

The background features a large, semi-transparent blue circle on the right side. To its left, there are overlapping geometric shapes in various shades of teal and green, including a grid pattern and diagonal stripes. The overall aesthetic is modern and tech-oriented.

MBRL

BASICS

- Intuition: Habitual vs Goal-directed
- Learns an environment/dynamics model/network
- Function Approximation: $P_\varphi: S \times A \rightarrow R \times S'$
- Model rollout: $S_t \rightarrow A_t | \pi_\theta(S_t) \rightarrow R_t | P_\varphi(S_t, A_t) \rightarrow S_{t+1} | P_\varphi(S_t, A_t) \rightarrow \dots$

MODEL-BASED CONTROL

- Model rollout from s_0
- Objective: $\min_{a_0 \dots a_{T-1}} \|s_{T-1} - s^*\|$ or $\max_{a_0 \dots a_{T-1}} \sum_{t=0}^{T-1} r_t$

MPC

- $s_0 \rightarrow a_0 | \pi_\theta(s_0) \rightarrow s_1 | P_\phi(s_0, a_0)$
- Backpropagate using objective $\max_{a_1 \dots a_{T-1}} \sum_{t=1}^{T-1} r_t$ from model unrolling
- $s_0 \leftarrow s_1$
- Repeat

ALEATORIC VS EPISTEMIC UNCERTAINTY

- Aleatoric def: Uncertainty due to system/environment(Latin: game of chance)
- Epistemic def: Uncertainty due to insufficient/biased data
- Compounding error(Long term)
- Probabilistic ensemble networks(Gaussian)
- E.g. PETS & MBPO

Algorithm 2 Model-Based Policy Optimization with Deep Reinforcement Learning

- 1: Initialize policy π_ϕ , predictive model p_θ , environment dataset \mathcal{D}_{env} , model dataset $\mathcal{D}_{\text{model}}$
 - 2: **for** N epochs **do**
 - 3: Train model p_θ on \mathcal{D}_{env} via maximum likelihood
 - 4: **for** E steps **do**
 - 5: Take action in environment according to π_ϕ ; add to \mathcal{D}_{env}
 - 6: **for** M model rollouts **do**
 - 7: Sample s_t uniformly from \mathcal{D}_{env}
 - 8: Perform k -step model rollout starting from s_t using policy π_ϕ ; add to $\mathcal{D}_{\text{model}}$
 - 9: **for** G gradient updates **do**
 - 10: Update policy parameters on model data: $\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi, \mathcal{D}_{\text{model}})$
-

The background features a complex design. On the left, there are overlapping geometric shapes in various shades of green and cyan, some with a grid-like pattern. On the right, a large circular gradient transitions from a light blue at the top to a dark blue at the bottom. The text 'OFF-POLICY RL' is centered in the white space of the circular gradient.

OFF-POLICY RL

DEFINITION

- Learns value of the optimal policy disregarding agent's actions
- Behavior policy(Correlated with current)
- On policy algos: SARSA, policy gradient methods
- Off policy algos: Q-learning, soft a2c, DDPG

DEEP DETERMINISTIC POLICY GRADIENT (DDPG)

- Objective: $\max_{\theta} \mathbb{E}_{\tau \sim P_{\theta}(\tau)} [\sum_{t=0}^{T-1} Q(s_t, a_t)]$
- $\nabla_{\theta} \mathbb{E}[\sum_{t=0}^{T-1} Q(s_t, a_t)] = \mathbb{E}[\sum_{t=0}^{T-1} \nabla_{a_t} Q(s_t, a_t) \nabla_{\theta} a_t]$
- Learns from experience tuples in a buffer(Not fixed)
- Extrapolation error: Distinguish good (s,a)s from the bad ones(states visited by policy).
- Not truly off-policy

