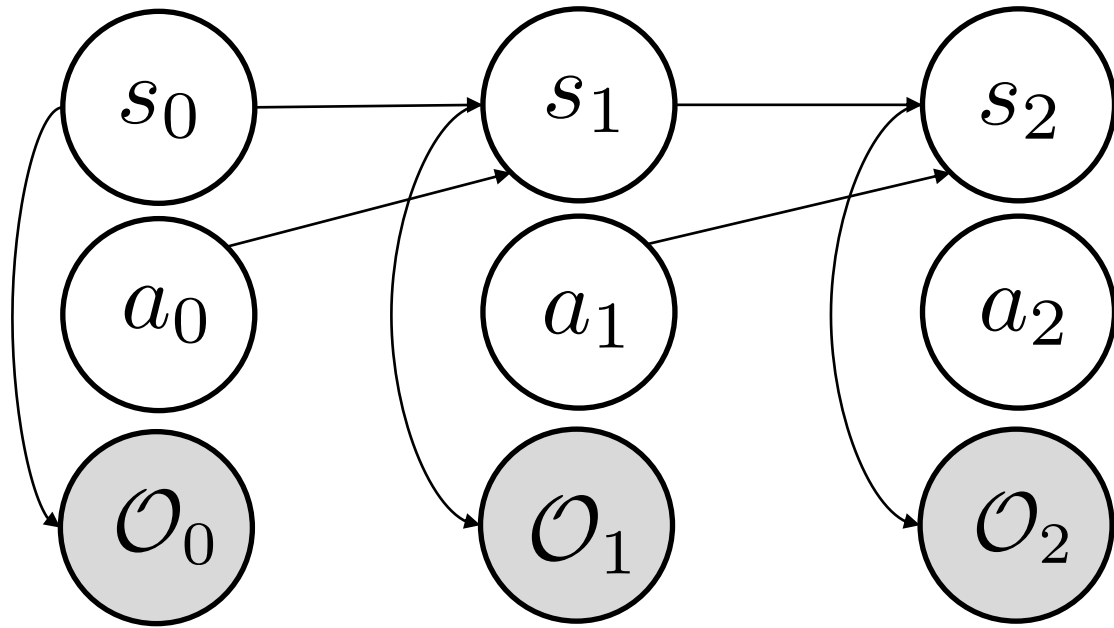


$$p(\mathcal{O}_t | s_t, a_t) = \exp(r(s_t, a_t))$$

Optimal Policy  $\rightarrow$  Posterior Inference

Approximate  $p(a_t | s_t, \mathcal{O}_{t:T} = 1)$  by  $q(a_t | s_t, \mathcal{O}_{t:T} = 1)$

# Ok so what did we show?



Optimal Policy  $\rightarrow$  Posterior Inference

$$p(a_t | s_t, \mathcal{O}_{t:T} = 1)$$

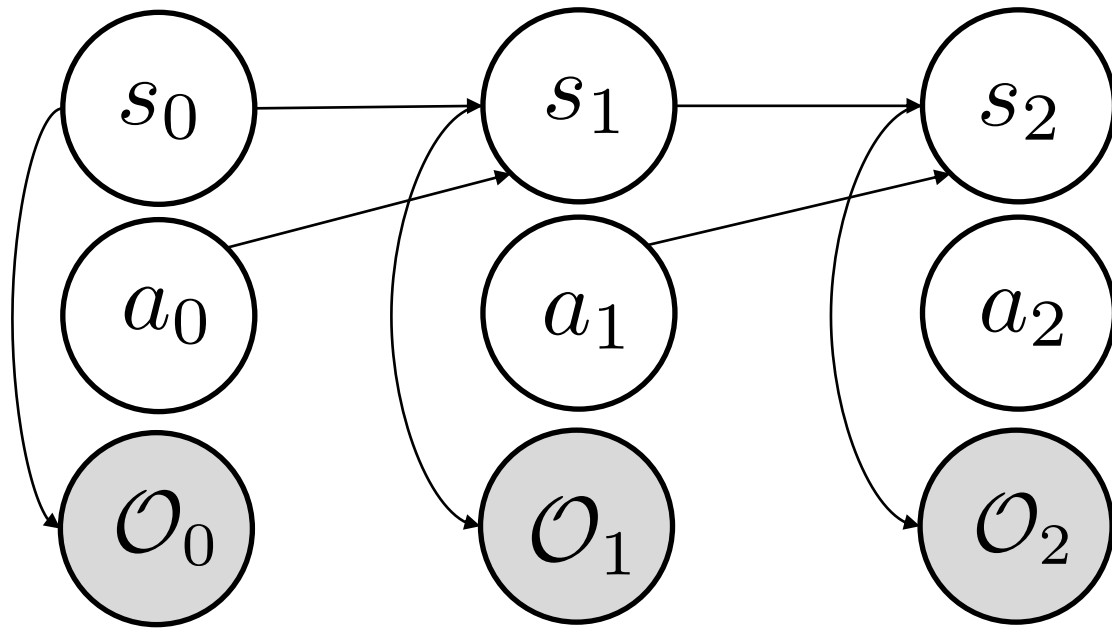
Difficult inference problem in closed form  
 $\rightarrow$  use variational inference

Find approximate posterior  $q(z|x)$  by optimizing the ELBO

$$\max_q \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z)) \right]$$

$$\begin{array}{ccc} x & & z \\ \updownarrow & & \updownarrow \\ (\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_T) & & (s_0, a_0, s_1, a_1, \dots, s_T, a_T) \end{array}$$

# Ok so what did we show?



Find approximate posterior  $q(z|x)$  by optimizing the ELBO

$$\max_q \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z)) \right]$$

$x$   $z$   
 $\updownarrow$   $\updownarrow$   
 $(\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_T)$   $(s_0, a_0, s_1, a_1, \dots, s_T, a_T)$

$$q(s_0, a_0, \dots, s_T, a_T | \mathcal{O}_0, \dots, \mathcal{O}_T) = p(s_0) \prod_{t=0}^T p(s_{t+1} | s_t, a_t) q(a_t | s_t)$$

$$\mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \sum_t \log p(\mathcal{O}_t | s_t, a_t) - \log q(a_t | s_t) \right] = \mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \sum_t r(s_t, a_t) + \mathcal{H}(q(\cdot | s_t)) \right]$$

Maximize ELBO with SGD = policy gradient!

