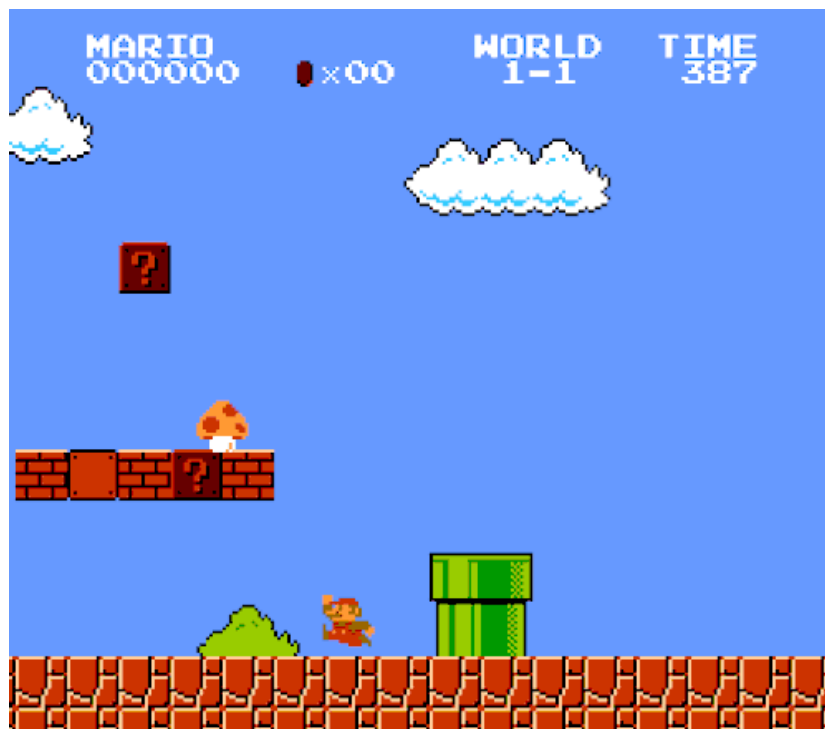# Reinforcement Learning
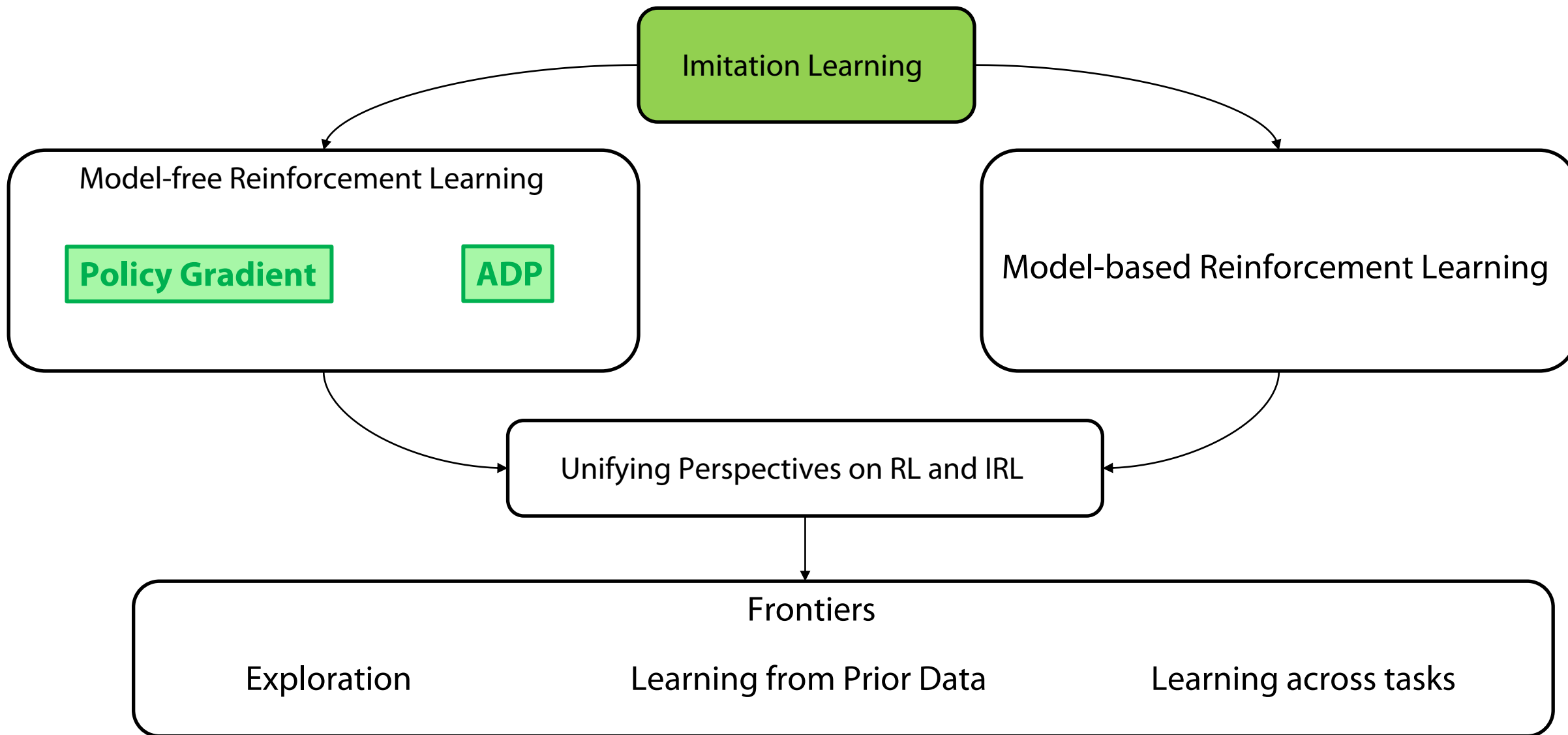# Spring 2024

Abhishek Gupta

TAs: Patrick Yin, Qiuyu Chen
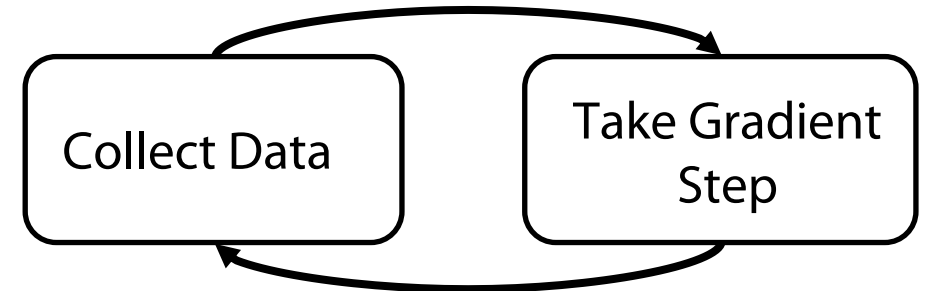
# Class Structure

# Resulting Algorithm (REINFORCE)

$$\nabla_\theta J(\theta) = \int p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) d\tau$$

$$\theta_{i+1} = \theta_i + \alpha \nabla_\theta J(\theta)|_{\theta=\theta_i}$$

Collect Data → Take Gradient Step

REINFORCE algorithm:

On-policy →

1. sample $\{\tau^i\}$ from $\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)$ (run it on the robot)
2. $\nabla_\theta J(\theta) \approx \sum_i \left( \sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i|\mathbf{s}_t^i) \right) \left( \sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$
3. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$
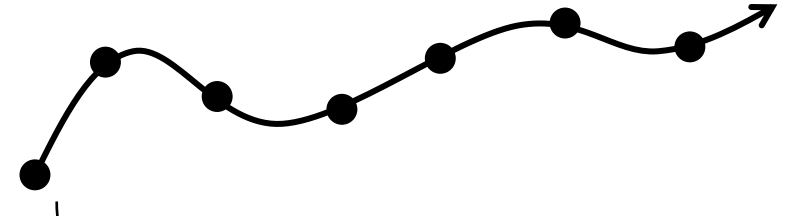
# Challenges in Policy Gradient

$$\nabla_\theta J(\theta) = \int p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) d\tau$$

$$\approx \frac{1}{N} \sum_{i=0}^{N} \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) \sum_{t'=0}^{T} r(s_{t'}^i, a_{t'}^i)$$
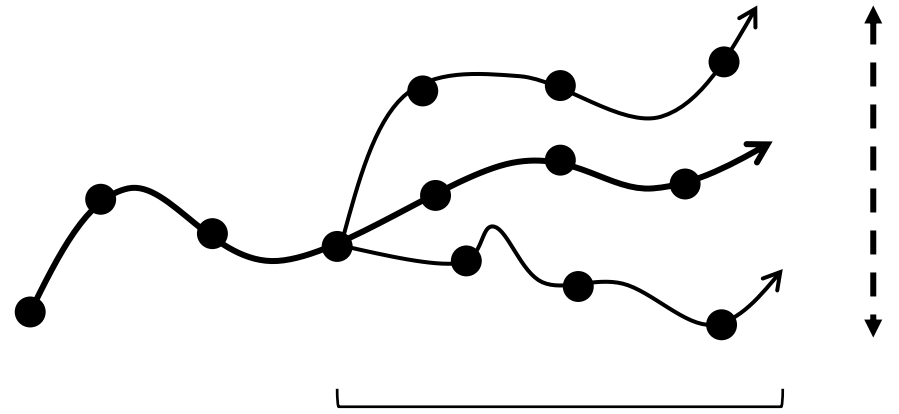
**High variance estimator!!**

Hard to tell what matters without many samples
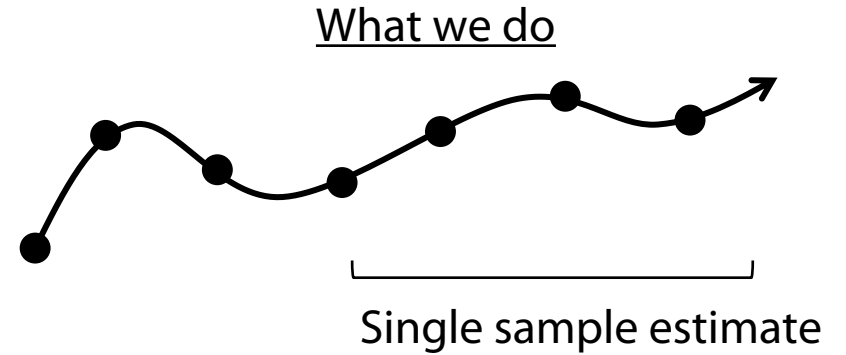
What we do

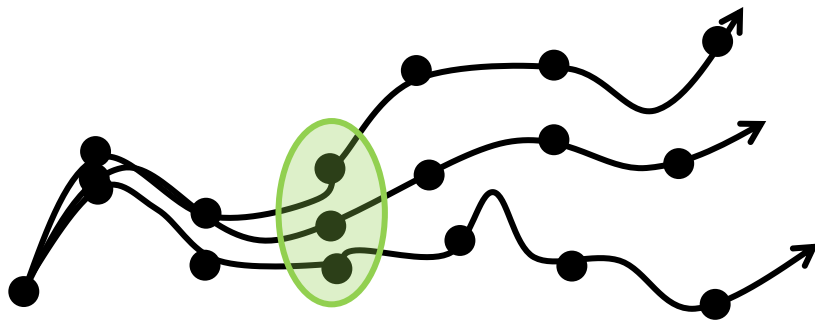Single sample estimate

What we actually want

Averaged return estimate
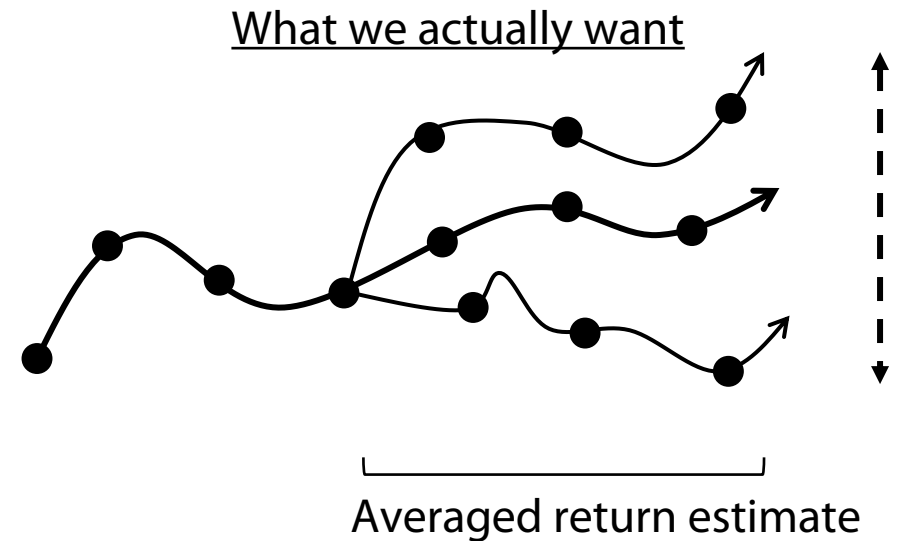
# What can we do to lower variance?

$$\nabla_\theta J(\theta) = \int p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) d\tau$$

$$\approx \frac{1}{N} \sum_{i=1}^{N} \sum_t \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) \sum_{t'=t}^{T} r(s_t^i, a_t^i)$$

Idea: bundle this across many (s, a) with a function approximator



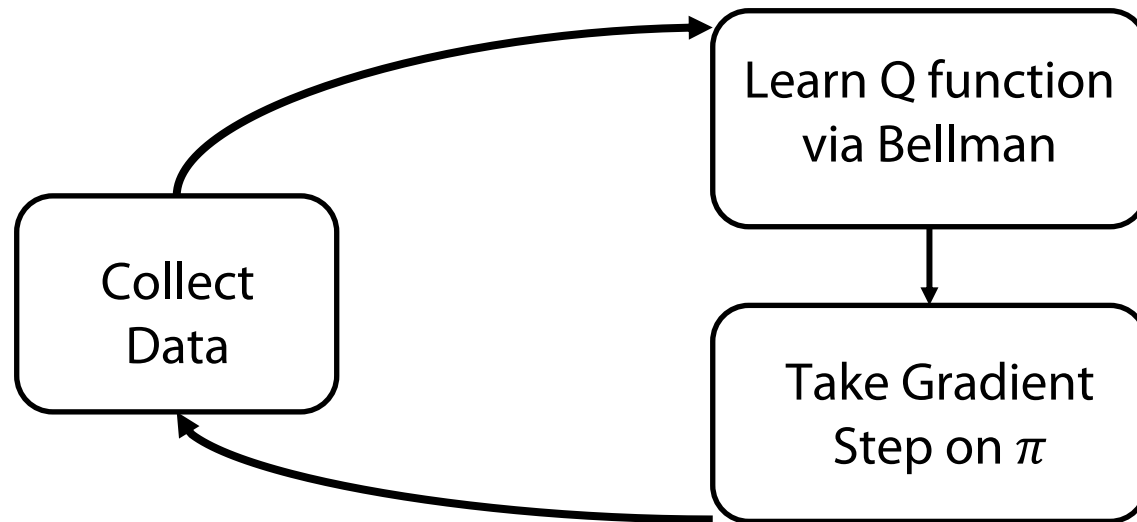Function approximator bundles return estimates across states

What we do



Single sample estimate

What we actually want



Averaged return estimate

# Recap of Off-Policy Reinforcement Learning

Critic: learned via the Bellman update (Policy Evaluation)

$$\min_\phi \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \mathcal{D}} \left( Q_\phi^\pi(s_t, a_t) - (r(s_t, a_t) + Q_{\hat{\phi}}^\pi(s_{t+1}, a_{t+1})) \right)^2 \quad a_{t+1} \sim \pi(\cdot | s_{t+1})$$

Learn Q function via Bellman

Collect Data

Take Gradient Step on $\pi$

Lowers variance and is off-policy!