Algorithm 1 Deep Deterministic Policy Gradient

- 1: Input: initial policy parameters θ , Q-function parameters ϕ , empty replay buffer \mathcal{D}
- 2: Set target parameters equal to main parameters $\theta_{\text{targ}} \leftarrow \theta$, $\phi_{\text{targ}} \leftarrow \phi$
- 3: **repeat**
- 4: Observe state s and select action $a = \text{clip}(\mu_{\theta}(s) + \epsilon, a_{\text{Low}}, a_{\text{High}})$, where $\epsilon \sim \mathcal{N}$
- 5: Execute a in the environment
- 6: Observe next state s' , reward r , and done signal d to indicate whether s' is terminal
- 7: Store (s, a, r, s', d) in replay buffer \mathcal{D}
- 8: If s' is terminal, reset environment state.
- 9: **if** it's time to update **then**
- 10: **for** however many updates **do**
- 11: Randomly sample a batch of transitions, $B = \{(s, a, r, s', d)\}$ from \mathcal{D}
- 12: Compute targets

$$y(r, s', d) = r + \gamma(1 - d)Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s'))$$

- 13: Update Q-function by one step of gradient descent using

$$\nabla_{\phi} \frac{1}{|B|} \sum_{(s,a,r,s',d) \in B} (Q_{\phi}(s, a) - y(r, s', d))^2$$

- 14: Update policy by one step of gradient ascent using

$$\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} Q_{\phi}(s, \mu_{\theta}(s))$$

- 15: Update target networks with

$$\begin{aligned}\phi_{\text{targ}} &\leftarrow \rho \phi_{\text{targ}} + (1 - \rho)\phi \\ \theta_{\text{targ}} &\leftarrow \rho \theta_{\text{targ}} + (1 - \rho)\theta\end{aligned}$$

- 16: **end for**
 - 17: **end if**
 - 18: **until** convergence
-

DDPG EXTENSIONS

- Hindsight Experience Replay
- Twin Delayed Deep Deterministic Policy Gradients (TD3)



VISUAL IMITATION LEARNING

IDEA

- Use visual feature matching as dense reward shaping from a single demonstration
- Grasp useful info from the real world(Computer vision)

SINGLE DEMONSTRATION AND MANY INTERACTIONS

- Self-supervised visual learning
- 3D body pose inference through SMPL(a 3D human shape low-parametric model), e.g. keypoint, segmentation and motion reprojections
- PPO & iLQR

LARGE SCALE DEMONSTRATION COLLECTION

- No/Little interactions
- Challenges: Diversity(Behavioural strategies) & Suboptimality(Longer trajectories)
- cVAE(conditional Variational Autoencoder): Generate all possible subgoals(Diversity), Low-level goal-conditioning → unimodal imitation(Suboptimality)