# Lecture outline

---

Control as Inference - Formulation



Variational Inference



Control as Inference to Derive Policy Gradient

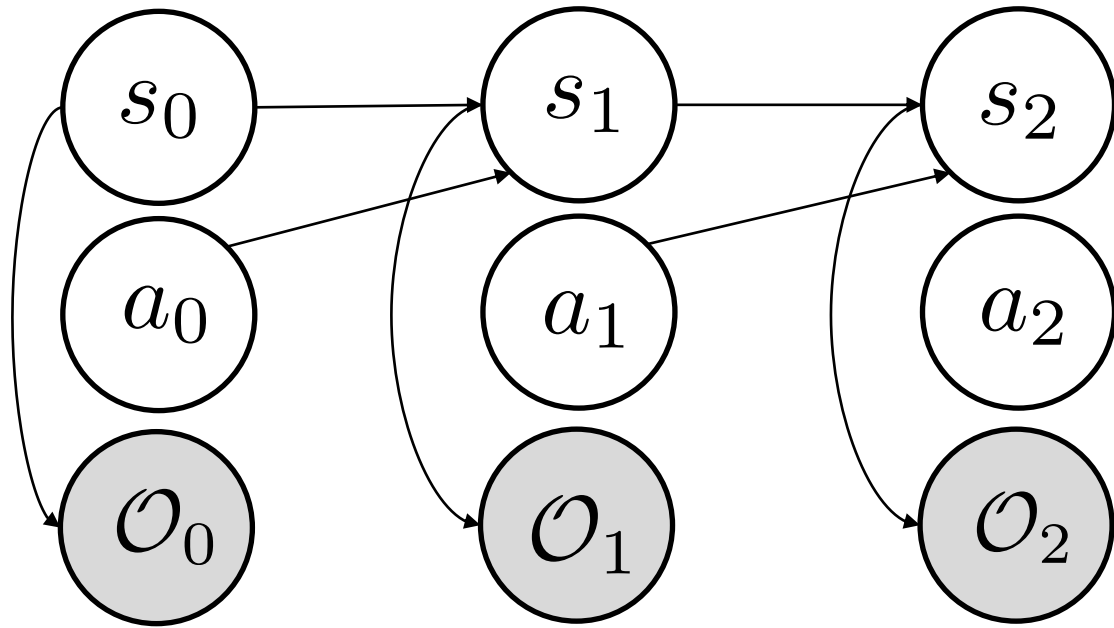


Control as Inference to Derive Q-learning



Control as Inference to Derive Model-Based RL

# Can we derive (soft) Q-learning from the ELBO?



Find approximate posterior  $q(z|x)$  by optimizing the ELBO

$$\max_q \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z)) \right]$$

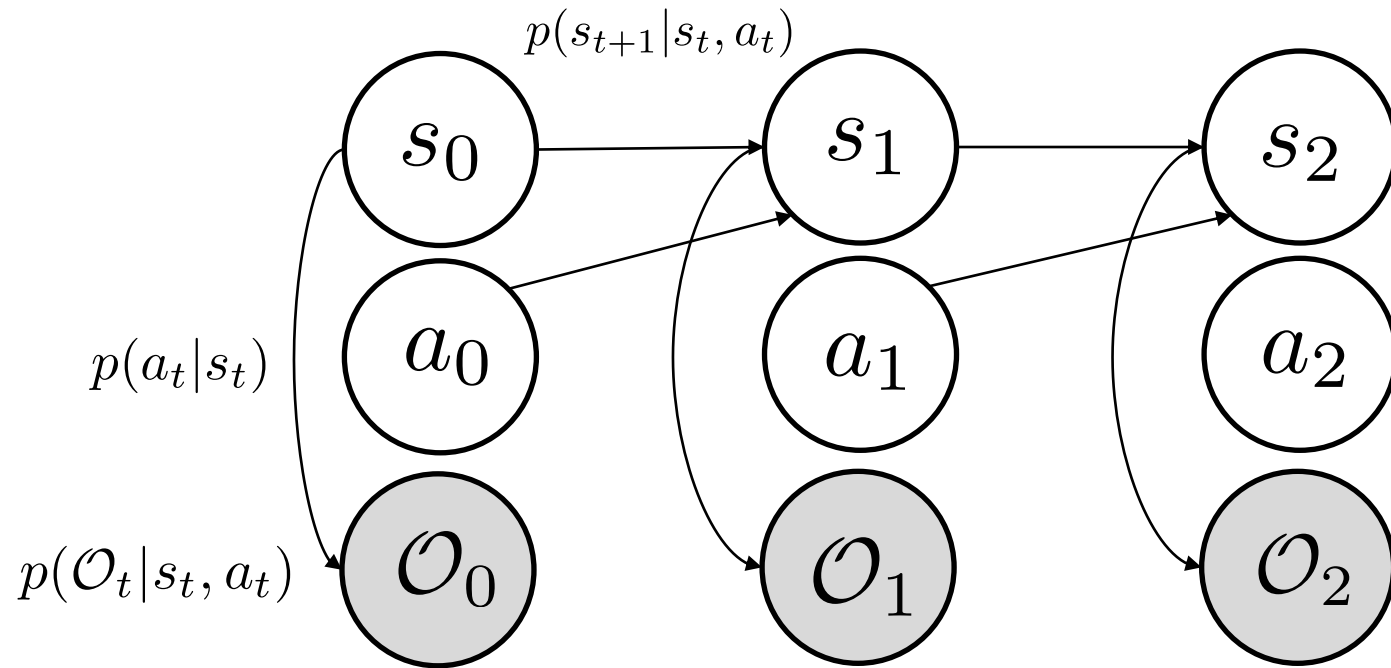
$$\begin{array}{ccc} & x & z \\ & \updownarrow & \updownarrow \\ (\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_T) & & (s_0, a_0, s_1, a_1, \dots, s_T, a_T) \end{array}$$

$$q(s_0, a_0, \dots, s_T, a_T | \mathcal{O}_0, \dots, \mathcal{O}_T) = p(s_0) \prod_{t=0}^T p(s_{t+1} | s_t, a_t) q(a_t | s_t)$$

$$\mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \sum_t \log p(\mathcal{O}_t | s_t, a_t) - \log q(a_t | s_t) \right] = \mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \sum_t r(s_t, a_t) + \mathcal{H}(q(\cdot | s_t)) \right]$$

Maximize ELBO with DP = Soft Q learning!