Actor: updated using learned critic (Policy Improvement)

$$\max_\pi \mathbb{E}_{s \sim \mathcal{D}} \mathbb{E}_{a \sim \pi(.|s)} \left[ Q^\pi(s, a) \right]$$

# Pros/Cons of Off-Policy Methods in Robotics

Pros:
1. Sample-efficient enough for real world
2. Can learn from images with suitable design choices
3. Off-policy, can incorporate prior data

Cons
1. Often unstable
2. Can achieve lower asymptotic performance
3. Requires significant storage

# Lecture outline

The Anatomy of Model-Based Reinforcement Learning

$\downarrow$

Model based RL v0 → random shooting + MPC

$\downarrow$

Model based RL v1 → MPPI + MPC

$\downarrow$

Model based RL v2 → uncertainty based models

$\downarrow$

Model based RL v3 → policy optimization with models

$\downarrow$

Model based RL v4 → latent space models with images

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

**Gradient Ascent**

**Dynamic Programming**

Model-Based Optimization

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

1. Learn a surrogate model of the transition dynamics from arbitrary off-policy data
2. Do reward maximization against this model

Intuitive: learn how the world works first and then plan in that model
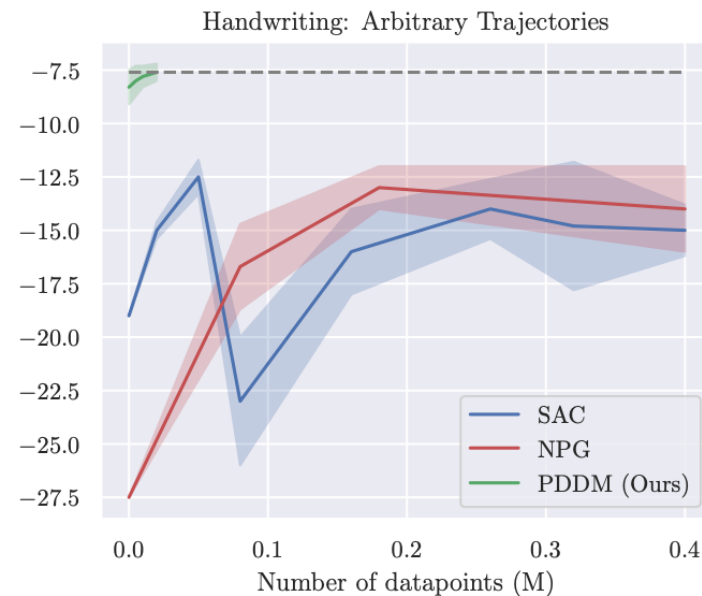
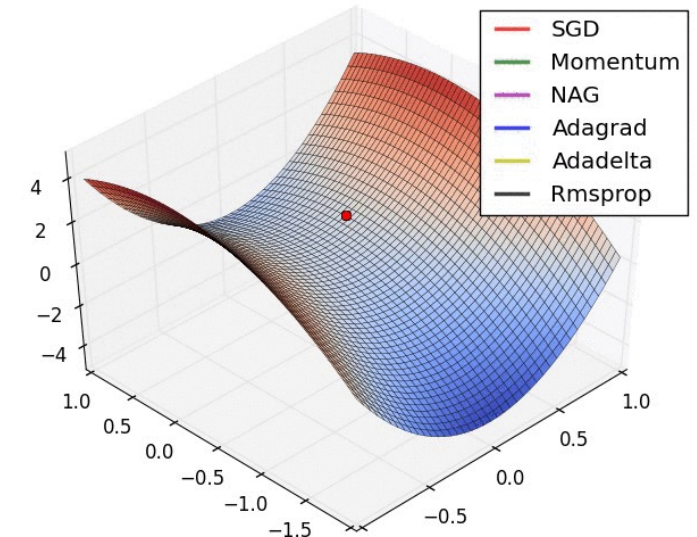# Why do model-based RL?
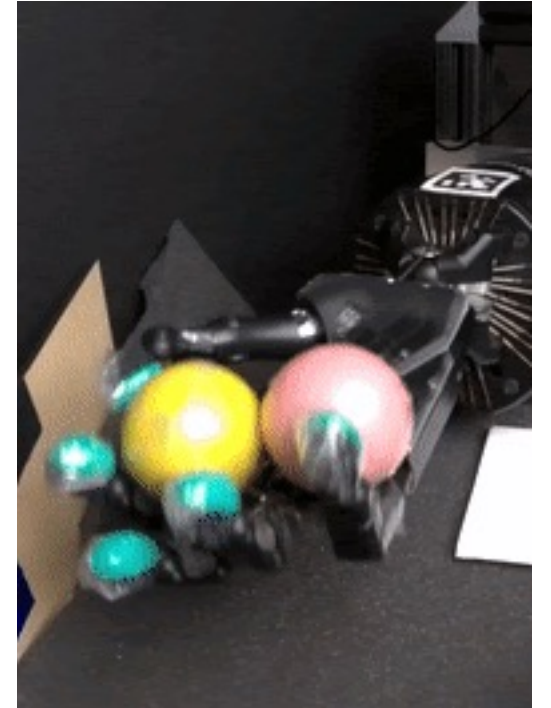
Why would we do this?

Transfer/Adaptive

Efficiency

Simplicity



Naturally off-policy!
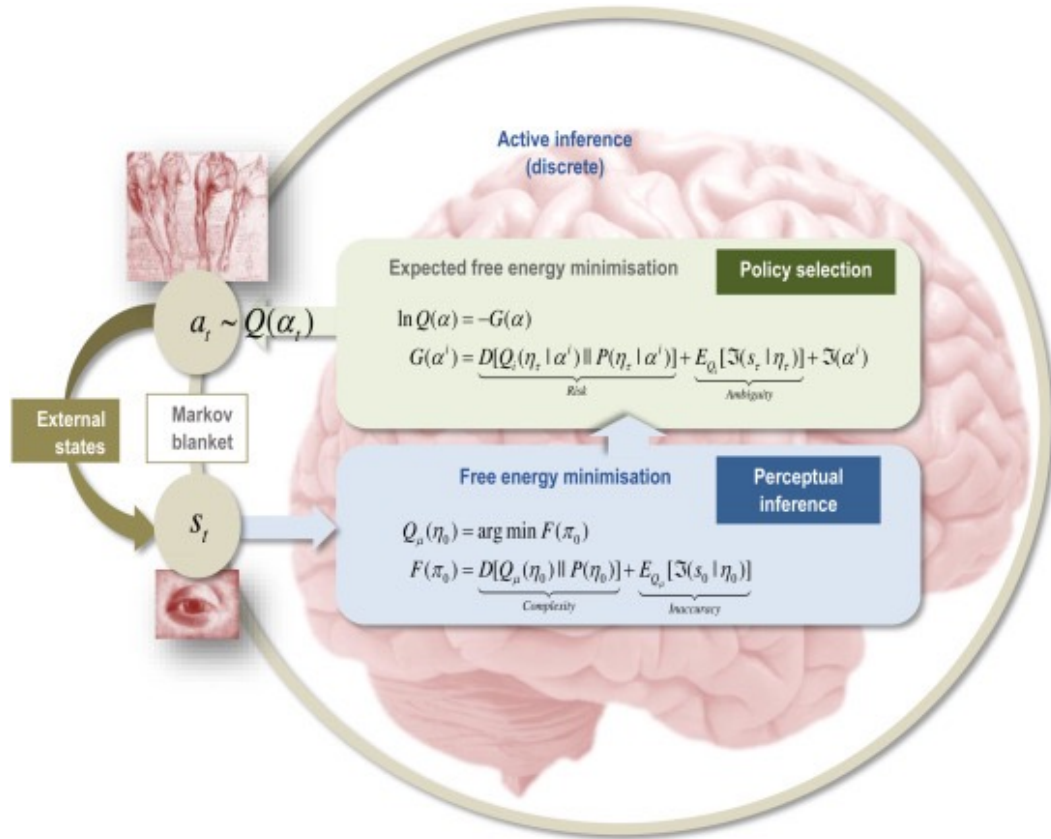
# Why do model-based RL?



Just 2 hours of real robot training

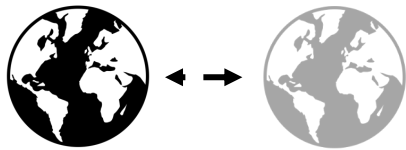Significant evidence for mechanisms for prediction of outcomes in neuro/cognitive science



**Reinforcement learning in the brain**

Yael Niv
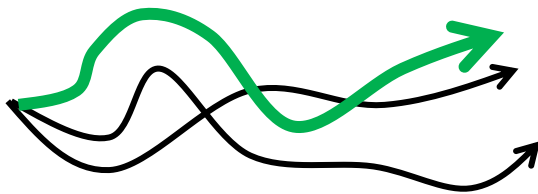Psychology Department & Princeton Neuroscience Institute, Princeton University

## Model Learning



$$\hat{p}_\theta \leftarrow \arg\min_{\hat{p}_\theta} \mathcal{L}(\mathcal{D}, \hat{p}_\theta)$$

## Planning



$$\arg\max_{\pi} \mathbb{E}_{\hat{p},\pi}\left[\sum_t r(s_t, a_t)\right]$$

Can also just be a single trajectory

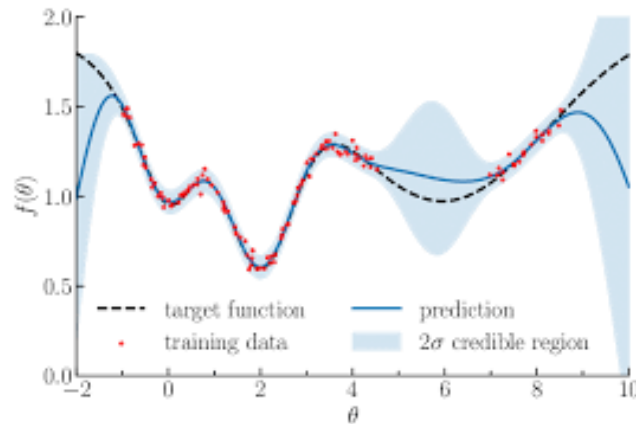How should we instantiate these?

# What will we not cover today?

### iLQR/iLQG

$$\min_{\mathbf{u}_1,\ldots,\mathbf{u}_T} \sum_{t=1}^{T} c(\mathbf{x}_t, \mathbf{u}_t) \text{ s.t. } \mathbf{x}_t = f(\mathbf{x}_{t-1}, \mathbf{u}_{t-1})$$
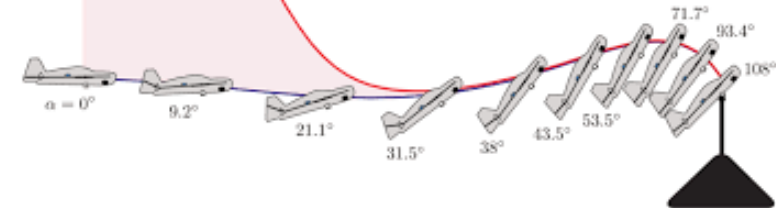
$$f(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{F}_t \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix} + \mathbf{f}_t$$

$$c(\mathbf{x}_t, \mathbf{u}_t) = \frac{1}{2} \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix}^T \mathbf{C}_t \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix} + \begin{bmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{bmatrix}^T \mathbf{c}_t$$

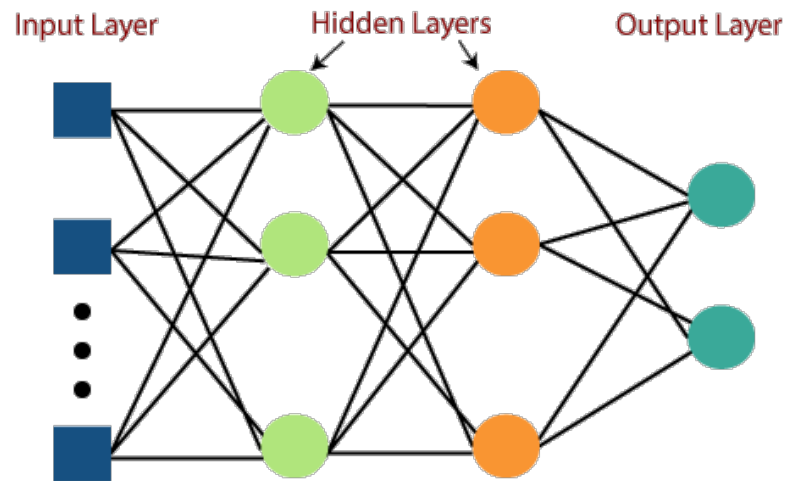### MBRL with GPs/Non-Parametrics



### Non-linear TrajOpt



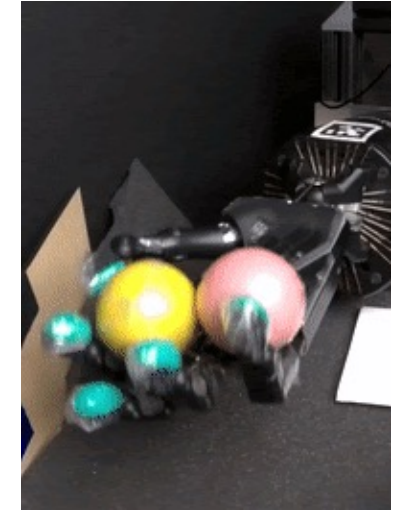Byron's lectures do a wonderful job, do go watch them!

# What will we cover today?
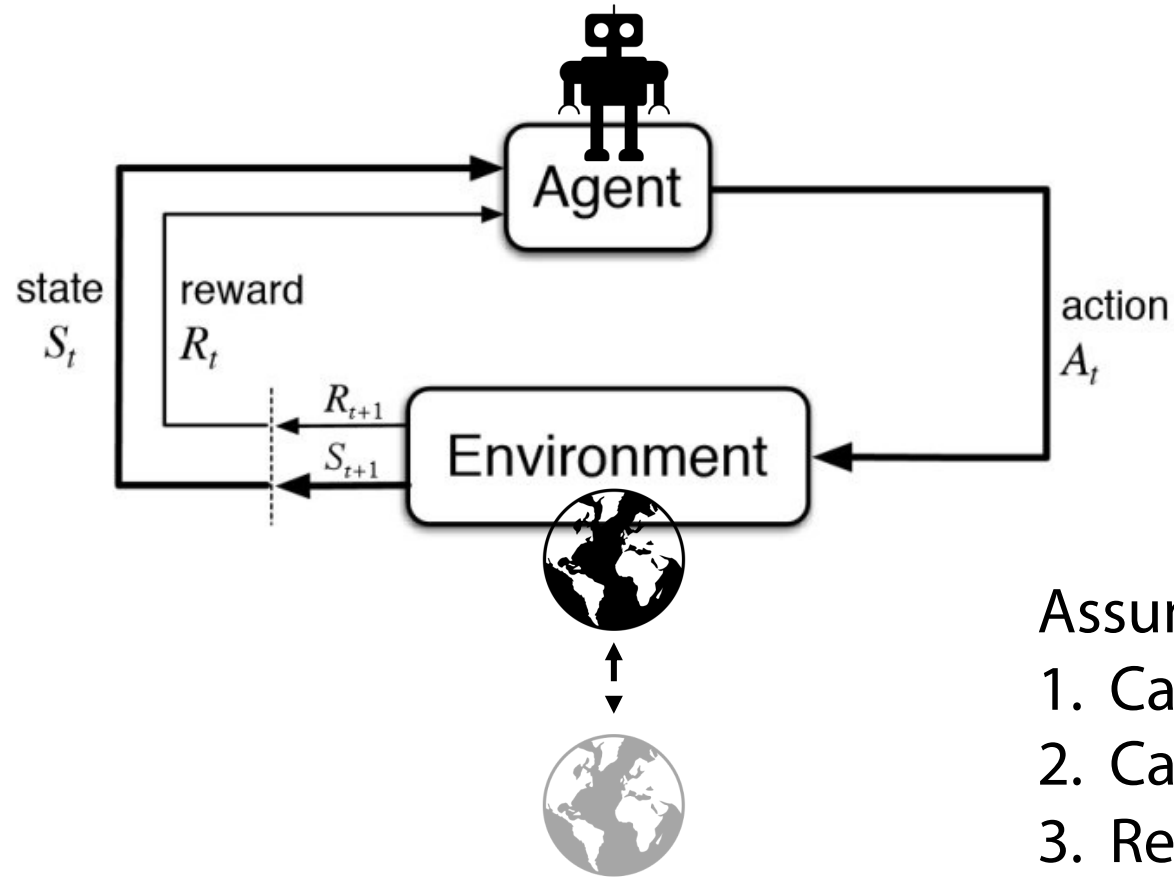
Use neural networks as our model!



$$\hat{p}_\theta \leftarrow \arg\min_{\hat{p}_\theta} \mathcal{L}(\mathcal{D}, \hat{p}_\theta)$$