Review: LQR, iLQR

Recitation 9 10-403

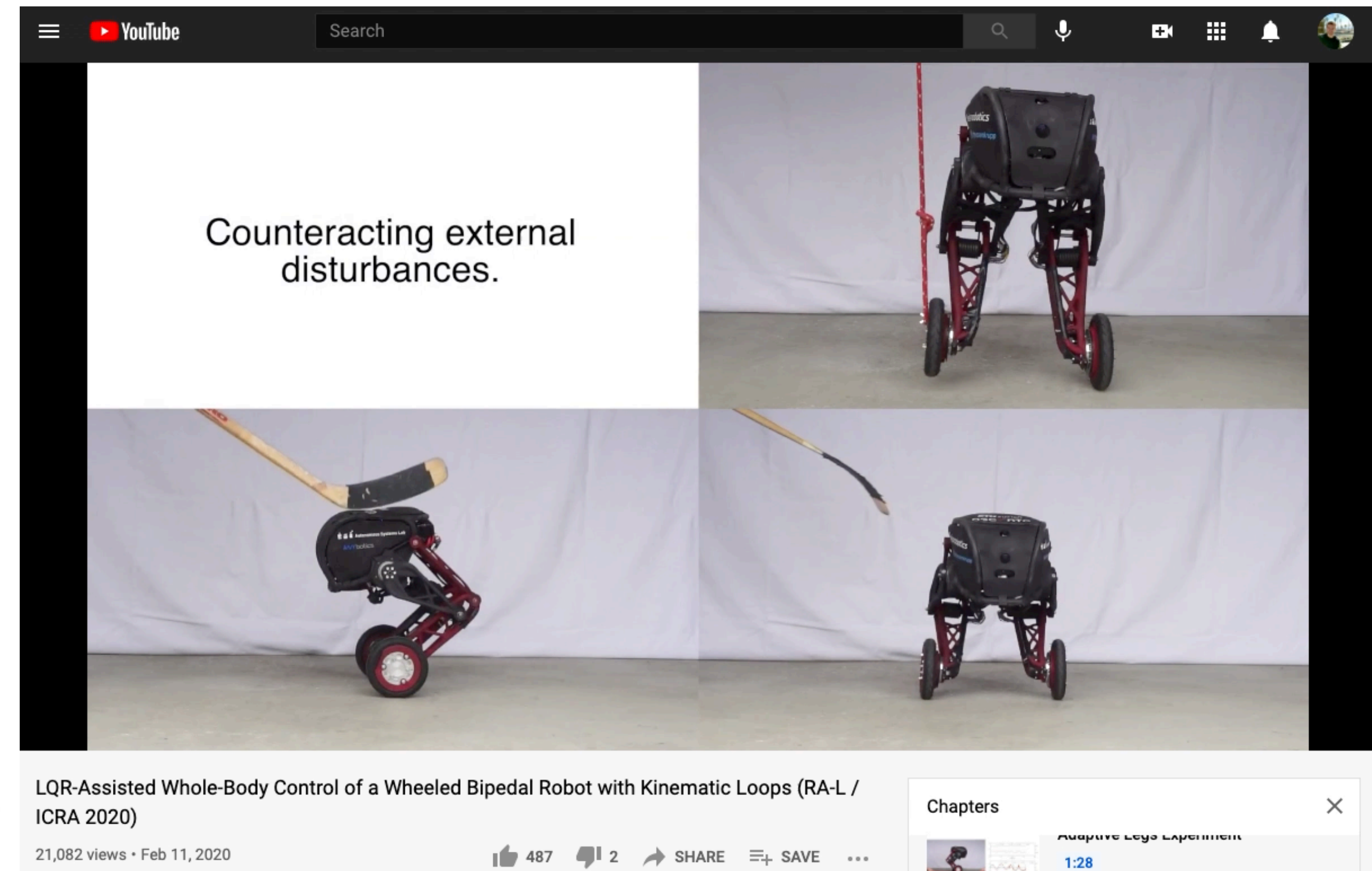
Review

LQR Overview

LQR-Assisted Whole-Body Control of a Wheeled Bipedal Robot with Kinematic Loops

Victor Klemm, Alessandro Morra, Lionel Gulich, Dominik Mannhart,
David Rohr, Mina Kamel, Yvain de Viragh, and Roland Siegwart

September 2019



The video player shows a sequence of four frames illustrating a robot's response to a disturbance. The top-left frame contains the text "Counteracting external disturbances." The top-right frame shows the robot being pulled by a red string. The bottom-left frame shows the robot tilted as it reacts to the pull. The bottom-right frame shows the robot returning to an upright position. The video title is "LQR-Assisted Whole-Body Control of a Wheeled Bipedal Robot with Kinematic Loops (RA-L / ICRA 2020)" and it has 21,082 views as of Feb 11, 2020. The video player interface includes a search bar, navigation icons, and a chapters list with a chapter titled "Adaptive Logo Experiment" at 1:28.