# So what are we doing inference over?



$$p(O_t|s_t, a_t) = \exp(r(s_t, a_t))$$

$$p(\tau|\mathcal{O}_{0:T} = 1) \propto p(\tau) \exp\left(\sum_{t=0}^T r(s_t, a_t)\right)$$

Use case 1:

Derive soft RL algorithms

Analogues for optimal Q and V

$$V(s_t) = \log p(\mathcal{O}_{t:T} = 1|s_t)$$

$$Q(s_t, a_t) = \log p(\mathcal{O}_{t:T} = 1|s_t, a_t)$$

“Likelihood of being optimal in the future at some state, action”

# Let's optimize the last step of the ELBO

$$\mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \sum_t r(s_t, a_t) + \mathcal{H}(q(\cdot | s_t)) \right]$$

Consider the last time step

$$\begin{aligned} \max \mathbb{E}_{s_T \sim q(s_T)} & \left[ \mathbb{E}_{a_T \sim q(a_T | s_T)} [r(s_T, a_T) - \log q(a_T | s_T)] \right] \\ & \text{(log-exp)} \\ & = \mathbb{E}_{s_T \sim q(s_T)} \left[ \mathbb{E}_{a_T \sim q(a_T | s_T)} [\log \exp(r(s_T, a_T)) - \log q(a_T | s_T)] \right] \\ & = \mathbb{E}_{s_T \sim q(s_T)} \left[ \mathbb{E}_{a_T \sim q(a_T | s_T)} \left[ \log \exp(r(s_T, a_T)) - \log \int \exp(r(s_T, a_T)) da_T - \log q(a_T | s_T) \right] + \log \int \exp(r(s_T, a_T)) da_T \right] \\ & \text{(Add subtract to normalize)} \\ & = \mathbb{E}_{s_T \sim q(s_T)} \left[ \mathbb{E}_{a_T \sim q(a_T | s_T)} \left[ \log \frac{\exp(r(s_T, a_T))}{\int \exp(r(s_T, a_T)) da_T} - \log q(a_T | s_T) \right] \right] \\ & = \mathbb{E}_{s_T \sim q(s_T)} \left[ \mathbb{E}_{a_T \sim q(a_T | s_T)} \left[ -\log \frac{q(a_T | s_T)}{\frac{\exp(r(s_T, a_T))}{\int \exp(r(s_T, a_T)) da_T}} \right] \right] = \mathbb{E}_{s_T \sim q(s_T)} \left[ -D_{KL}(q(a_T | s_T) \parallel \frac{\exp(r(s_T, a_T))}{\int \exp(r(s_T, a_T)) da_T}) \right] \\ & \text{(Definition of KL divergence)} \end{aligned}$$

# Let's optimize the last step of the ELBO

$$\mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \sum_t r(s_t, a_t) + \mathcal{H}(q(\cdot | s_t)) \right]$$

Consider the last time step

$$\max \mathbb{E}_{s_T \sim q(s_T)} \left[ \mathbb{E}_{a_T \sim q(a_T | s_T)} [r(s_T, a_T) - \log q(a_T | s_T)] \right]$$

$$= \mathbb{E}_{s_T \sim q(s_T)} \left[ -D_{KL}(q(a_T | s_T) \parallel \frac{\exp(r(s_T, a_T))}{\int \exp(r(s_T, a_T)) da_T}) \right]$$

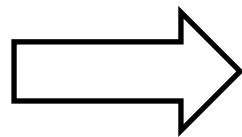
KL-divergence  $D(p, q)$  is always non negative and is minimized when  $p = q$

$$q(a_T | s_T) = \frac{\exp(r(s_T, a_T))}{\int \exp(r(s_T, a_T)) da_T}$$

Ok let's simplify

$$Q(s_T, a_T) = r(s_T, a_T)$$

$$V(s_T) = \log \int \exp(r(s_T, a_T)) da_T$$



$$q(a_T | s_T) = \exp(Q(s_T, a_T) - V(s_T))$$

Optimal policy is proportional to exponential advantage  
(soft-max)



# Let's optimize the step before

$$\mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \sum_t r(s_t, a_t) + \mathcal{H}(q(\cdot | s_t)) \right]$$

Consider the second last time step

$$\arg \max_q \mathbb{E}_{s_{T-1} \sim q(s_{T-1})} \left[ \mathbb{E}_{a_{T-1} \sim q(a_{T-1} | s_{T-1})} \left[ r(s_{T-1}, a_{T-1}) - \log q(a_{T-1} | s_{T-1}) + \mathbb{E}_{\substack{s_T \sim p(s_T | s_{T-1}, a_{T-1}) \\ a_T \sim q(a_T | s_T)}} [r(s_T, a_T) - \log q(a_T | s_T)] \right] \right]$$

Exactly what we computed in the last step

$$\arg \max_q \mathbb{E}_{s_{T-1} \sim q(s_{T-1})} \left[ \mathbb{E}_{a_{T-1} \sim q(a_{T-1} | s_{T-1})} \left[ r(s_{T-1}, a_{T-1}) - \log q(a_{T-1} | s_{T-1}) + \mathbb{E}_{s_T \sim p(s_T | s_{T-1}, a_{T-1})} [V(s_T)] \right] \right]$$