$$\arg\max_{\pi} \mathbb{E}_{\hat{p},\pi} \left[ \sum_t r(s_t, a_t) \right]$$

Input Layer   Hidden Layers   Output Layer
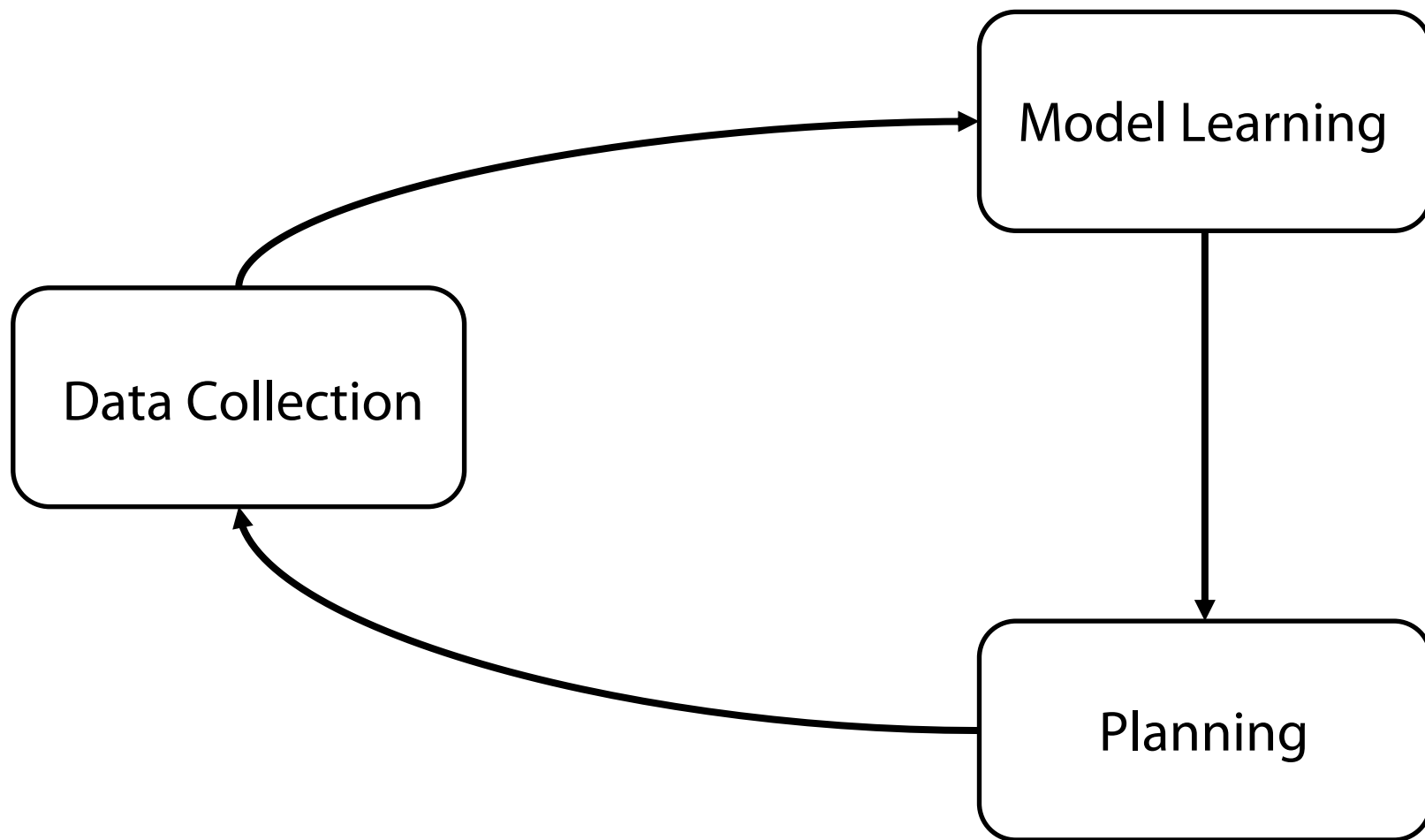


Testing lap time: 9.45 s

Trajectory
Rollouts

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

Assumptions:
1. Can only **sample** from dynamics
2. Can **reset** the environment
3. Reward function is **known**
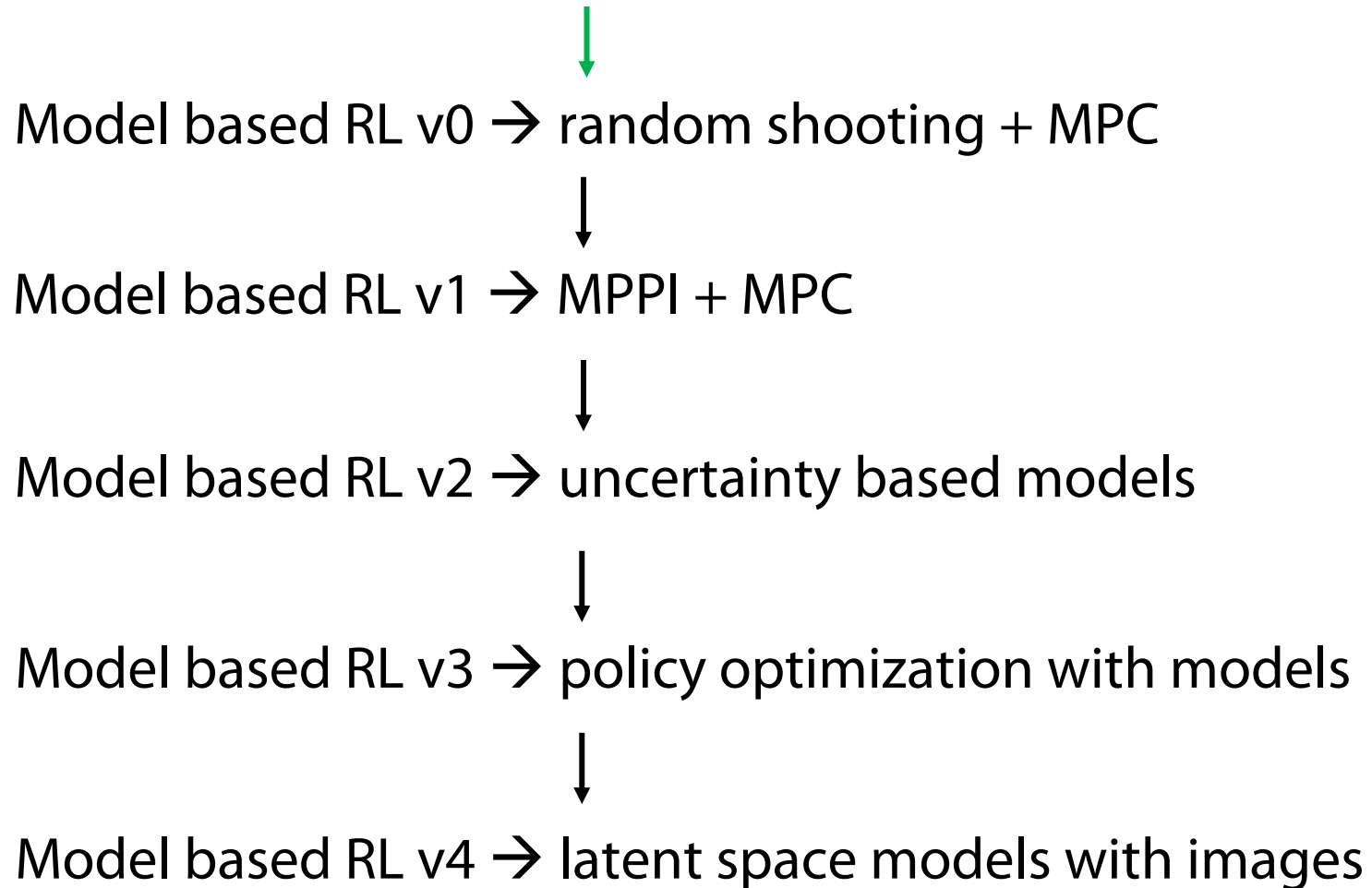
We will get into this in a later lecture!
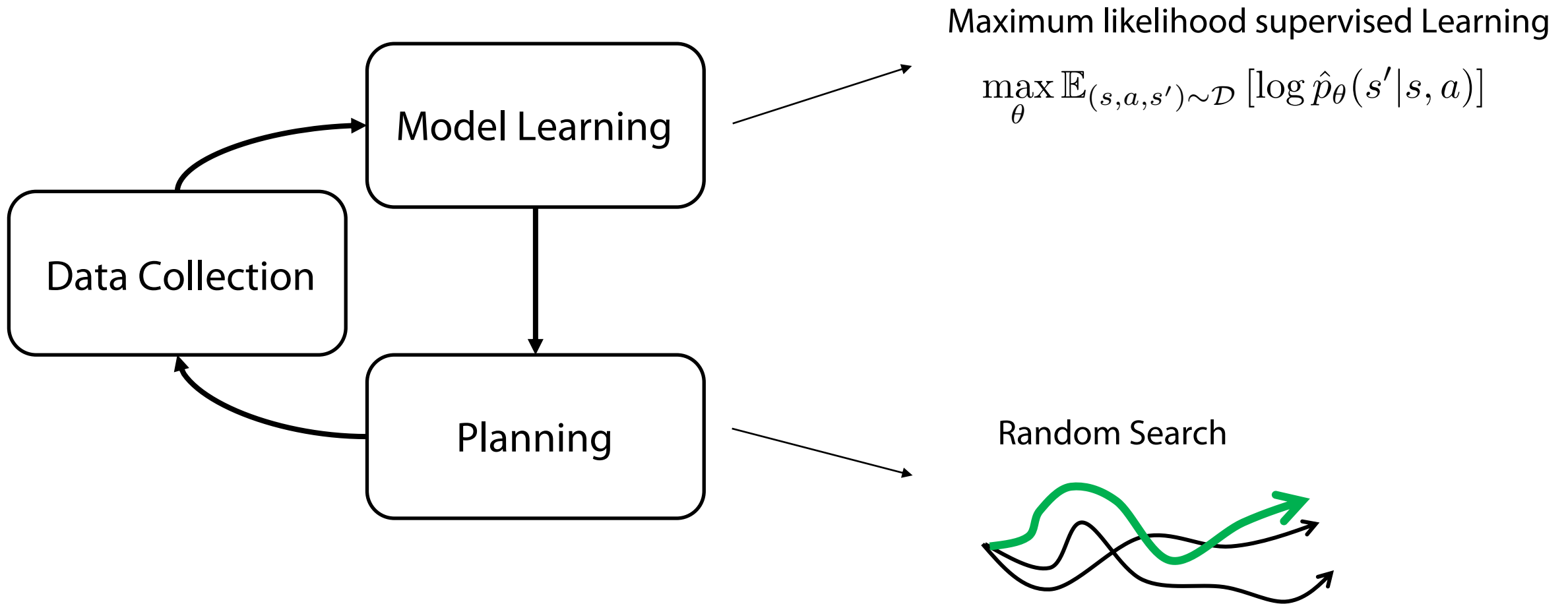
# Model Based RL – A template

# Lecture outline
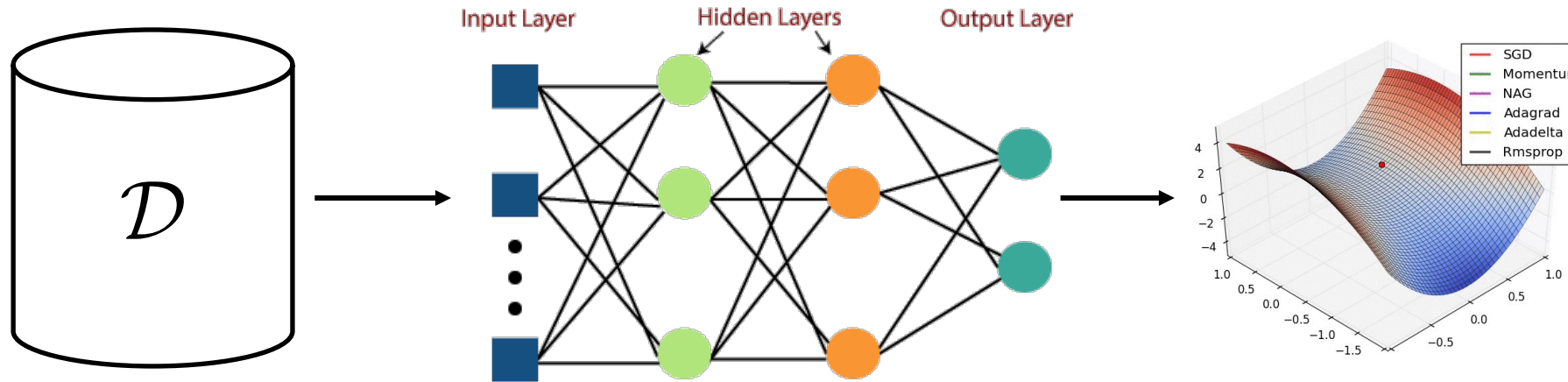
The Anatomy of Model-Based Reinforcement Learning

↓

Model based RL v0 → random shooting + MPC

↓

Model based RL v1 → MPPI + MPC

↓

Model based RL v2 → uncertainty based models

↓

Model based RL v3 → policy optimization with models

↓

Model based RL v4 → latent space models with images

Maximum likelihood supervised Learning

$$\max_{\theta} \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[ \log \hat{p}_{\theta}(s'|s,a) \right]$$

Random Search

$$\max_{\theta} \mathbb{E}_{(s,a,s')\sim\mathcal{D}} \left[\log \hat{p}_{\theta}(s'|s,a)\right]$$

Fit 1-step models



Trick: Model Residual's (s' –s)

Choice of $\hat{p}_{\theta}$ distribution determines the loss function:
1. Gaussian → $L_2$
2. Energy Based Model → Contrastive Divergence
3. Diffusion Model → Score Matching

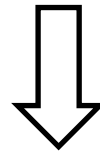More expressive may be better, at the risk of overfitting

# Model Based RL – Naïve Algorithm (Planning)

Planning

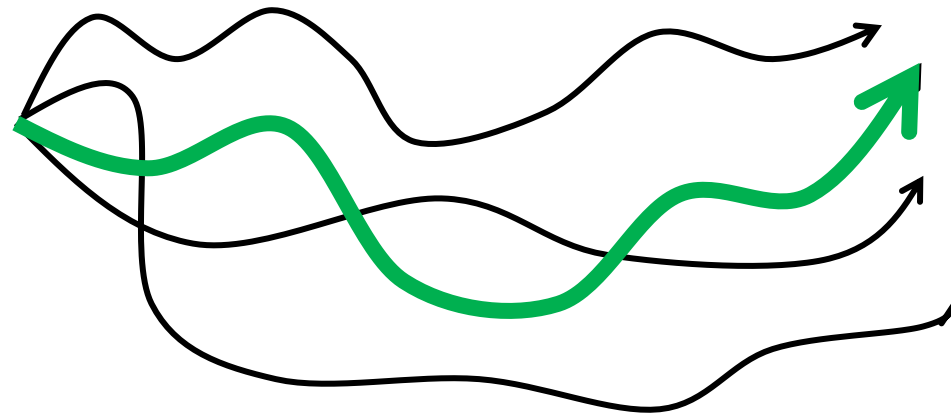$$\max_{a_0, a_1, \ldots, a_T} \sum_{t=0}^{T} r(\hat{s}_t, a_t)$$

$$\hat{s}_{t+1} \sim \hat{p}_\theta(s_{t+1} | \hat{s}_t, a_t)$$

$$\hat{s}_1 \sim \hat{p}_\theta(s_{t+1} | s_0, a_0)$$

Just do random search!

$$\arg \max_{a_0^j, a_1^j, \ldots, a_T^j} \sum_{t=0}^{T} r(\hat{s}_t^j, a_t^j)$$

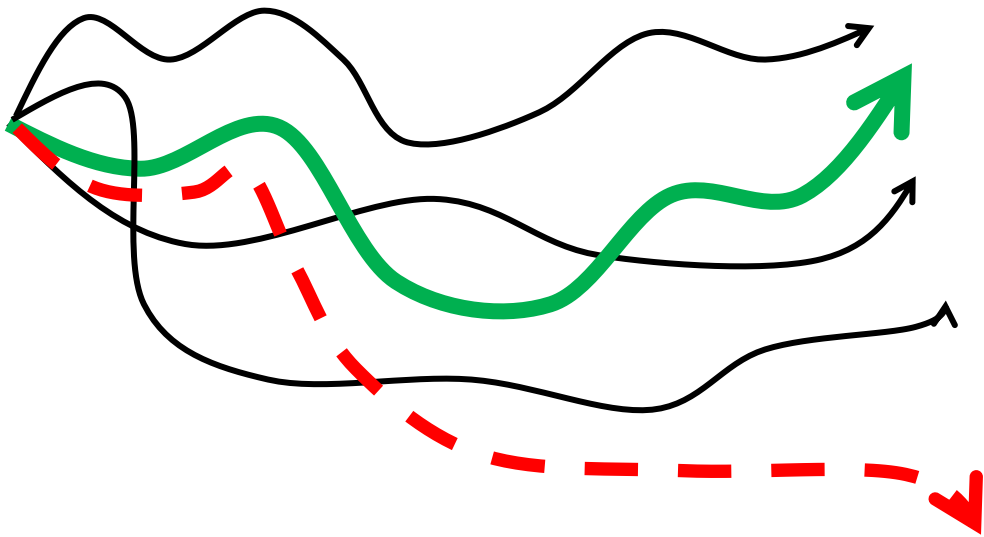$$\hat{s}_{t+1}^j \sim \hat{p}_\theta(. | \hat{s}_t^j, a_t^j)$$
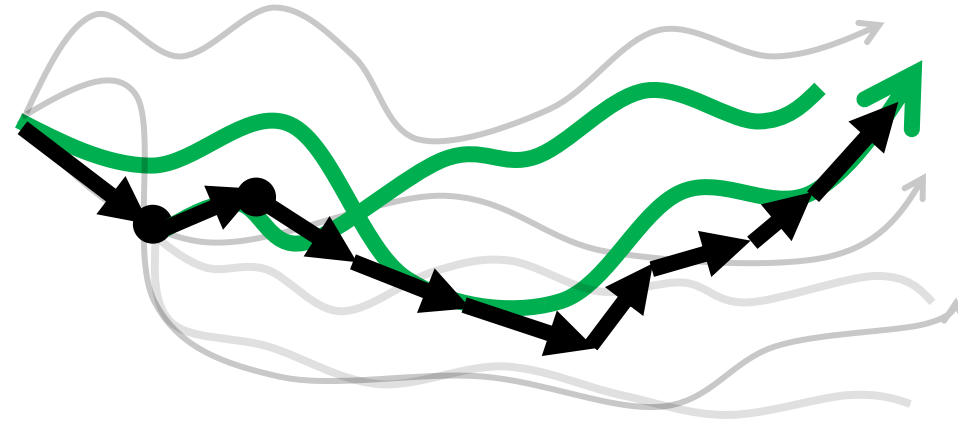
Just execute actions open loop!