Review

LQR Overview

- $x, u = s, a$; $x \in \mathbb{R}^d, u \in \mathbb{R}^k$
- $x_{t+1} = Ax_t + Bu_t$ Linear dynamics
- $c(x_t, u_t) = x_t^T Q x_t + u_t^T R u_t$ Quadratic costs (can think of as negative of MDP rewards)
- Bellman Equation for deterministic policy π (in cost form):
 - $V^\pi(x) = x^T Q x + \pi(x)^T R \pi(x) + V^\pi(Ax + B\pi(x))$
 - Note: discounting not necessary provided the system is “controllable”
- Bellman Optimality Equation (in cost form):
 - $V^\star(x) = \min_u x^T Q x + u^T R u + V^\star(Ax + Bu)$
- As usual, solving the Bellman Optimality Equation gives us the optimal policy
 - Often this involves solving something known as an ***Algebraic Ricatti Equation***

Review

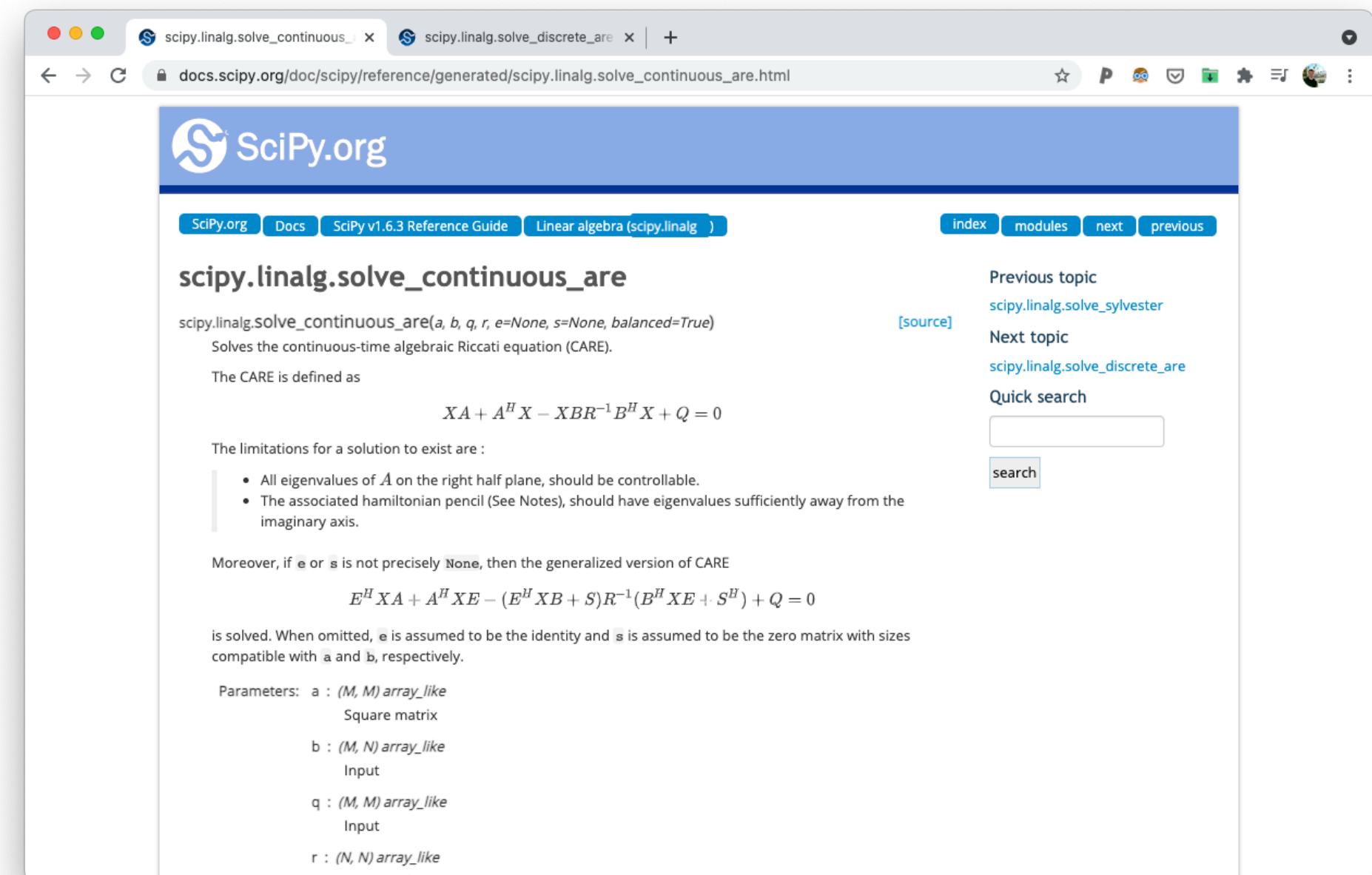
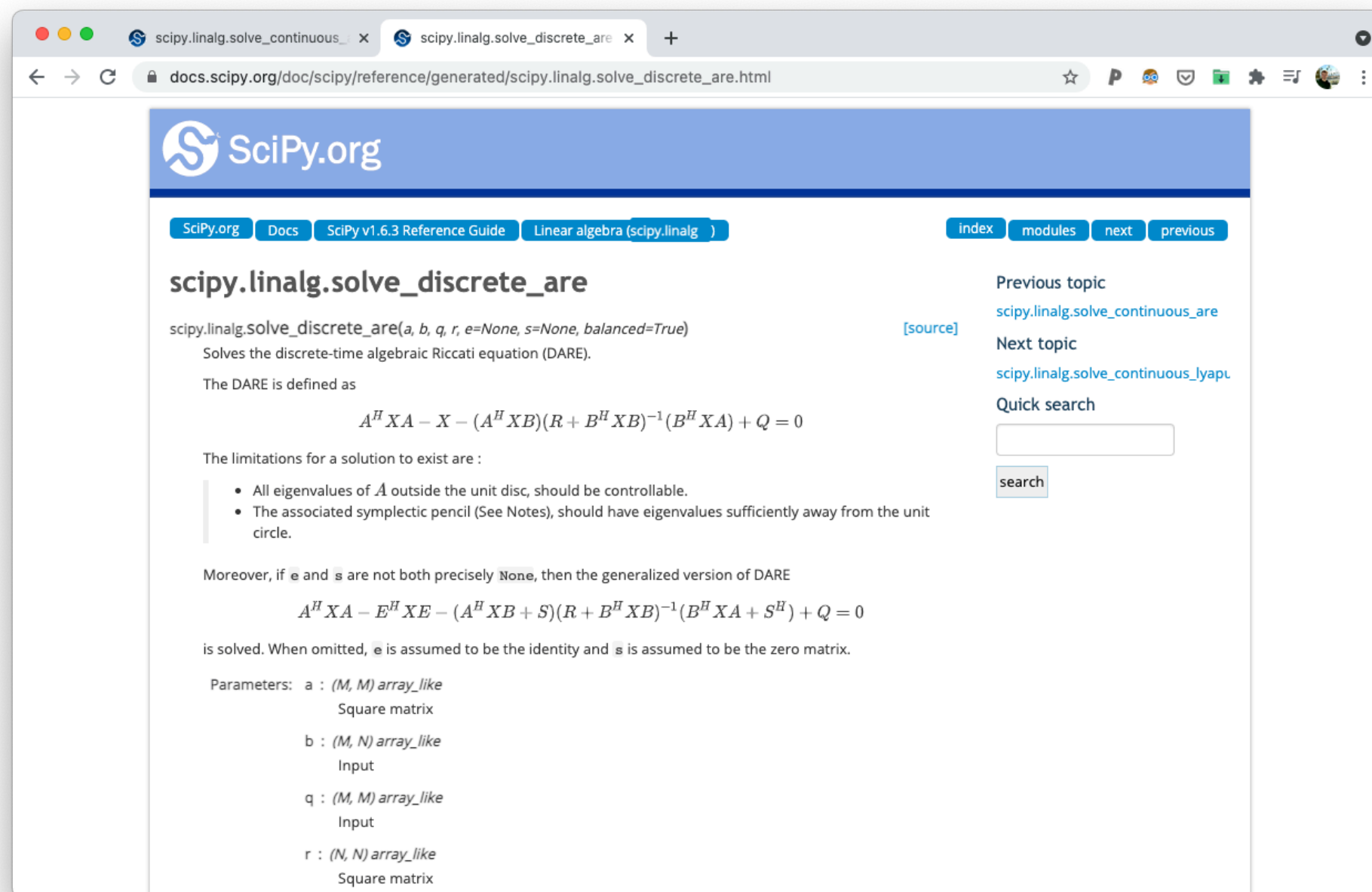
Algebraic Ricatti Equation Overview