From the last slide

Let us call this  $Q(s_{T-1}, a_{T-1})$

$$Q(s_{T-1}, a_{T-1}) = r(s_{T-1}, a_{T-1}) + \mathbb{E}_{s_T \sim p(s_T | s_{T-1}, a_{T-1})} [V(s_T)] \quad \text{(Looks like Bellman!)}$$

$$\arg \max_q \mathbb{E}_{s_{T-1} \sim q(s_{T-1})} \left[ \mathbb{E}_{a_{T-1} \sim q(a_{T-1} | s_{T-1})} \left[ Q(s_{T-1}, a_{T-1}) - \log q(a_{T-1} | s_{T-1}) \right] \right]$$

Looks a lot like the previous time-step

# Let's optimize the step before

$$\mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \sum_t r(s_t, a_t) + \mathcal{H}(q(\cdot | s_t)) \right]$$

Consider the second last time step

$$\arg \max_q \mathbb{E}_{s_{T-1} \sim q(s_{T-1})} \left[ \mathbb{E}_{a_{T-1} \sim q(a_{T-1} | s_{T-1})} \left[ Q(s_{T-1}, a_{T-1}) - \log q(a_{T-1} | s_{T-1}) \right] \right]$$

$$Q(s_{T-1}, a_{T-1}) = r(s_{T-1}, a_{T-1}) + \mathbb{E}_{s_T \sim p(s_T | s_{T-1}, a_{T-1})} [V(s_T)]$$

Referring back to the last time step math and pattern matching

$$V(s_{T-1}) = \log \int \exp(Q(s_{T-1}, a_{T-1})) da_{T-1}$$

$$Q(s_{T-1}, a_{T-1}) = r(s_{T-1}, a_{T-1}) + \mathbb{E}_{s_T \sim p(s_T | s_{T-1}, a_{T-1})} [V(s_T)]$$

$$q(a_{T-1} | s_{T-1}) = \exp(Q(s_{T-1}, a_{T-1}) - V(s_{T-1}))$$

Optimal policy is proportional to exponential advantage  
(soft-max)

# Let's make it recursive

This suggests a recursive dynamic programming algorithm!

$$Q(s_T, a_T) = r(s_T, a_T)$$

$$V(s_T) = \log \int \exp(r(s_T, a_T)) da_T$$

For  $t = T-1$  to 1:

$$Q_t(s_t, a_t) = r(s_t, a_t) + \mathbb{E}_{s_{t+1} \sim p(s_{t+1}|s_t, a_t)} [V_{t+1}(s_{t+1})] \quad \text{(Bellman update)}$$

$$V_t(s_t) = \log \int \exp(Q_t(s_t, a_t)) da_t \quad \text{(Soft-max)}$$

$$q(a_t|s_t) = \exp(Q_t(s_t, a_t) - V_t(s_t)) \quad \text{(Soft-max)}$$

Very similar to the “soft” (entropy) Q-learning procedure from earlier lectures!

# What does this suggest as an algorithm?

Optimize a "soft" Bellman equation

$$Q(s_t, a_t) \leftarrow r_t + \gamma \mathbb{E}_{s_{t+1} \sim p_s} [V(s_{t+1})]$$

$$Q_{\text{soft}}(s_t, a_t) \leftarrow r_t + \gamma \mathbb{E}_{s_{t+1} \sim p_s} [V_{\text{soft}}(s_{t+1})]$$

$$V(s_t) \leftarrow \max_a Q(s_t, a)$$

$$V_{\text{soft}}(s_t) \leftarrow \alpha \log \int_{\mathcal{A}} \exp \left( \frac{1}{\alpha} Q_{\text{soft}}(s_t, a') \right) da'$$

$$\pi(a|s_t) \leftarrow \arg \max_a Q(s_t, a)$$

$$\pi_{\text{soft}}(a|s_t) = \exp \left( \frac{1}{\alpha} (Q_{\text{soft}}(s_t, a) - V_{\text{soft}}(s_t)) \right)$$

Go from max to "softmax" (imagine if  $\alpha$  goes to 0, it becomes a max)

Prevents premature collapse of exploration while smoothing out optimization landscape!

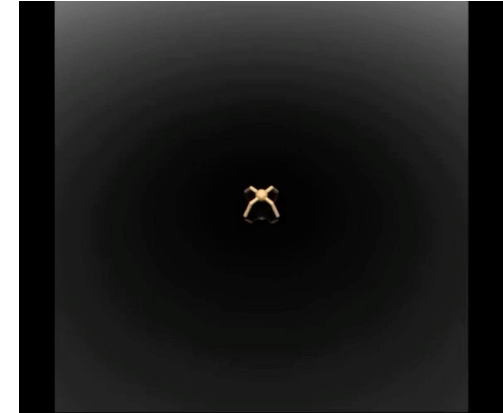
# Why should we ever do soft-Q learning?

## Optimization benefits

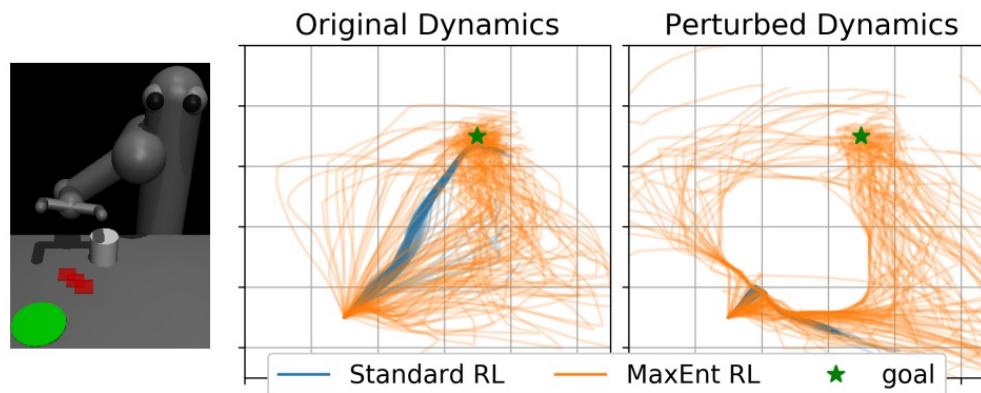
**Corollary 5.1.** (Iteration complexity with log barrier regularization) Let  $\beta_\lambda := \frac{8\gamma}{(1-\gamma)^3} + \frac{2\lambda}{|\mathcal{S}|}$ . Starting from any initial  $\theta^{(0)}$ , consider the updates (13) with  $\lambda = \frac{\epsilon(1-\gamma)}{2\left\|\frac{d\rho^*}{\mu}\right\|_\infty}$  and  $\eta = 1/\beta_\lambda$ . Then for all starting state distributions  $\rho$ , we have

$$\min_{t < T} \{V^*(\rho) - V^{(t)}(\rho)\} \leq \epsilon \quad \text{whenever} \quad T \geq \frac{320|\mathcal{S}|^2|\mathcal{A}|^2}{(1-\gamma)^6 \epsilon^2} \left\|\frac{d\rho^*}{\mu}\right\|_\infty^2.$$