Can soften by taking softmax rather than argmax

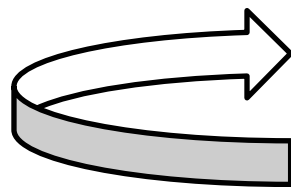# Model Based RL – Naïve Algorithm (MPC)

Without feedback, an open loop controller can diverge even for minimal noise
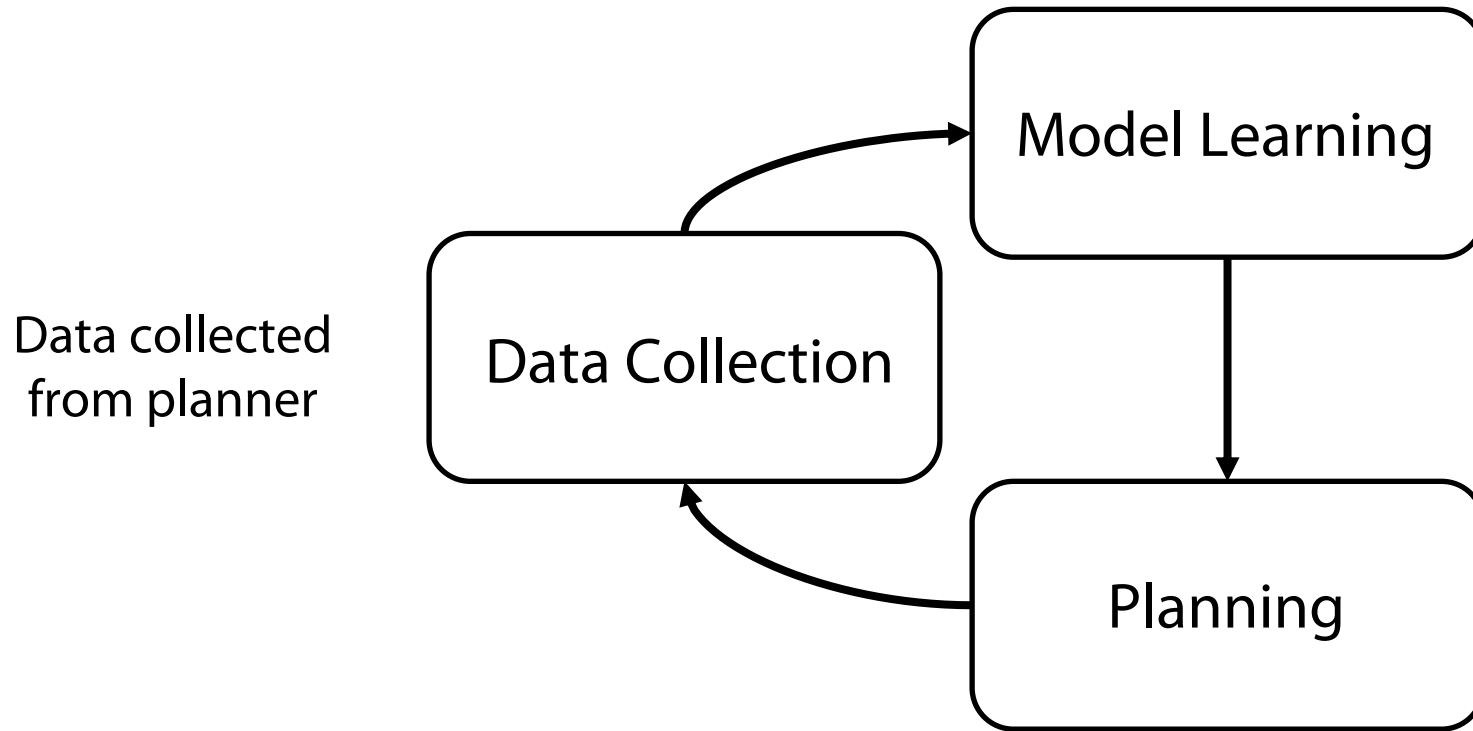
Replanning can help with divergence

Model-Predictive/Receding Horizon Control

1. Plan with random shooting from $s_t$
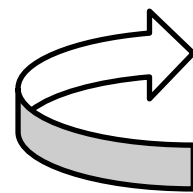2. Execute the first action $a_0$ and reach $s_{t+1}$

# Model Based RL – Naïve Algorithm (v0)



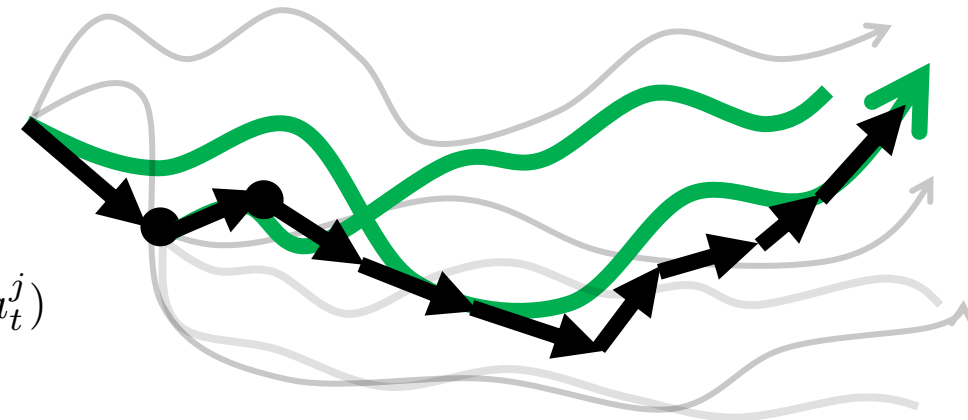Maximum likelihood supervised Learning

$$\max_\theta \mathbb{E}_{(s,a,s')\sim\mathcal{D}}\left[\log\hat{p}_\theta(s'|s,a)\right]$$

Data collected from planner

Planning with Shooting + MPC

Better than open loop planning because of feedback
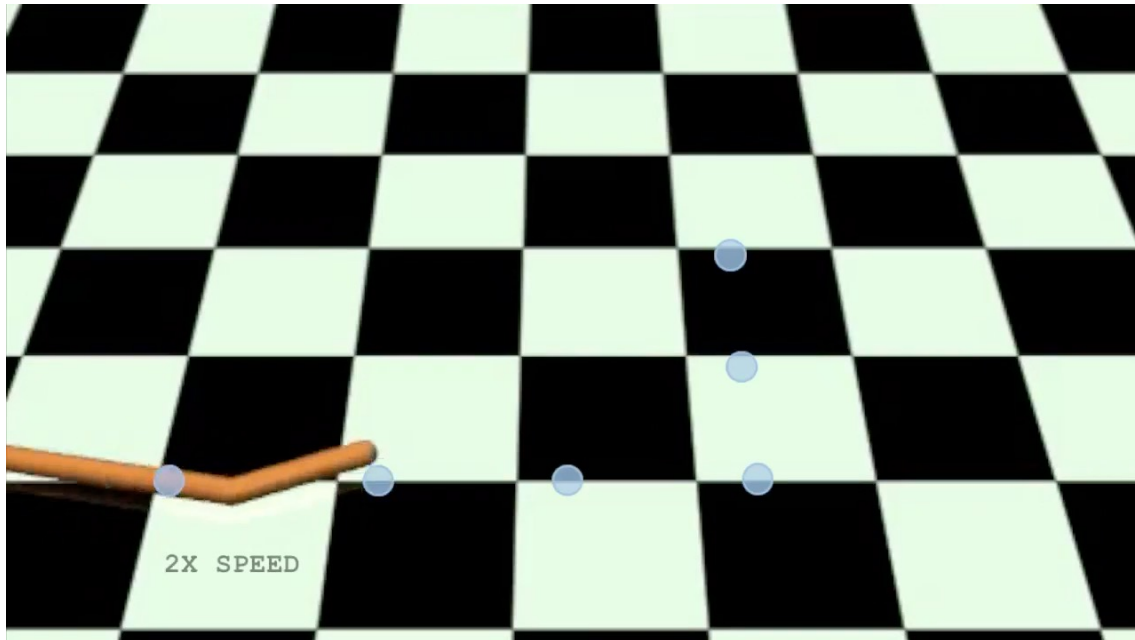
$$\arg\max_{a_0^j,a_1^j,...,a_T^j}\sum_{t=0}^{T}r(\hat{s}_t^j,a_t^j)$$

$$\hat{s}_{t+1}^j \sim \hat{p}_\theta(.|\hat{s}_t^j,a_t^j)$$

# Does it work?



Just 20 minutes of training time with random data!

# Does it work?



Cheetah

Significant gap from MFRL

Cumulative Reward

- Mb
- Mf
- Mb-Mf (ours)

# Lecture outline

The Anatomy of Model-Based Reinforcement Learning

↓

Model based RL v0 → random shooting + MPC

↓

Model based RL v1 → MPPI + MPC
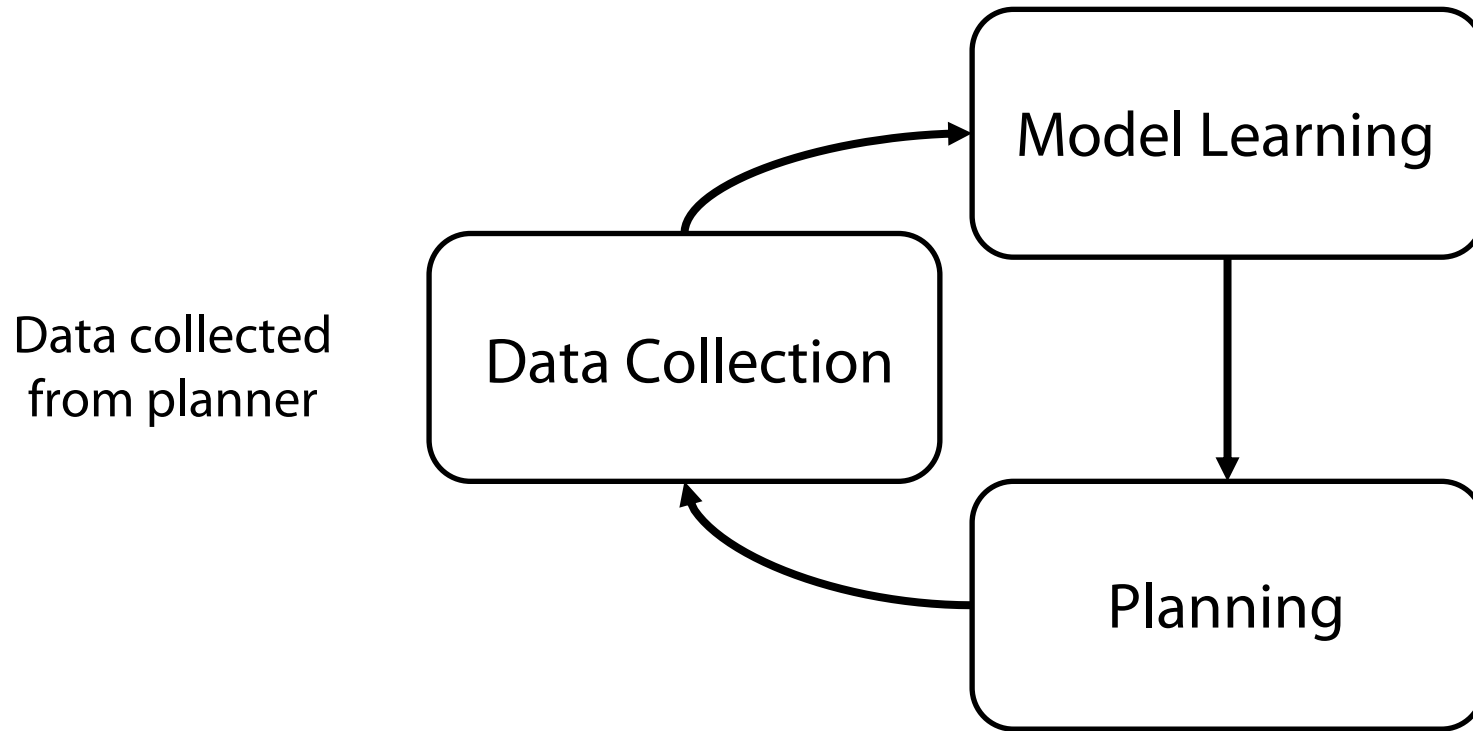
↓

Model based RL v2 → uncertainty based models

↓

Model based RL v3 → policy optimization with models

↓

Model based RL v4 → latent space models with images
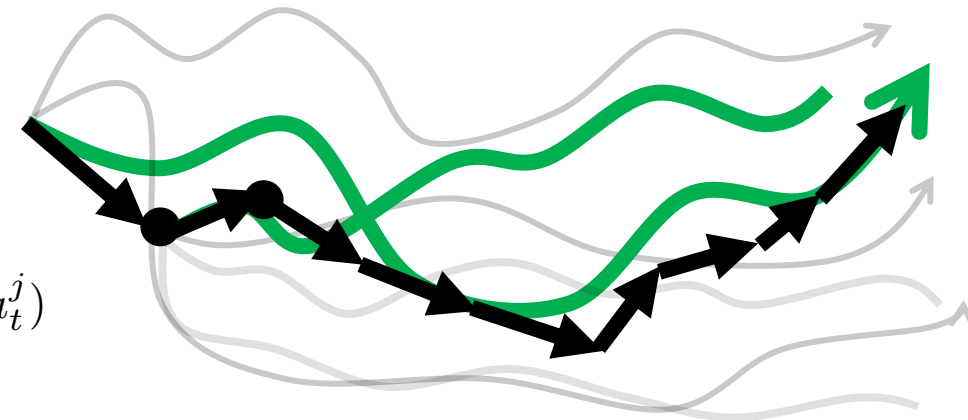
# What might be the issue?