$$\begin{aligned}x_{t+1} &= Ax_t + Bu_t \\c(x_t, u_t) &= x_t^T Qx_t + u_t^T Ru_t\end{aligned}$$

- Suppose that V^* can be expressed as a quadratic form (this turns out to be true):
 - $V^*(x) = x^T P^* x = \min_u x^T Qx + u^T Ru + (Ax + Bu)^T P^*(Ax + Bu) = \min_u f(x, u)$
- Notice that $f(x, u)$ is convex in u (this follows from assumptions on Q, R), so we can take gradients and set to zero to solve RHS:
 - $\nabla_u f(x, u) = 2Ru + 2B^T P^*(Ax + Bu) \stackrel{\text{set}}{=} 0 \implies u^* = -(R + B^T P^* B)^{-1} B^T P^* Ax = -Kx$
- Plugging back into the Bellman Optimality Equation:
 - $$\begin{aligned}V^*(x) &= x^T P^* x = \min_u x^T Qx + u^T Ru + (Ax + Bu)^T P^*(Ax + Bu) \\&= x^T Qx + u^{*T} Ru^* + (Ax + Bu^*)^T P^*(Ax + Bu^*) \\&= x^T Qx + x^T K^T R Kx + (Ax - BKx)^T P^*(Ax - BKx)\end{aligned}$$
- This yields the following matrix equation (if we solve this then we solve the Bellman Optimality Equation):
 - $P^* = Q + K^T R K + (A - BK)^T P^*(A - BK)$ (where $K = (R + B^T P^* B)^{-1} B^T P^* A$)
- Substituting K for the full expression and manipulating we arrive at the **Discrete-time Algebraic Ricatti Equation (DARE)**:
 - $P^* = Q + A^T P^* A - (A^T P^* B)(R + B^T P^* B)^{-1} (B^T P^* A)$
- Once we solve for P^* we can then easily find K (also known as the gain matrix) to find the optimal policy $\pi^*(x) = -Kx$

Review

LQR Extensions



- Many variants on the LQR setting, including:
 - Continuous-time LQR (requires solving **CARE** instead of **DARE**)
 - Time-varying LQR with varying dynamics A_t, B_t (e.g. imagine an ageing motor)
 - Requires a different solution method and yields $\{K_t\}$ instead of constant K as in vanilla LQR
 - Having non-zero, time-varying regulation points (vanilla LQR assumes we always wish to keep the state at the origin)
 - If at time t the regulation point is x_t^* then we take action $-K(x_t - x_t^*)$ where K is from vanilla LQR

Review

iLQR Overview

- In practice, few systems are truly linear throughout all state-action space
- Popular algorithm: iLQR (similar to Differential Dynamic Programming, DDP). Note that iLQR is a form of Model Predictive Control. Note also that iLQR assumes that the dynamics model is known perfectly (various papers have addressed model misspecification issues and robustness; however we will just focus on the vanilla iLQR algorithm)