## Transfer



## Deals better with misspecification



# Ok so what did we show?

Find approximate posterior  $q(z|x)$  by optimizing the ELBO using dynamic programming

$$\mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_t \log p(\mathcal{O}_t | s_t, a_t) - \log q(a_t | s_t) \right] = \max_q \mathbb{E}_{x \sim p(x)} \left[ \mathbb{E}_{z \sim q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z)) \right] \\ = \mathbb{E}_{\substack{s_0 \sim p(s_0) \\ a_t \sim q(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_t r(s_t, a_t) + \mathcal{H}(q(\cdot | s_t)) \right]$$

Can derive a "soft" dynamic programming Q-learning update

For  $t = T-1$  to 1:

$$Q_t(s_t, a_t) = r(s_t, a_t) + \mathbb{E}_{s_{t+1} \sim p(s_{t+1}|s_t, a_t)} [V_{t+1}(s_{t+1})] \quad (\text{Bellman update})$$

$$V_t(s_t) = \log \int \exp(Q(s_t, a_t)) da_t \quad (\text{Soft-max})$$

$$q(a_t | s_t) = \exp(Q_t(s_t, a_t) - V_t(s_t)) \quad (\text{Soft-max})$$

# Lecture outline

---

Control as Inference - Formulation



Variational Inference



Control as Inference to Derive Policy Gradient

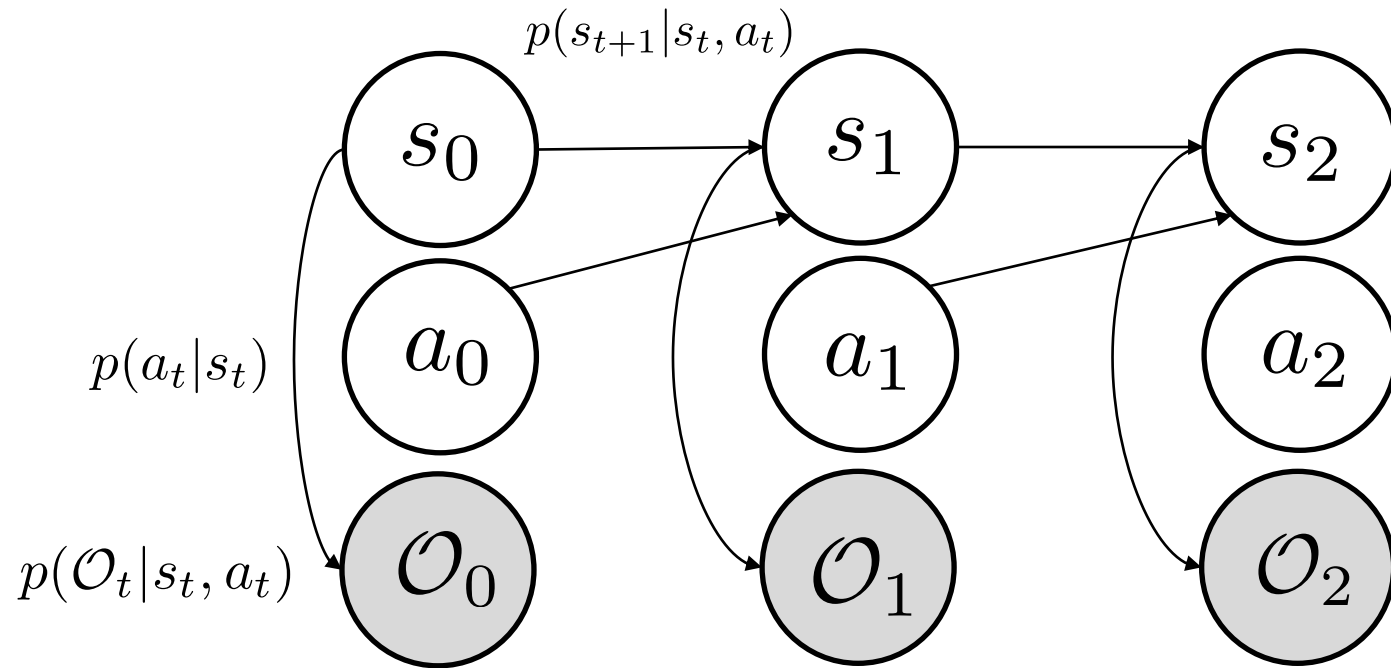


Control as Inference to Derive Q-learning



Control as Inference to Derive Model-Based RL

# Let's back up from VI to max likelihood



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

$$p(\tau | \mathcal{O}_{0:T} = 1) \propto p(\tau) \exp\left(\sum_{t=0}^T r(s_t, a_t)\right)$$

Let us assume we get a bunch of data of  $(s, a, s', r)$  from the true system  $p$

We will try to learn a surrogate model  $\hat{p}$  to approximate  $p$ , use it for posterior inference

# Model Learning via Maximum Likelihood

$$\min_{\hat{p}} D_{KL}(p(s_0, \dots, s_T, a_0, \dots, a_T, \mathcal{O}_0, \dots, \mathcal{O}_T) || \hat{p}(s_0, \dots, s_T, a_0, \dots, a_T, \mathcal{O}_0, \dots, \mathcal{O}_T))$$