Model Learning

Data Collection

Planning

Data collected from planner

Maximum likelihood supervised Learning

$$\max_\theta \mathbb{E}_{(s,a,s')\sim\mathcal{D}}\left[\log \hat{p}_\theta(s'|s,a)\right]$$

Planning with Shooting + MPC

**Searching for a needle in a haystack by random shooting, high variance!**

$$\arg\max_{a_0^j,a_1^j,\ldots,a_T^j}\sum_{t=0}^{T} r(\hat{s}_t^j, a_t^j)$$

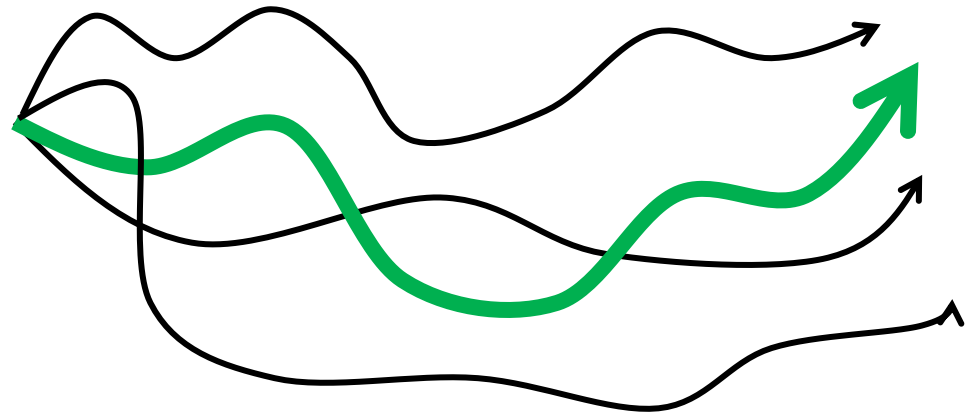$$\hat{s}_{t+1}^j \sim \hat{p}_\theta(.|\hat{s}_t^j, a_t^j)$$

# Better Sampling Techniques for Shooting

Sampled from stationary uniform/gaussian distribution

Can we inform the sampling function with the reward function?

$$\arg \max_{a_0^j, a_1^j, \ldots, a_T^j} \sum_{t=0}^{T} r(\hat{s}_t^j, a_t^j)$$

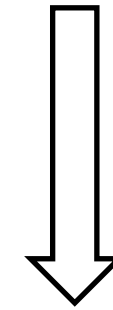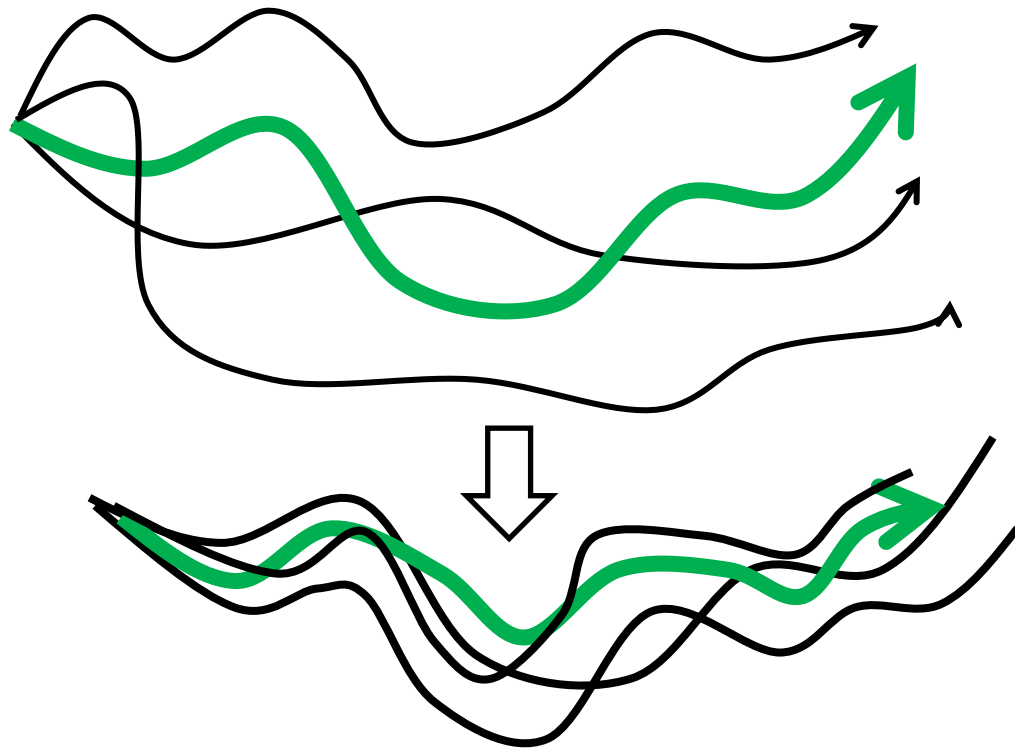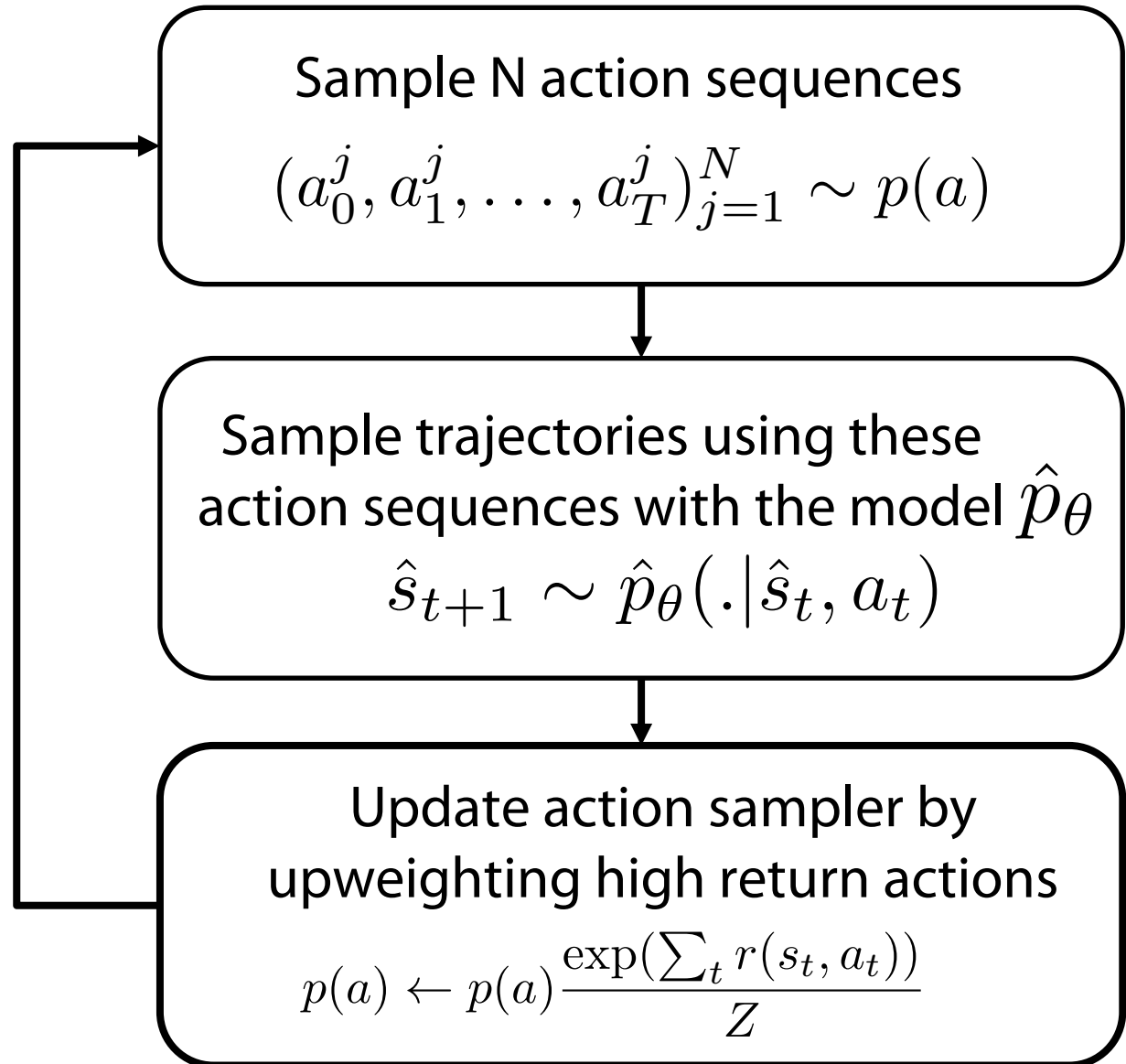$$\hat{s}_{t+1}^j \sim \hat{p}_\theta(.|\hat{s}_t^j, a_t^j)$$

Idea: Iteratively upweight sampling distribution around the things that are higher returns
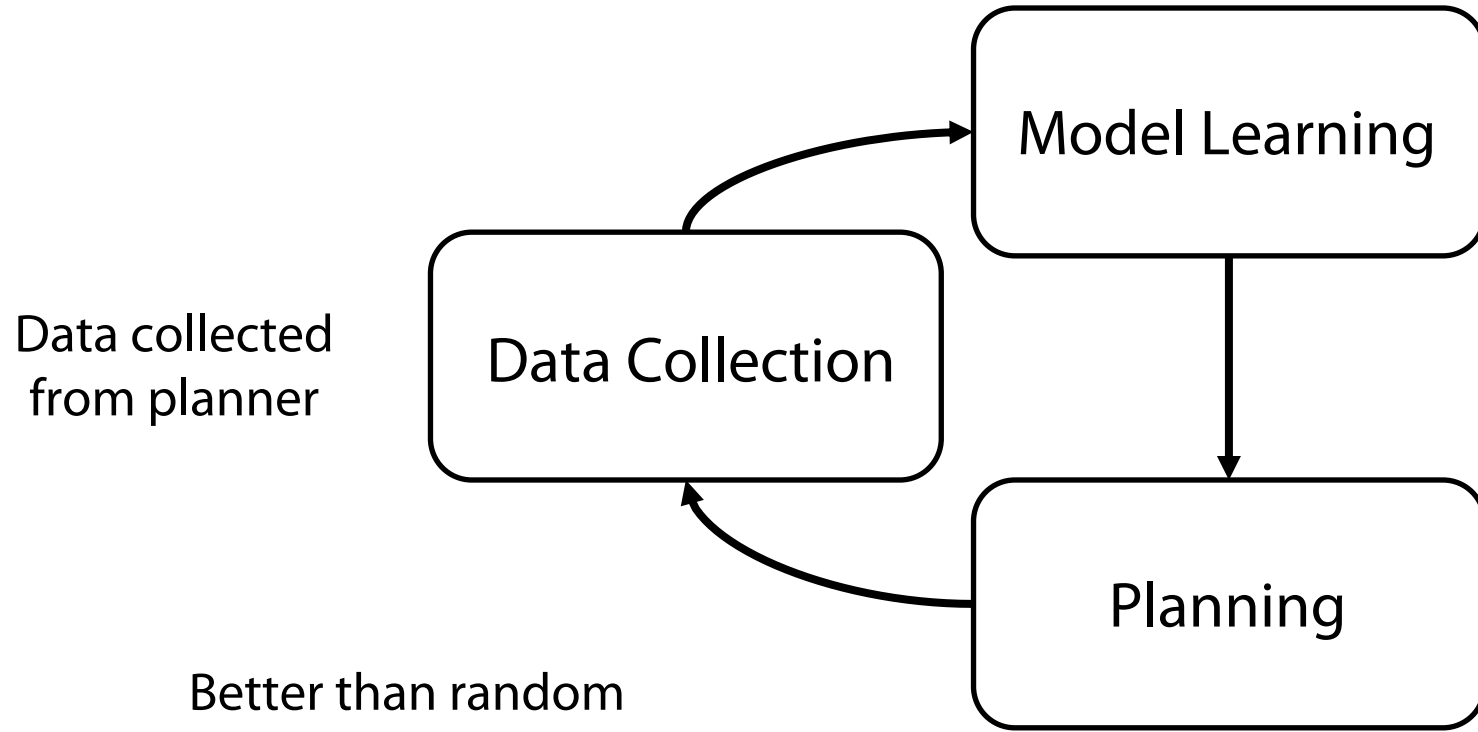
# Better Sampling Techniques for Shooting - MPPI

Idea: Iteratively upweight sampling distribution around the things that are higher returns



Referred to as **MPPI**, lower variance!

Sample N action sequences
$$(a_0^j, a_1^j, \ldots, a_T^j)_{j=1}^N \sim p(a)$$

Sample trajectories using these action sequences with the model $\hat{p}_\theta$
$$\hat{s}_{t+1} \sim \hat{p}_\theta(.|\hat{s}_t, a_t)$$

Update action sampler by upweighting high return actions
$$p(a) \leftarrow p(a) \frac{\exp(\sum_t r(s_t, a_t))}{Z}$$

# Model Based RL – Better Sampling Methods (v1)



**Model Learning**

**Data Collection**

**Planning**

Data collected from planner

Better than random shooting + MPC, since lower variance!

Aside: Can derive this update trying to bring sampling distribution close to optimal distribution

Maximum likelihood supervised Learning

$$\max_{\theta} \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[ \log \hat{p}_{\theta}(s'|s,a) \right]$$

Planning with MPPI + MPC

$$\arg \max_{a_0^j, a_1^j, \ldots, a_T^j} \sum_{t=0}^{T} r(\hat{s}_t^j, a_t^j)$$

$$\hat{s}_{t+1}^j \sim \hat{p}_{\theta}(.|\hat{s}_t^j, a_t^j)$$

$$p(a) \leftarrow p(a) \frac{\exp(\sum_t r(s_t, a_t))}{Z}$$

# Does it work?



Testing lap time: 9.45 s

Trajectory Rollouts

# Lecture outline

The Anatomy of Model-Based Reinforcement Learning

$\downarrow$

Model based RL v0 → random shooting + MPC

$\downarrow$

Model based RL v1 → MPPI + MPC

$\downarrow$

Model based RL v2 → uncertainty based models

$\downarrow$

Model based RL v3 → policy optimization with models

$\downarrow$

Model based RL v4 → latent space models with images

# What might be the issue?

Rollouts under learned model != Rollouts under true model

Model bias/compounding error

True Rollout

Why does this happen? → lack of data

1. Errors in state go to OOD next states
2. Deviations in actions go to OOD next states

Predicted Rollout Under Model

Model is bad on OOD states!

Most trained deep models can only roll out for 5-10 steps maximum!

Idea 1: Change the training objective of the model to directly account for this!

## Equation error – 1 step prediction error

$$\max_\theta \mathbb{E}_{(s,a,s')\sim\mathcal{D}}\left[\log \hat{p}_\theta(s'|s,a)\right]$$

Model error under learned mode $\hat{p}_\theta$ rather under true model

Can be a challenging non-convex optimization!

## Simulation error – K step prediction error
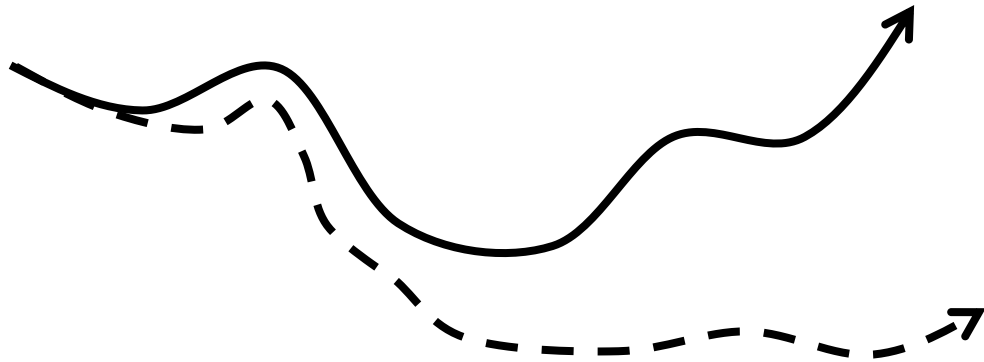
$$\max_\theta \sum_t \log \hat{p}_\theta(s_{t+1}|\hat{s}_t, a_t)$$

$$\hat{s}_t \sim \hat{p}_\theta(.|\hat{s}_{t-1}, a_{t-1})$$

# How might we deal with compounding error?

Idea 2: Estimate when OOD and account for it

Measure uncertainty!

Maximum likelihood models

Uncertainty-aware models



Being aware of uncertainty allows us to account for the effects of model bias!
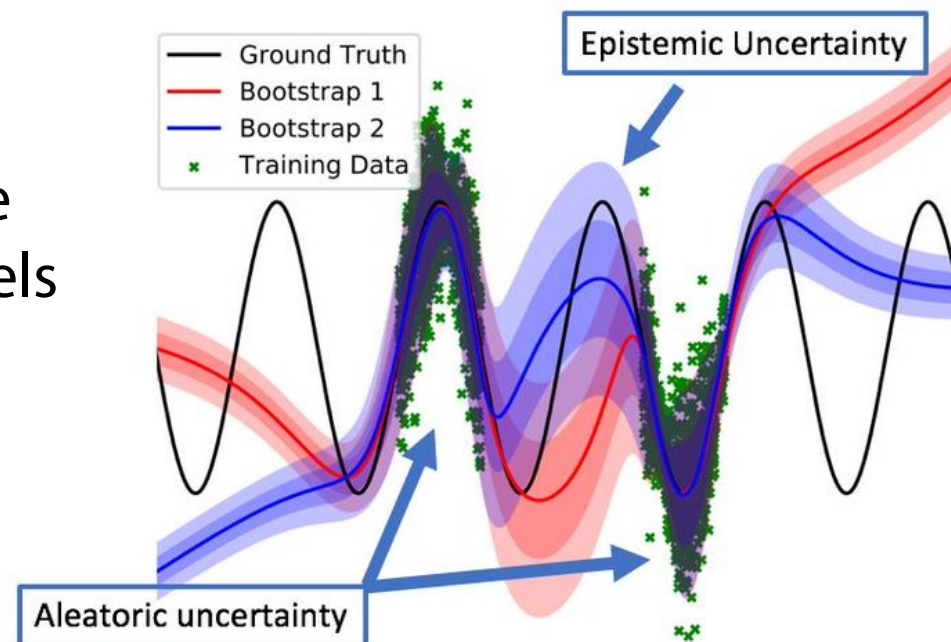
# What is uncertainty?