↓ Definition of KLD

$$\max_{\hat{p}} \mathbb{E}_{p(s_0, \dots, s_T, a_0, \dots, a_T, \mathcal{O}_0, \dots, \mathcal{O}_T)} [\log \hat{p}(s_0, \dots, s_T, a_0, \dots, a_T, \mathcal{O}_0, \dots, \mathcal{O}_T)]$$

↓ Expansion of joint

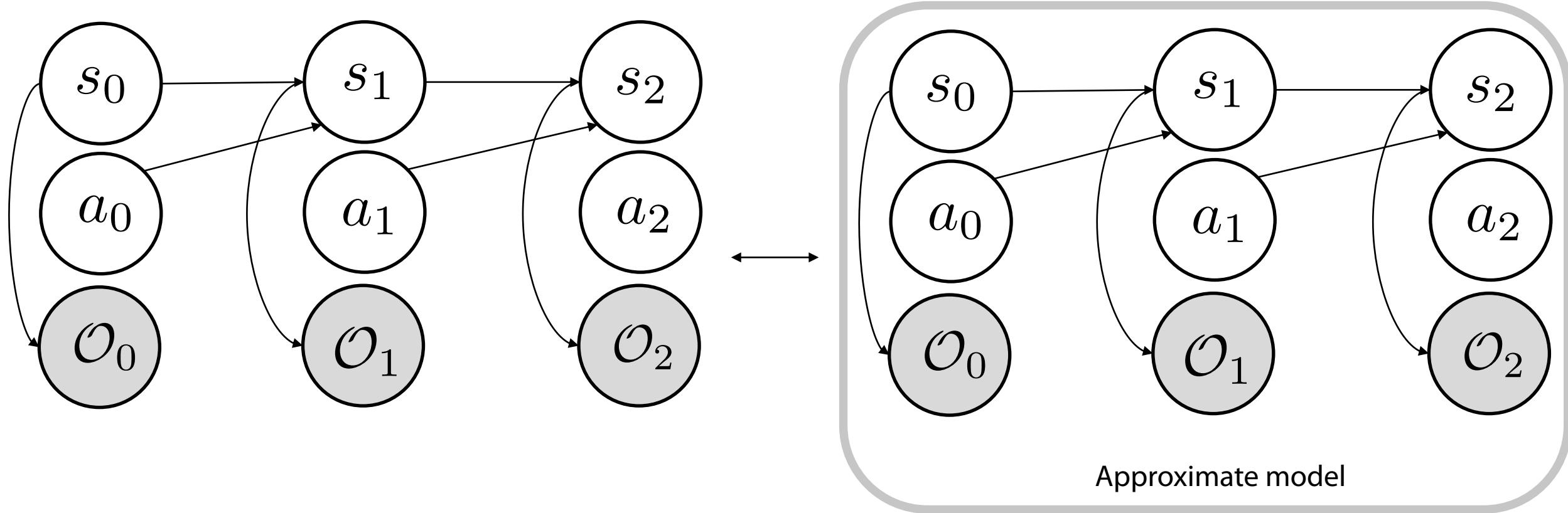
$$\max_{\hat{p}} \mathbb{E}_{p(s_0, \dots, s_T, a_0, \dots, a_T, \mathcal{O}_0, \dots, \mathcal{O}_T)} \left[ \log \hat{p}(s_0) + \sum_t [\log \hat{p}(s_{t+1} | s_t, a_t) + \log \hat{p}(\mathcal{O}_t | s_t, a_t)] \right]$$

Model learning

Reward learning

Fitting  $\hat{p}$  amounts to supervised learning on dynamics and rewards

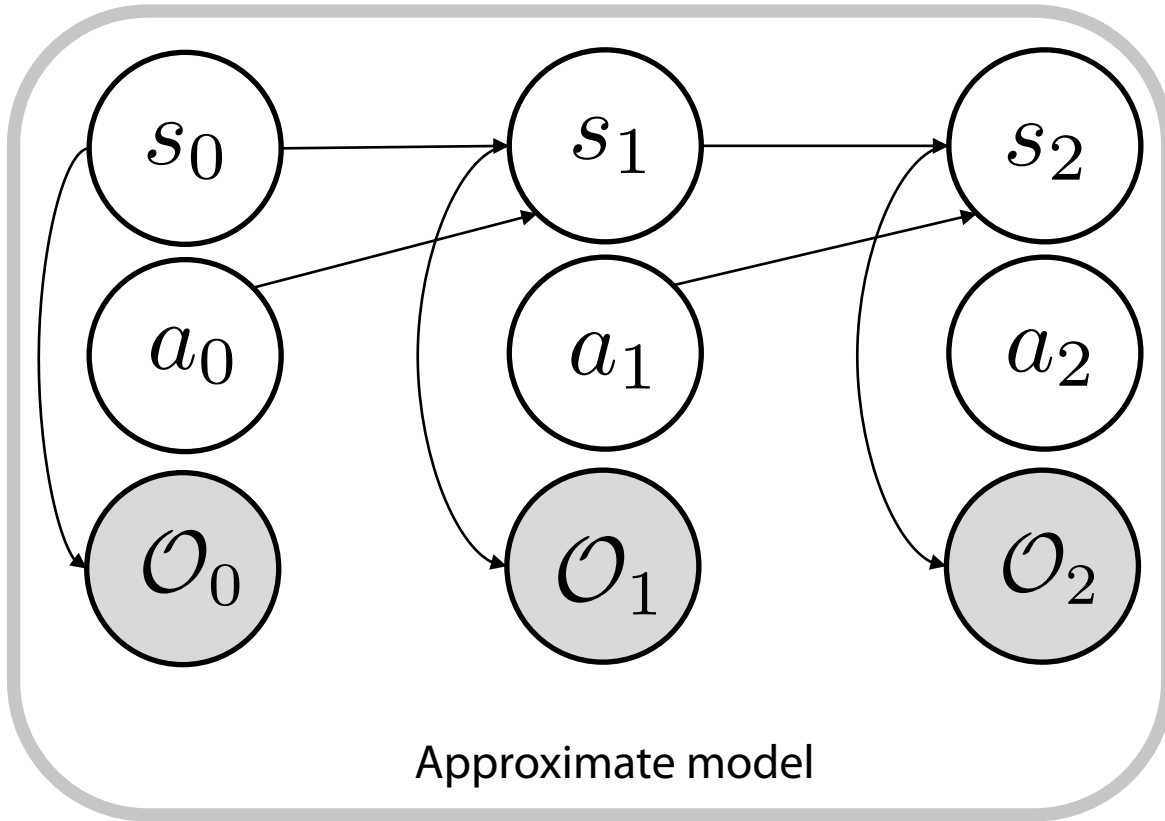
# Model Learning via Maximum Likelihood



Fitting  $\hat{p}$  amounts to supervised learning on dynamics and rewards

How do we actually use this approximate model to obtain optimal actions?

# Policy Extraction via Posterior Inference



Key idea: pretend that approximate model  $\hat{p}$  is the true model



$$\hat{p}(a_t | s_t, \mathcal{O}_{t:T} = 1)$$

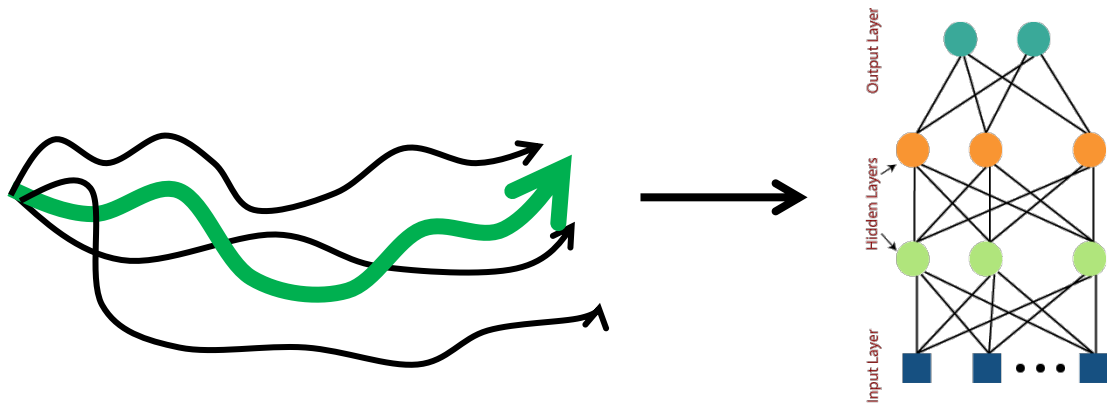
Just like in MFRL  $\rightarrow$   
perform posterior inference

Certainty equivalence  $\longleftarrow$  But pretend that the model were true

# Ok so how we do perform this inference?

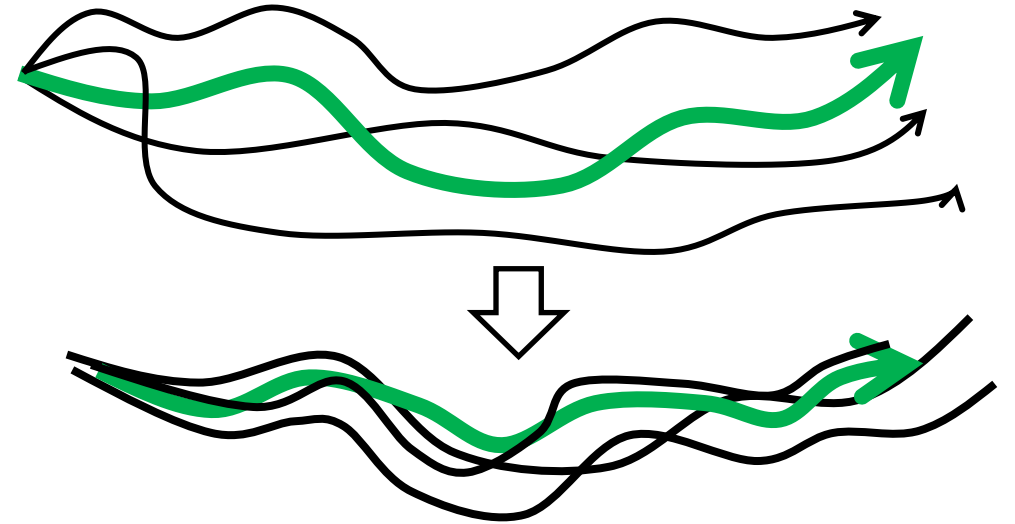
$$\hat{p}(a_t | s_t, \mathcal{O}_{t:T} = 1)$$

Idea 1: Variational inference in  $\hat{p}$



Model-based policy optimization methods (Dyna ++)

Idea 2: Use Monte-Carlo Sampling for Inference



MPPI-style planning methods

# Equivalence between posterior inference and MPPI

$$\hat{p}(a_t | s_t, \mathcal{O}_{t:T} = 1)$$

Let's expand out the nasty integrals with Bayes rule

$$= \frac{\hat{p}(a_t, s_t, \mathcal{O}_{t:T} = 1)}{\hat{p}(s_t, \mathcal{O}_{t:T} = 1)}$$

$$= \frac{\int \int \cdots \int \hat{p}(a_t, s_t, a_{t+1}, s_{t+1}, \dots, a_T, s_T, \mathcal{O}_{t:T} = 1) ds_{t+1} da_{t+1} \dots ds_T da_T}{\hat{p}(s_t, \mathcal{O}_{t:T} = 1)}$$

$$\propto \int \int \cdots \int \hat{p}(a_t, s_t, a_{t+1}, s_{t+1}, \dots, a_T, s_T, \mathcal{O}_{t:T} = 1) ds_{t+1} da_{t+1} \dots ds_T da_T$$