**Alleatoric Uncertainty**

(environment stochasticity)

Easier, can use
stochastic models



**Epistemic Uncertainty**

(Lack of data)

More challenging, need
to compute posterior

Let's largely focus on epistemic uncertainty

# How might we measure uncertainty?

$$p(\theta|\mathcal{D})$$

**Difficult to estimate directly!**

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{\int p(\mathcal{D}|\theta')p(\theta')d\theta'}$$

1. Bayesian neural networks
2. Ensemble methods
3. …

Directly model posterior distribution

Use variational inference to avoid computing partition function

$$\min_{q(\theta|\mathcal{D})} D_{KL}(q(\theta|\mathcal{D}) \;||\; p(\theta|\mathcal{D}))$$

Challenge: can be difficult to express rich distributions



output:

weights:
(with distribution)

hidden layer:

weights:
(with distribution)

input:

y

X

$$p(\theta|\mathcal{D})$$

Difficult to estimate directly!

Learn an ensemble of models

1. Bayesian neural networks
2. Ensemble methods
3. …



Low data regime → high ensemble variance

Approximate posterior

Easier and more expressive than BNNs!

Learn ensembles of dynamics models with MLE rather than a single model



$$\max_\theta \mathbb{E}_{(s,a,s')\sim\mathcal{D}} \left[\log \hat{p}_\theta(s'|s,a)\right]$$

$$\max_\theta \mathbb{E}_{(s,a,s')\sim\mathcal{D}} \left[\log \hat{p}_\theta(s'|s,a)\right]$$

$$\max_\theta \mathbb{E}_{(s,a,s')\sim\mathcal{D}} \left[\log \hat{p}_\theta(s'|s,a)\right]$$

Learn ensembles by either subsampling the data or having different initializations

Take expected value under the uncertain dynamics



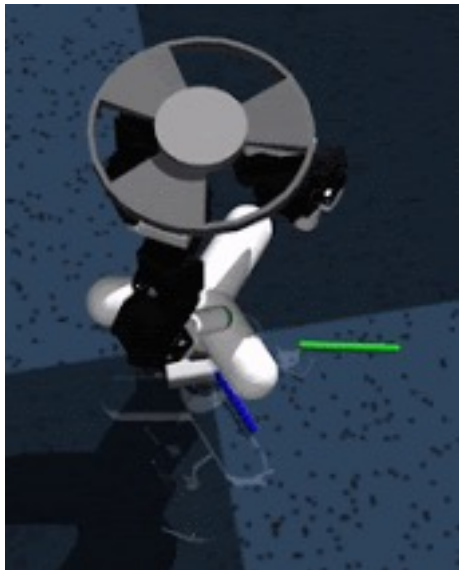Low uncertainty

High uncertainty

Expected value over ensemble

$$\arg \max_{(a_0^j, a_1^j, \dots, a_T^j)_{j=1}^{N}} \sum_{i=1}^{K} \sum_{t=0}^{T} r((\hat{s}_t^j)^i, a_t^j)$$

$$(\hat{s}_{t+1}^j)^i \sim \hat{p}_{\theta_i}(.|(\hat{s}_t^j)^i, a_t^j)$$

Can also swap which ensemble element is propagated at every step or just pick randomly amongst them

Avoids overly OOD settings since the expected reward is affected by uncertainty

Take **pessimistic** value under the uncertain dynamics

Low uncertainty

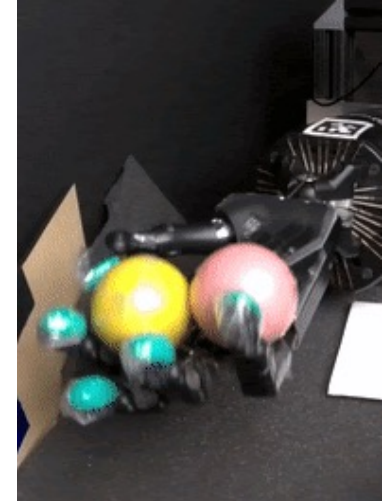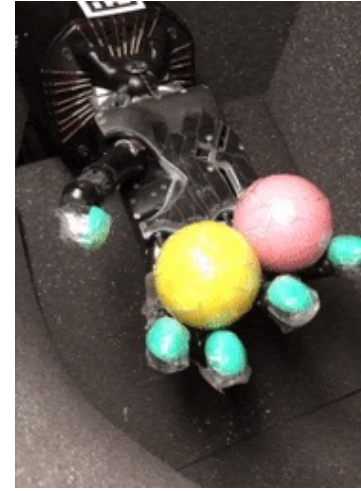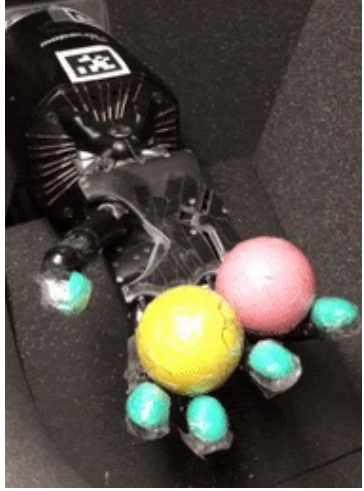Penalize ensemble variance

High uncertainty

$$\arg \max_{(a_0^j, a_1^j, \ldots, a_T^j)_{j=1}^N} \sum_{i=1}^{K} \sum_{t=0}^{T} r((\hat{s}_t^j)^i, a_t^j) - \lambda \mathrm{Var}((\hat{s}_t^j)^i)$$

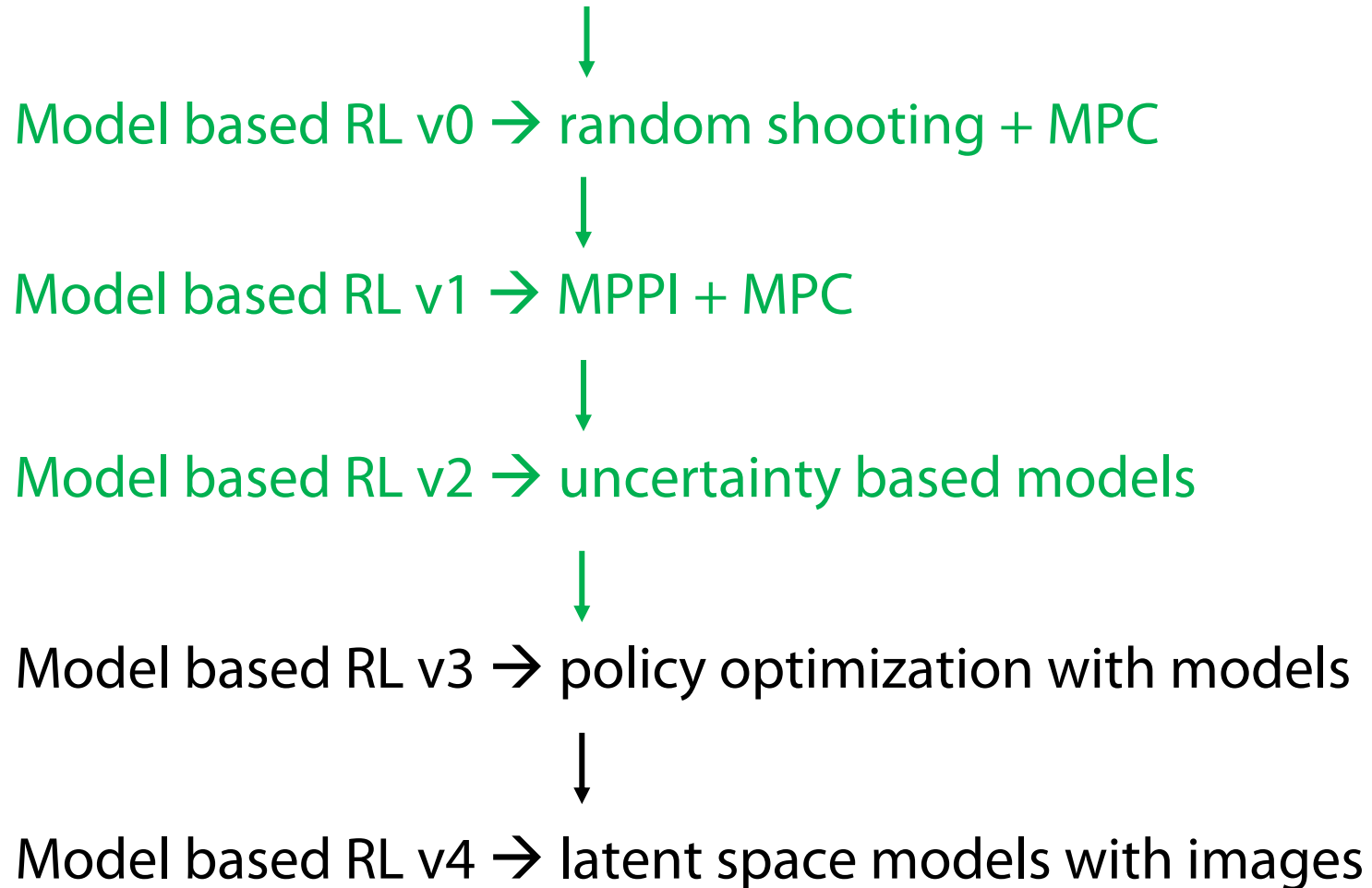$$(\hat{s}_{t+1}^j)^i \sim \hat{p}_{\theta_i}(. | (\hat{s}_t^j)^i, a_t^j)$$

Avoids overly OOD settings since these states are explicitly penalized

# Does this work?

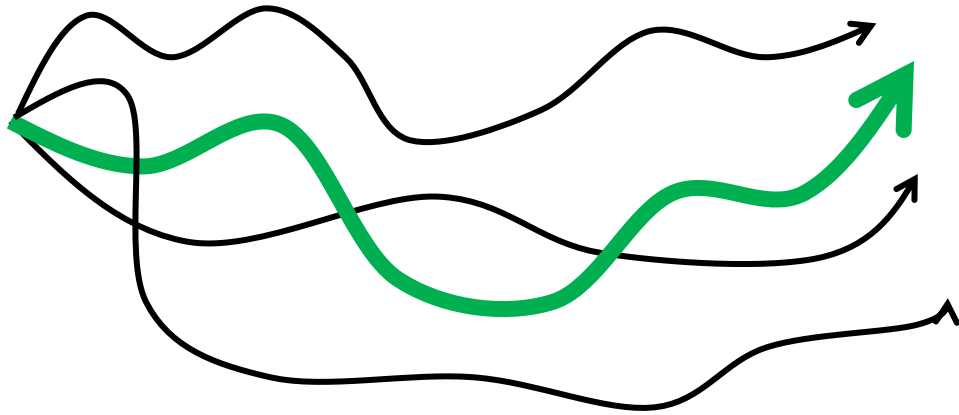# Lecture outline
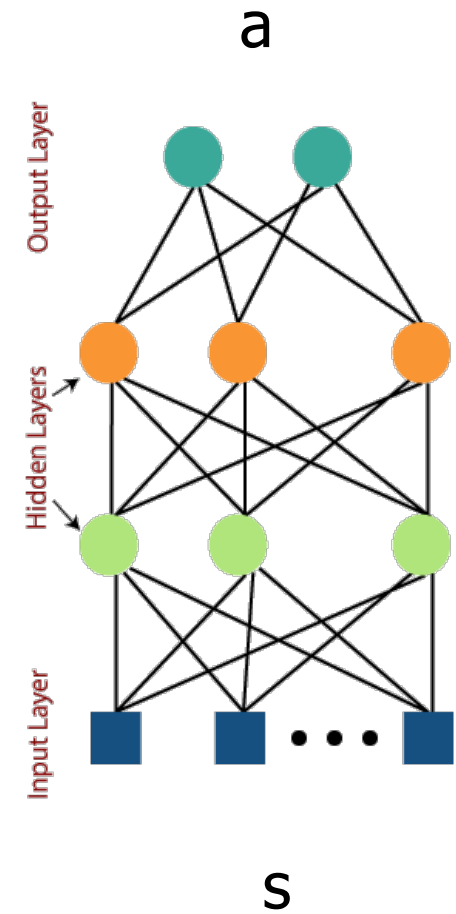
The Anatomy of Model-Based Reinforcement Learning

$\downarrow$

Model based RL v0 → random shooting + MPC

$\downarrow$

Model based RL v1 → MPPI + MPC

$\downarrow$

Model based RL v2 → uncertainty based models

$\downarrow$

Model based RL v3 → policy optimization with models

$\downarrow$

Model based RL v4 → latent space models with images

# What might be the issue?

Huge number of samples
needed to reduce variance

Amortize planning
into a policy

a

Output Layer

Hidden Layers

Input Layer

s
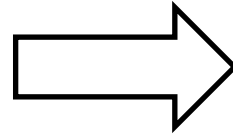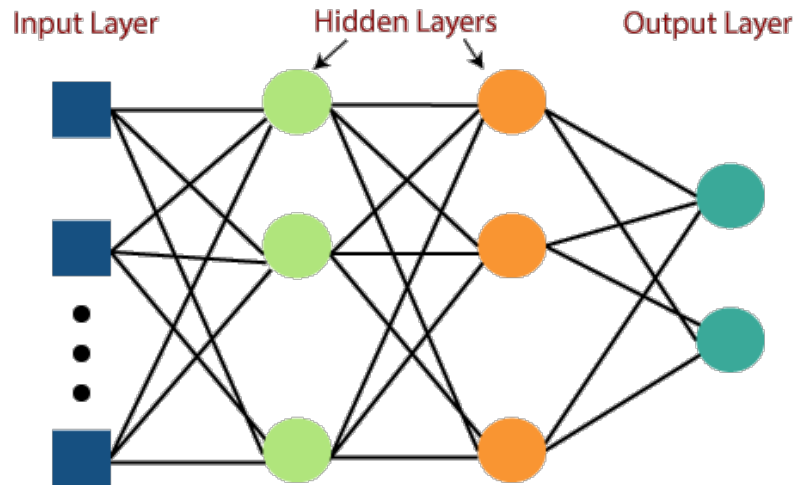
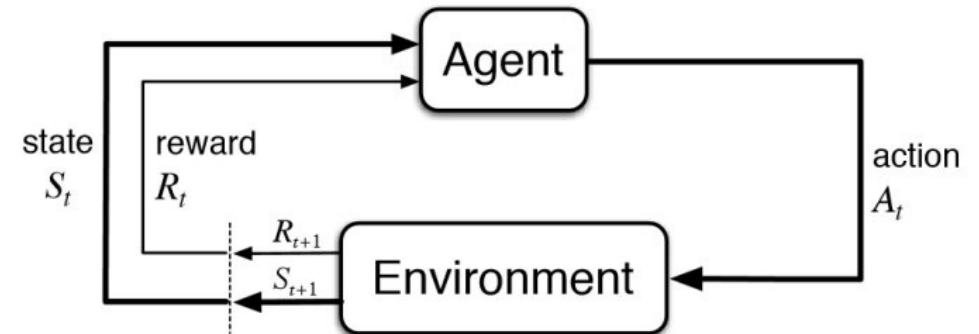Extremely slow, hard to run in real time

# Speeding Up Model-Based Planning

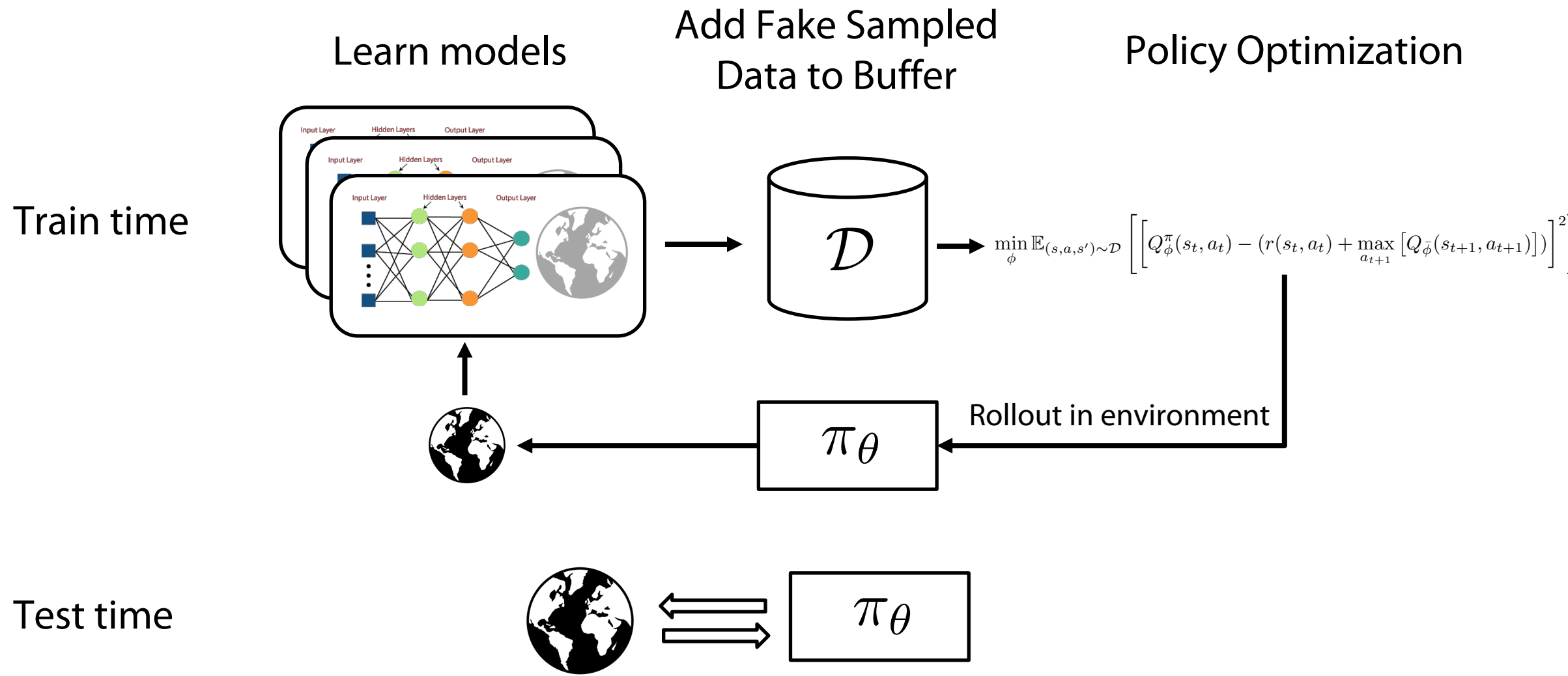$$\max_{\theta} \mathbb{E}_{(s,a,s')\sim\mathcal{D}} \left[\log \hat{p}_\theta(s'|s,a)\right]$$

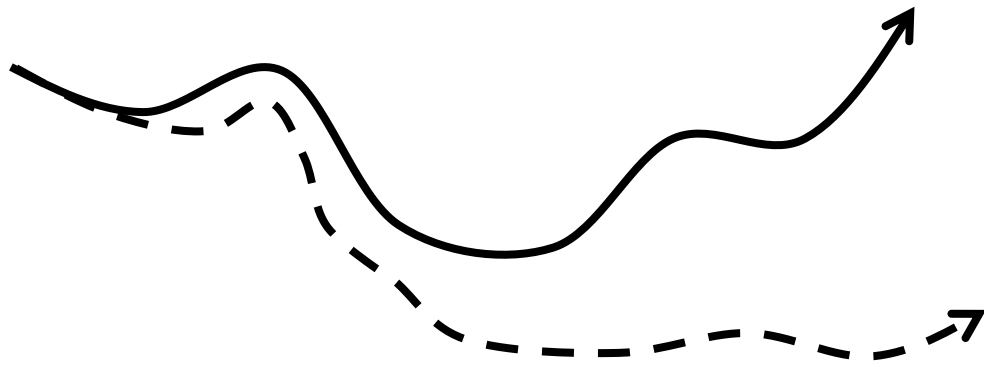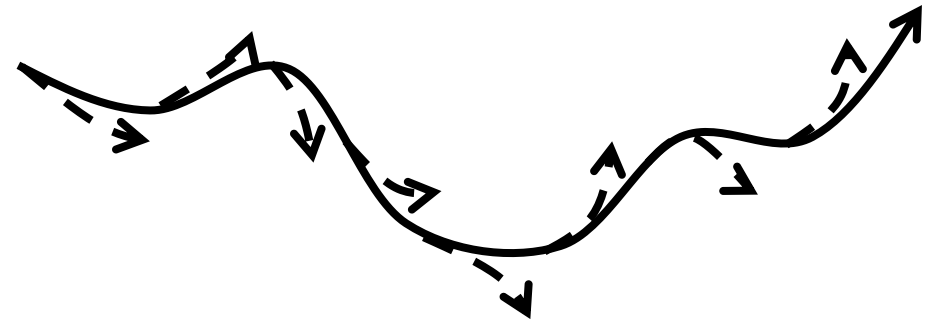Use model(s) to generate data for policy optimization



Input Layer    Hidden Layers    Output Layer

state $S_t$   reward $R_t$   Agent   action $A_t$

$R_{t+1}$

$S_{t+1}$

Environment

Can use PG or off-policy!

# Generating Data for Policy Optimization
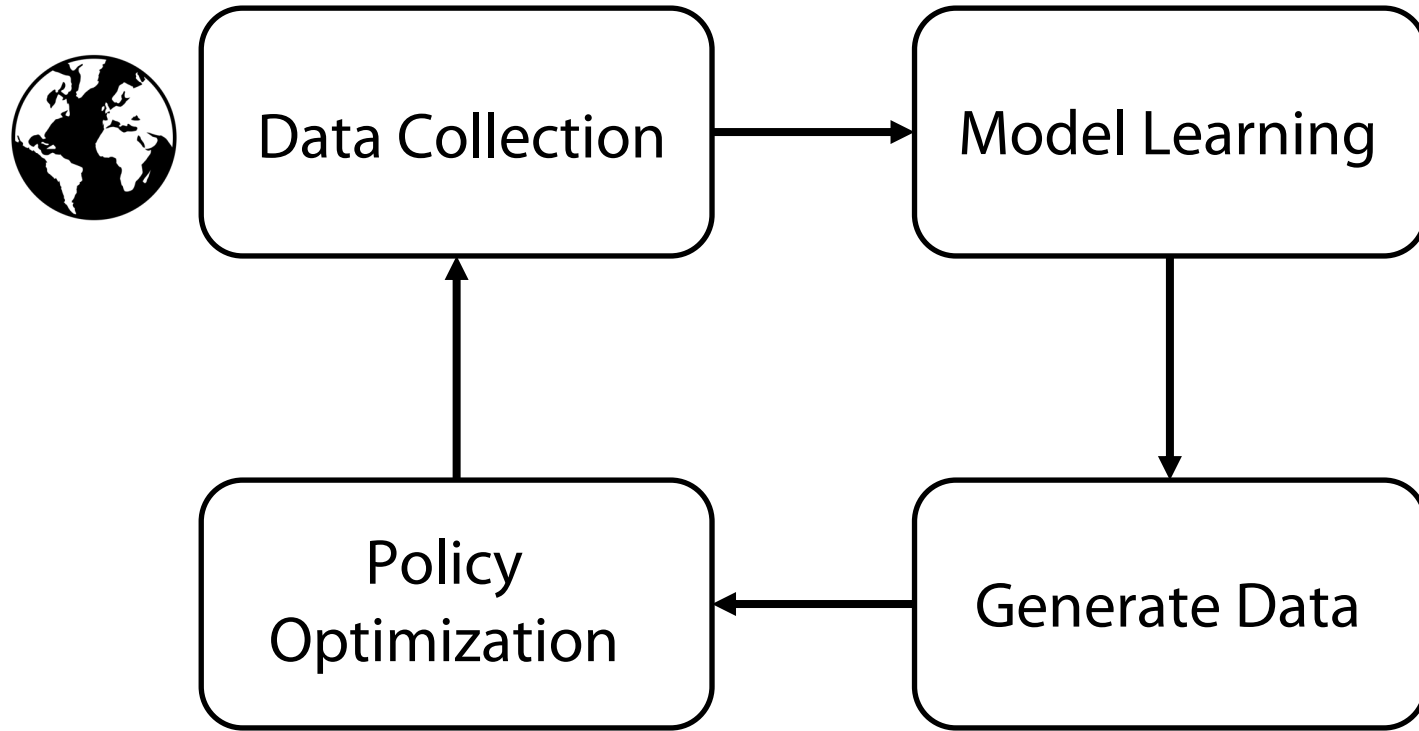
Long horizon rollouts can deviate

Short horizon rollouts deviate far less

Balance between off-policy coverage and compounding error
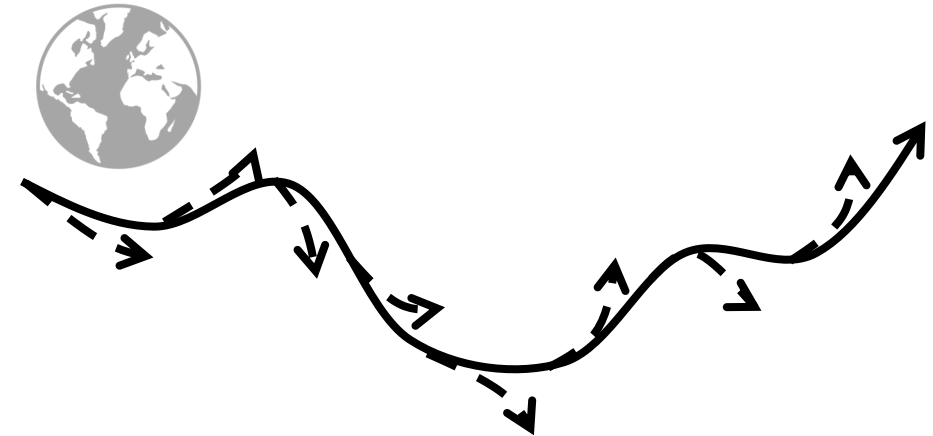
More in the readings!

# Model Based RL – Using Models for Policy Optimization (v3)

Data Collection → Model Learning

Policy Optimization ← Generate Data

Maximum likelihood supervised Learning

$$\max_\theta \mathbb{E}_{(s,a,s')\sim\mathcal{D}} \left[\log \hat{p}_\theta(s'|s,a)\right]$$

$$\min_\phi \mathbb{E}_{(s,a,s')\sim\mathcal{D}} \left[\left[Q_\phi^\pi(s_t,a_t) - (r(s_t,a_t) + \max_{a_{t+1}}\left[Q_{\bar{\phi}}(s_{t+1},a_{t+1})\right])\right]^2\right]$$

More expensive/harder at training time, faster at test time

# Does this work?

# Lecture outline

The Anatomy of Model-Based Reinforcement Learning

↓

Model based RL v0 → random shooting + MPC

↓

Model based RL v1 → MPPI + MPC

↓

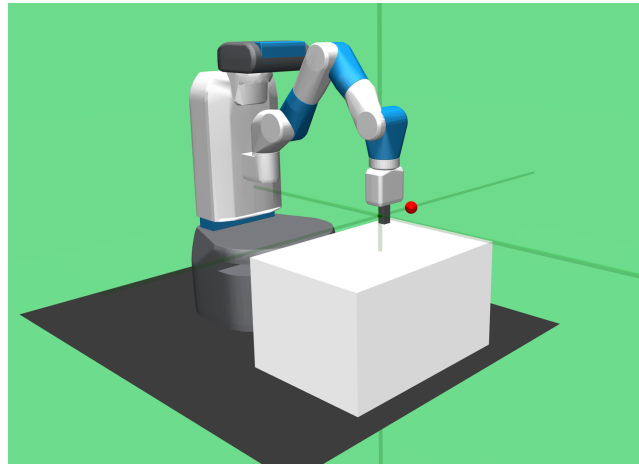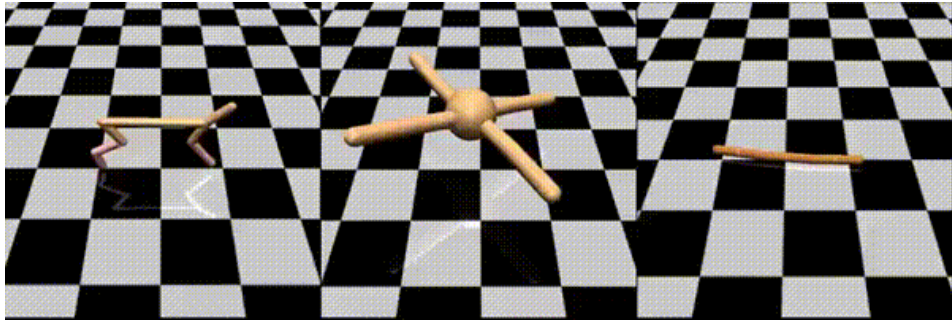Model based RL v2 → uncertainty based models

↓

Model based RL v3 → policy optimization with models

↓

Model based RL v4 → latent space models with images
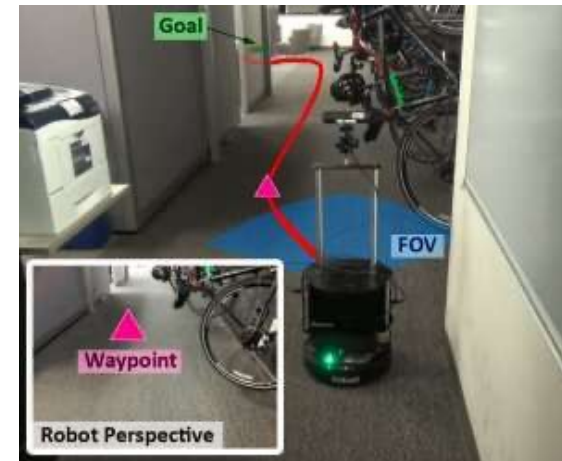
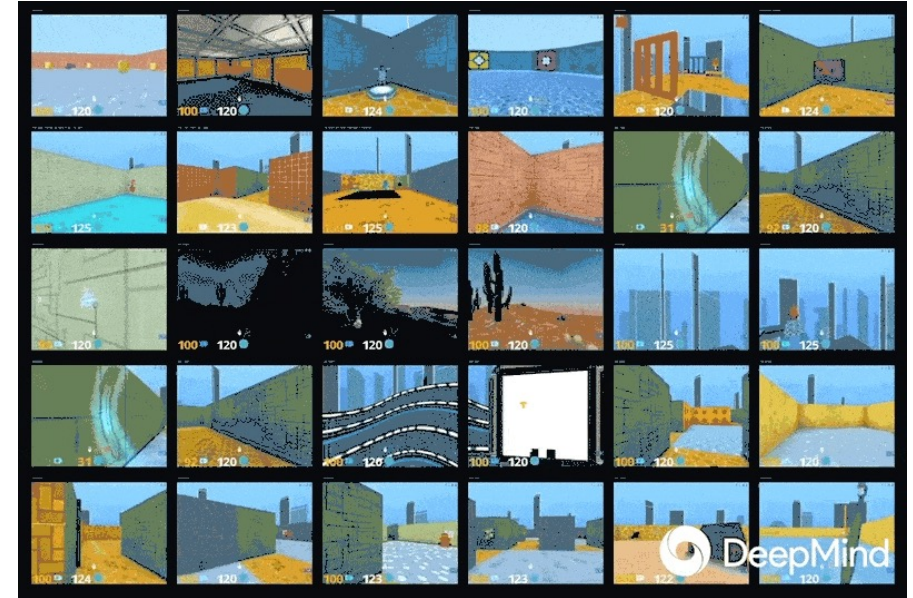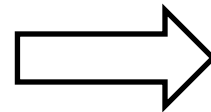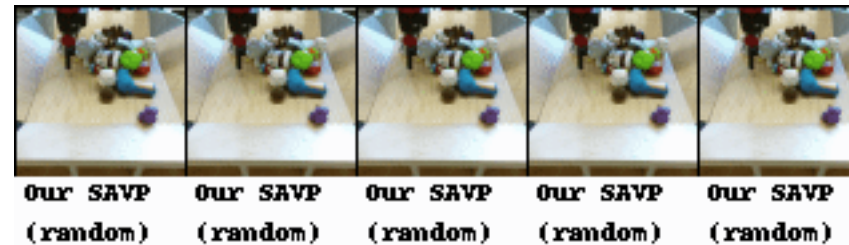# What about images?