# Equivalence between posterior inference and MPPI

$$\hat{p}(a_t | s_t, \mathcal{O}_{t:T} = 1)$$

$$\propto \int \int \cdots \int \hat{p}(a_t, s_t, a_{t+1}, s_{t+1}, \dots, a_T, s_T, \mathcal{O}_{t:T} = 1) ds_{t+1} da_{t+1} \dots ds_T da_T$$

$$\propto \int \int \cdots \int \hat{p}(s_0) \Pi_t \left[ \begin{array}{ccc} \text{Dynamics} & \text{Action prior} & \text{Optimality} \\ \hat{p}(s_{t+1} | s_t, a_t) & p(a_t | s_t) & p(\mathcal{O}_t | s_t, a_t) \end{array} \right] ds_{t+1} da_{t+1} \dots ds_T da_T$$

$$\propto \int \int \cdots \int \hat{p}(s_0) \Pi_t \left[ \hat{p}(s_{t+1} | s_t, a_t) p(a_t | s_t) \right] \exp \left[ \sum_t r(s_t, a_t) \right] ds_{t+1} da_{t+1} \dots ds_T da_T$$

Substituting optimality definition  $p(\mathcal{O}_t | s_t, a_t) = \exp(r(s_t, a_t))$

$$\propto \mathbb{E}_{\substack{s_0 \sim \hat{p}(s_0) \\ a_t \sim \hat{p}(a_t | s_t) \\ s_{t+1} \sim \hat{p}(s_{t+1} | s_t, a_t)}} \left[ \exp \left[ \sum_t r(s_t, a_t) \right] \right] \quad \text{Just using definition of expectation}$$

# Equivalence between posterior inference and MPPI

$$\hat{p}(a_t | s_t, \mathcal{O}_{t:T} = 1)$$

$$\propto \mathbb{E}_{\substack{s_0 \sim \hat{p}(s_0) \\ a_t \sim \hat{p}(a_t | s_t) \\ s_{t+1} \sim \hat{p}(s_{t+1} | s_t, a_t)}} \left[ \exp \left[ \sum_t r(s_t, a_t) \right] \right]$$

Taking a bunch of samples through model  $\rightarrow$

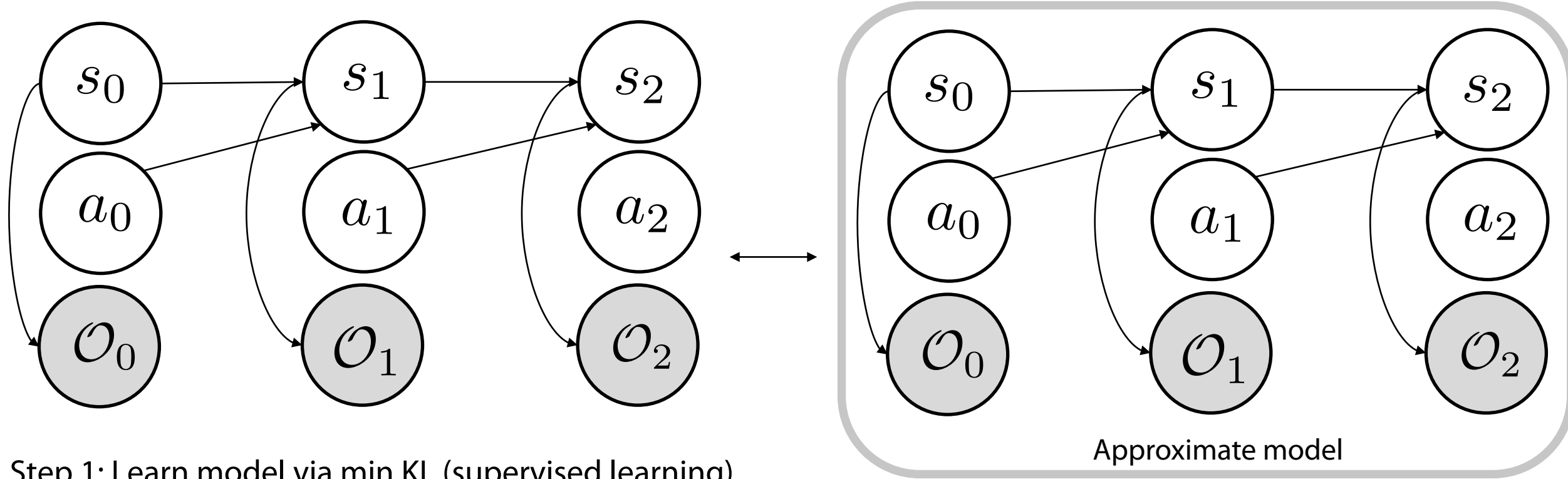
choose actions proportional to the expected sum of rewards

Can keep repeating with updated action prior

$$\propto \mathbb{E}_{\substack{s_0 \sim \hat{p}(s_0) \\ a_t \sim \hat{p}(a_t | s_t) \\ s_{t+1} \sim \hat{p}(s_{t+1} | s_t, a_t)}} \left[ \exp \left[ \sum_t r(s_t, a_t) \right] \right]$$

Can be thought of as a sampling-based Monte-Carlo approximation to posterior

# Ok so what did we show?



Step 1: Learn model via min KL (supervised learning)

$$\max_{\hat{p}} \mathbb{E}_{p(s_0, \dots, s_T, a_0, \dots, a_T, \mathcal{O}_0, \dots, \mathcal{O}_T)} \left[ \log \hat{p}(s_0) + \sum_t [\log \hat{p}(s_{t+1} | s_t, a_t) + \log \hat{p}(\mathcal{O}_t | s_t, a_t)] \right]$$

Step 2: Obtain posterior actions via Monte-Carlo approximation (approx MPPI)

$$\propto \mathbb{E}_{\substack{s_0 \sim \hat{p}(s_0) \\ a_t \sim \hat{p}(a_t | s_t) \\ s_{t+1} \sim \hat{p}(s_{t+1} | s_t, a_t)}} \left[ \exp \left[ \sum_t r(s_t, a_t) \right] \right]$$

# Lecture outline

---

Control as Inference - Formulation



Variational Inference



Control as Inference to Derive Policy Gradient



Control as Inference to Derive Q-learning



Control as Inference to Derive Model-Based RL

# Class Structure