State based domains

Image based domains

# Why is learning from images hard?

Generative modeling is videos, challenging to model multimodal correlated predictions
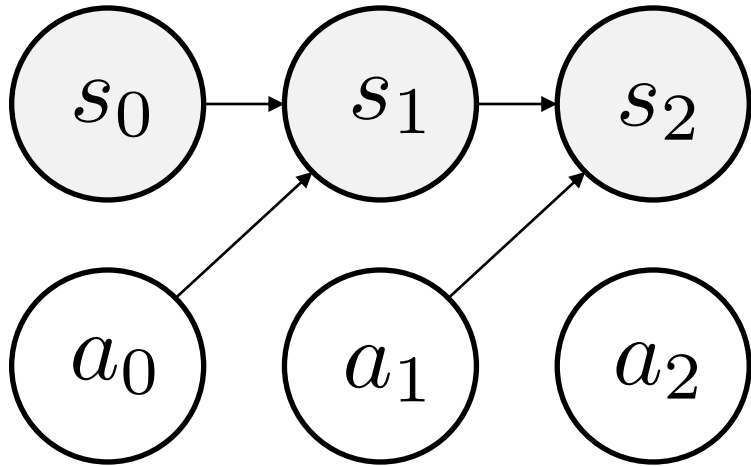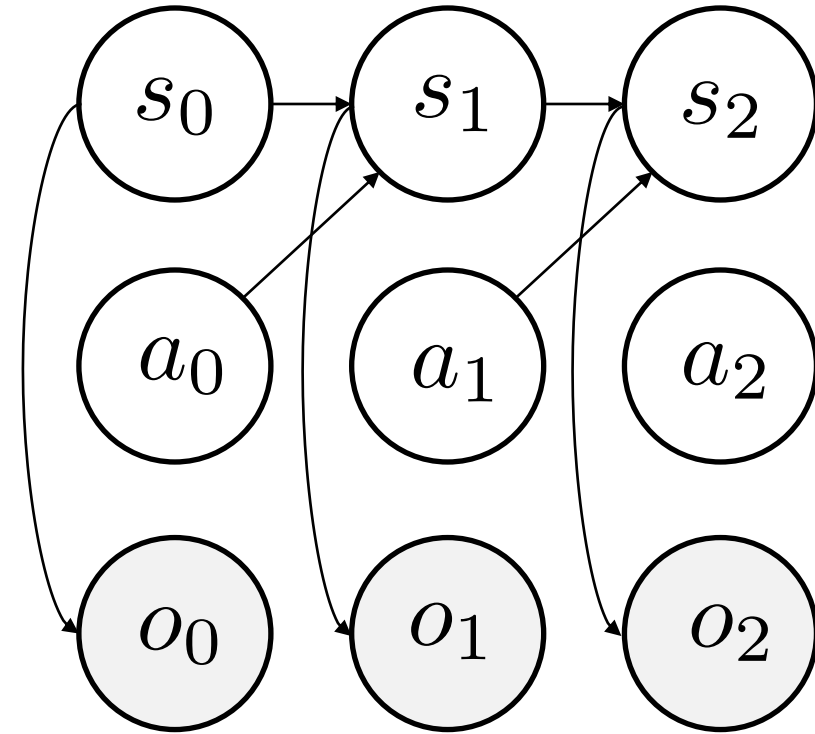


Partially observable!

Long horizon predictions in video space can be challenging!

# Model Based RL – Latent Space Models for Image Based RL (v4)
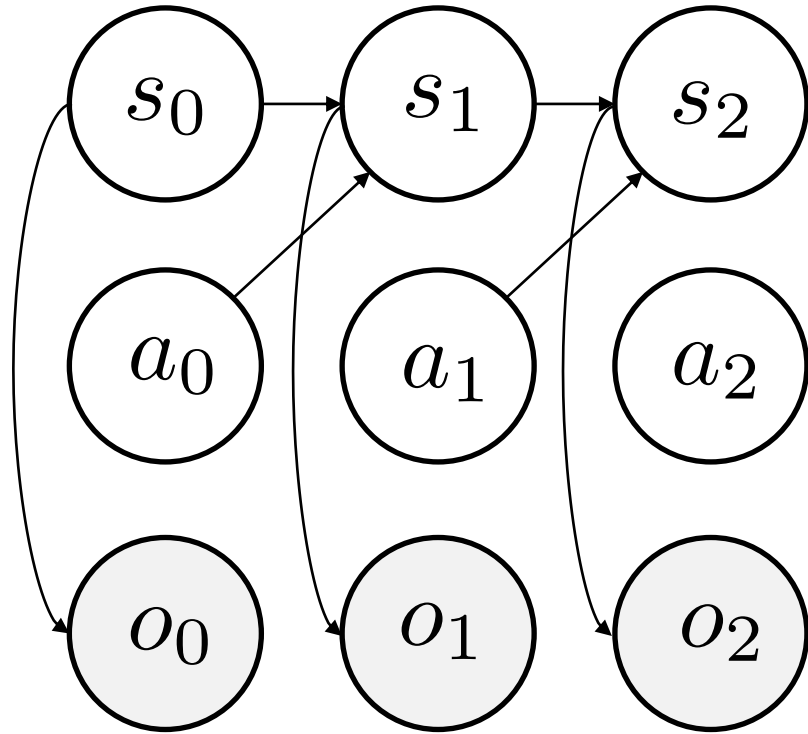
Fully observed – Markovian case

Partially observed – Non-Markovian case

If we can infer latent state and learn dynamics, then we can plan in a much smaller space

How do we infer latent state and learn dynamics in this space?

Learn latent encoder to infer latent state from observations $q_\phi(s_t|o_{1:t})$

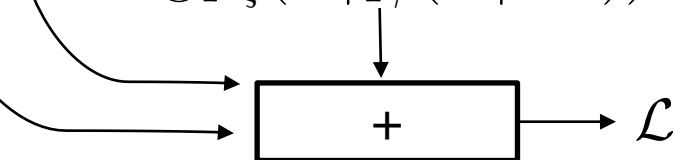Learn action conditioned latent transition model $p_\eta(s_{t+1}|s_t, a_t)$

$$\log p_\eta(q_\phi(s_{t+1}|o_{1:t+1})|q_\phi(s_t|o_{1:t}), a_t)$$

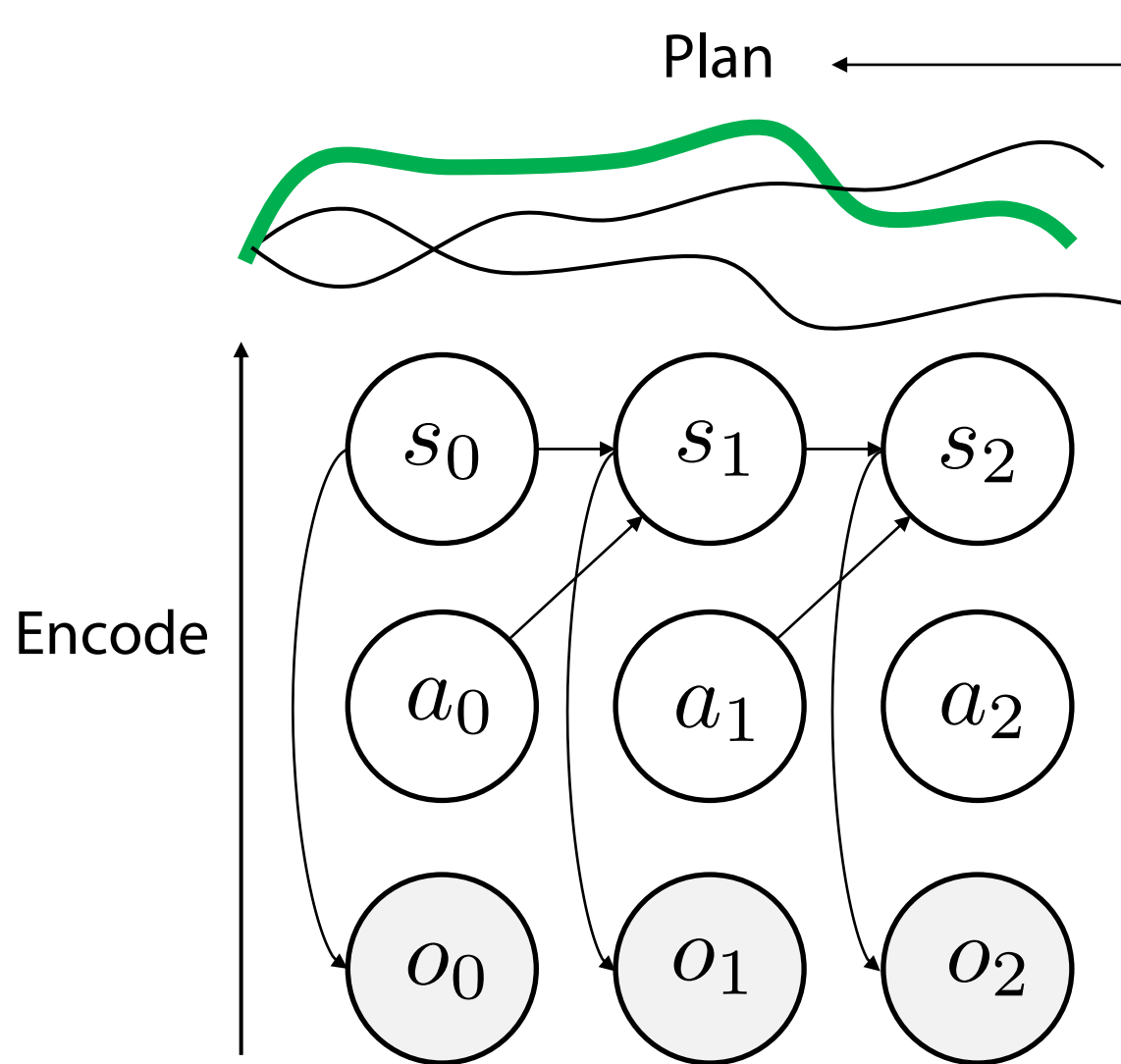Learn latent decoder to reconstruct observations $p_\psi(o_t|s_t)$

$$\log p_\psi(o_t|s_t)$$

Learn reward predictor from latent state $p_\zeta(r_t|s_t)$

$$\log p_\zeta(r_t|q_\phi(s_t|o_{1:t}))$$

Can derive the whole thing from first principles using variational inference!

# How do we **use** latent space models?



Plan

Encode

$s_0 \to s_1 \to s_2$

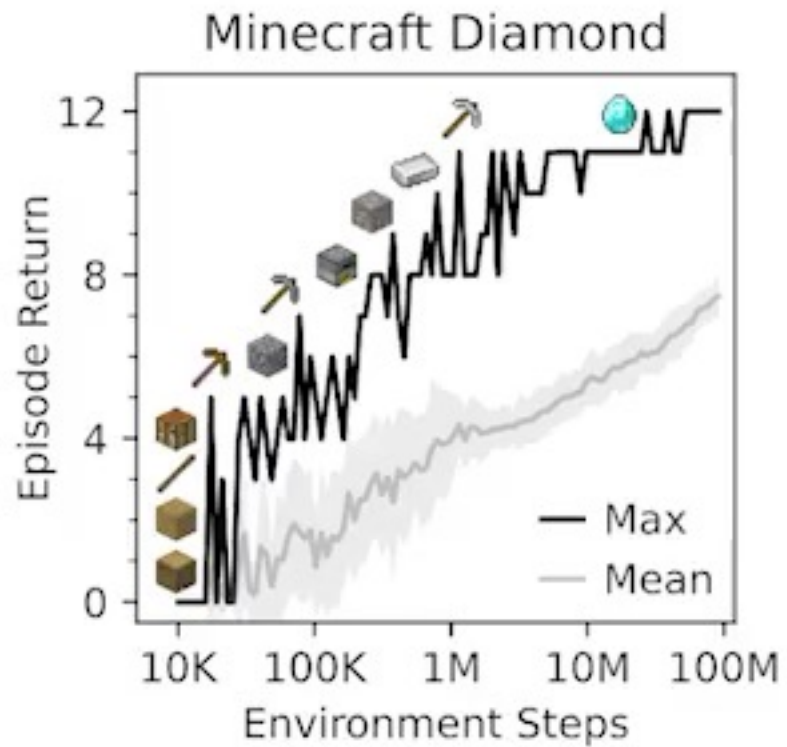$a_0 \quad a_1 \quad a_2$

$o_0 \quad o_1 \quad o_2$

Apply any of the methods from this lecture, just in latent space!

1. Avoids predicting image frames at planning time
2. Scales much better than image prediction
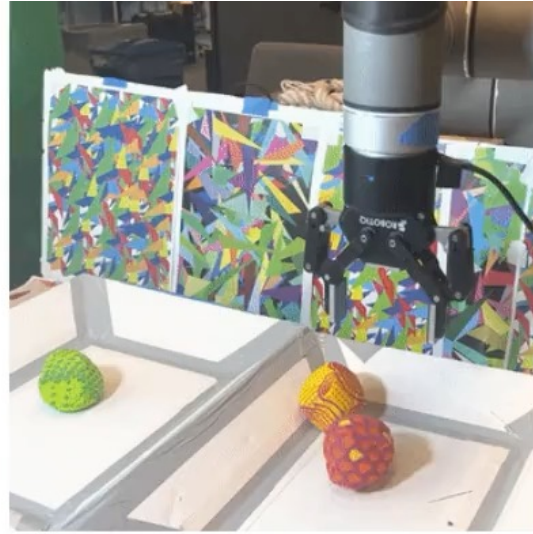3. Allows for longer horizon predictions

Minecraft Diamond
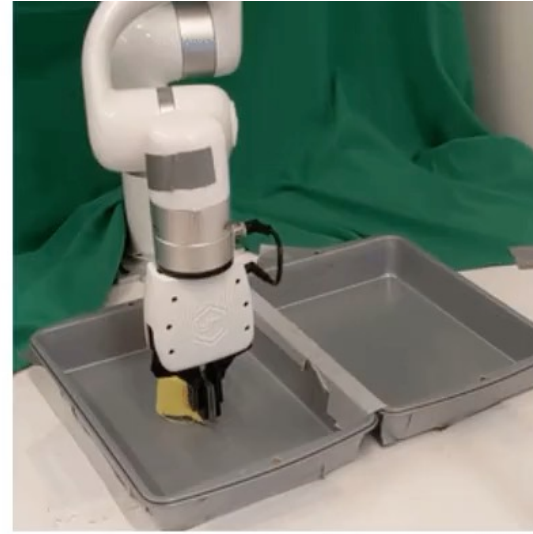
DreamerV3 First Diamond from Scratch

# Does this work?



A1 Quadruped
Walking

UR5 Multi-Object
Visual Pick Place

XArm Visual Pick
and Place

Sphero Ollie Visual
Navigation

Training from images in < 1 hour!

# Lecture outline

The Anatomy of Model-Based Reinforcement Learning

↓

Model based RL v0 → random shooting + MPC

↓

Model based RL v1 → MPPI + MPC

↓

Model based RL v2 → uncertainty based models

↓

Model based RL v3 → policy optimization with models

↓

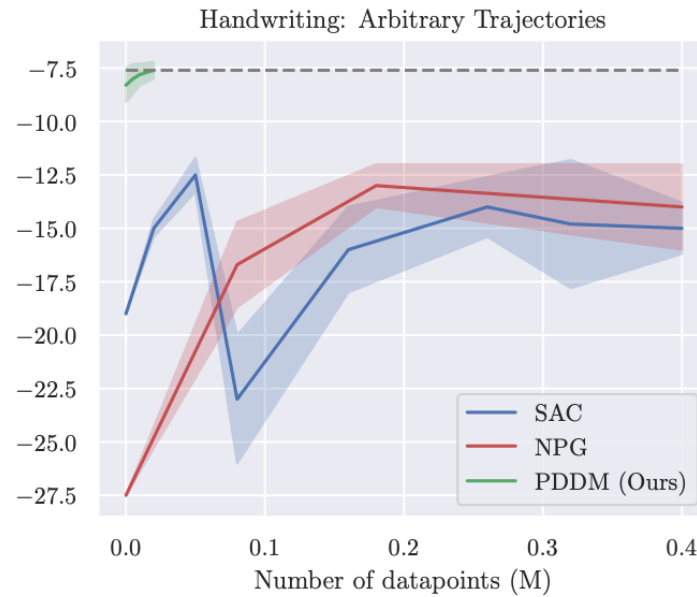Model based RL v4 → latent space models with images

# Why should you care?

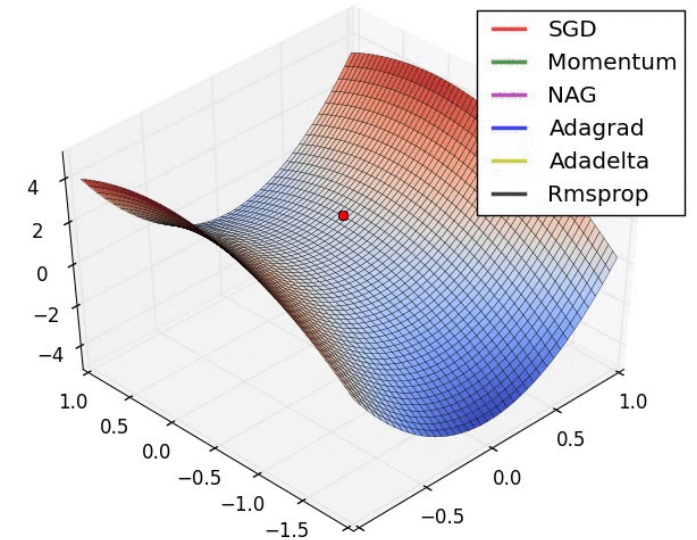Model based RL **<u>may be</u>** a much more practical path to real world robotics



Transfer/Adaptive

Efficiency

Simplicity

Likely to be the most future proof one!

# Are models really that different than Q-functions?

Models                                                    Q-functions

**Similar**

1. Off-policy
2. Models the future

Very different than PG methods → on-policy, models current given future

**Different**

| Models | Q-functions |
|---|---|
| 1. 1-step modeling | 1. Cumulative modeling |
| 2. Models states | 2. Models returns |
| 3. Can evaluate arbitrary policies | 3. Can evaluate only policy $\pi$ |
| 4. Parametric storage of training data | 4. Non-parametric storage of data |

# Class Structure