Input: initial action sequence $\{u_0^0, \dots, u_{T-1}^0\}$ and state trajectory $\{x_0^0, \dots, x_{T-1}^0\}$, goal state trajectory $\{x_0^\star, \dots, x_{T-1}^\star\}$

For $i = 0, 1, \dots, I$:

1. Linearise f at each (x_t^i, u_t^i) yielding $\{(A_t^i = D_x f(x_t^i, u_t^i), B_t^i = D_u f(x_t^i, u_t^i))\}_{t=0}^{T-1}$ (use finite differences or auto-diff software)
2. Solve a time-varying LQR problem using $\{(A_t^i, B_t^i)\}$ and Q, R yielding $\{K_t^i\}_{t=0}^{T-1}$ (see lecture notes for backward pass equations)
3. Set α to some positive value and line search over α until we find an action sequence with better total cost:
 - 3.1. Rollout $\{\alpha K_t^i\}_{t=0}^{T-1}$ from x_0 using $\{x_0^\star, \dots, x_{T-1}^\star\}$ to obtain $\{u_0^{i+1}, \dots, u_{T-1}^{i+1}\}$ and $\{x_0^{i+1}, \dots, x_{T-1}^{i+1}\}$: $u_t^{i+1} = u_t^i - \alpha K_t^i (x_t^{i+1} - x_t^i - x_t^\star)$
 - 3.2. $\alpha \leftarrow \frac{\alpha}{2}$
4. Break if stopping criteria met (e.g. changes in $\{u_0^{i+1}, \dots, u_{T-1}^{i+1}\}$ small)

Return u_0^I

Review

iLQR Overview

- In practice, few systems are truly linear throughout all state-action space
- Popular algorithm: iLQR (similar to Differential Dynamic Programming, DDP). Note that iLQR is a form of Model Predictive Control. Note also that iLQR assumes that the dynamics model is known perfectly (various papers have addressed model misspecification issues and robustness; however we will just focus on the vanilla iLQR algorithm)

Input: initial action sequence $\{u_0^0, \dots, u_{T-1}^0\}$ and state trajectory $\{x_0^0, \dots, x_{T-1}^0\}$, goal state trajectory $\{x_0^*, \dots, x_{T-1}^*\}$

For $i = 0, 1, \dots, I$:

1. Linearise f at each (x_t^i, u_t^i) yielding $\{(A_t^i = D_x f(x_t^i, u_t^i), B_t^i = D_u f(x_t^i, u_t^i))\}_{t=0}^{T-1}$ (use finite differences or auto-diff software)
2. Solve a time-varying LQR problem using $\{(A_t^i, B_t^i)\}$ and Q, R yielding $\{K_t^i\}_{t=0}^{T-1}$ (see lecture notes for backward pass equations)
3. Set α to some positive value and line search over α until we find an action sequence with better total cost:
 - 3.1. Rollout $\{\alpha K_t^i\}_{t=0}^{T-1}$ from x_0 using $\{x_0^*, \dots, x_{T-1}^*\}$ to obtain $\{u_0^{i+1}, \dots, u_{T-1}^{i+1}\}$ and $\{x_0^{i+1}, \dots, x_{T-1}^{i+1}\}$: $u_t^{i+1} = u_t^i - \alpha K_t^i (x_t^{i+1} - x_t^i - x_t^*)$
 - 3.2. $\alpha \leftarrow \frac{\alpha}{2}$
4. Break if stopping criteria met (e.g. changes in $\{u_0^{i+1}, \dots, u_{T-1}^{i+1}\}$ small)

Return u_0^I

1. $-x_t^*$ because not interested in regulating to origin
2. $-x_t^i$ because LQR solution is w.r.t. linearisation point
3. α because our linearisation may not hold well and we need to act more closely to the linearisation point
4. u_t^i because LQR solution is w.r.t. linearisation point

Review

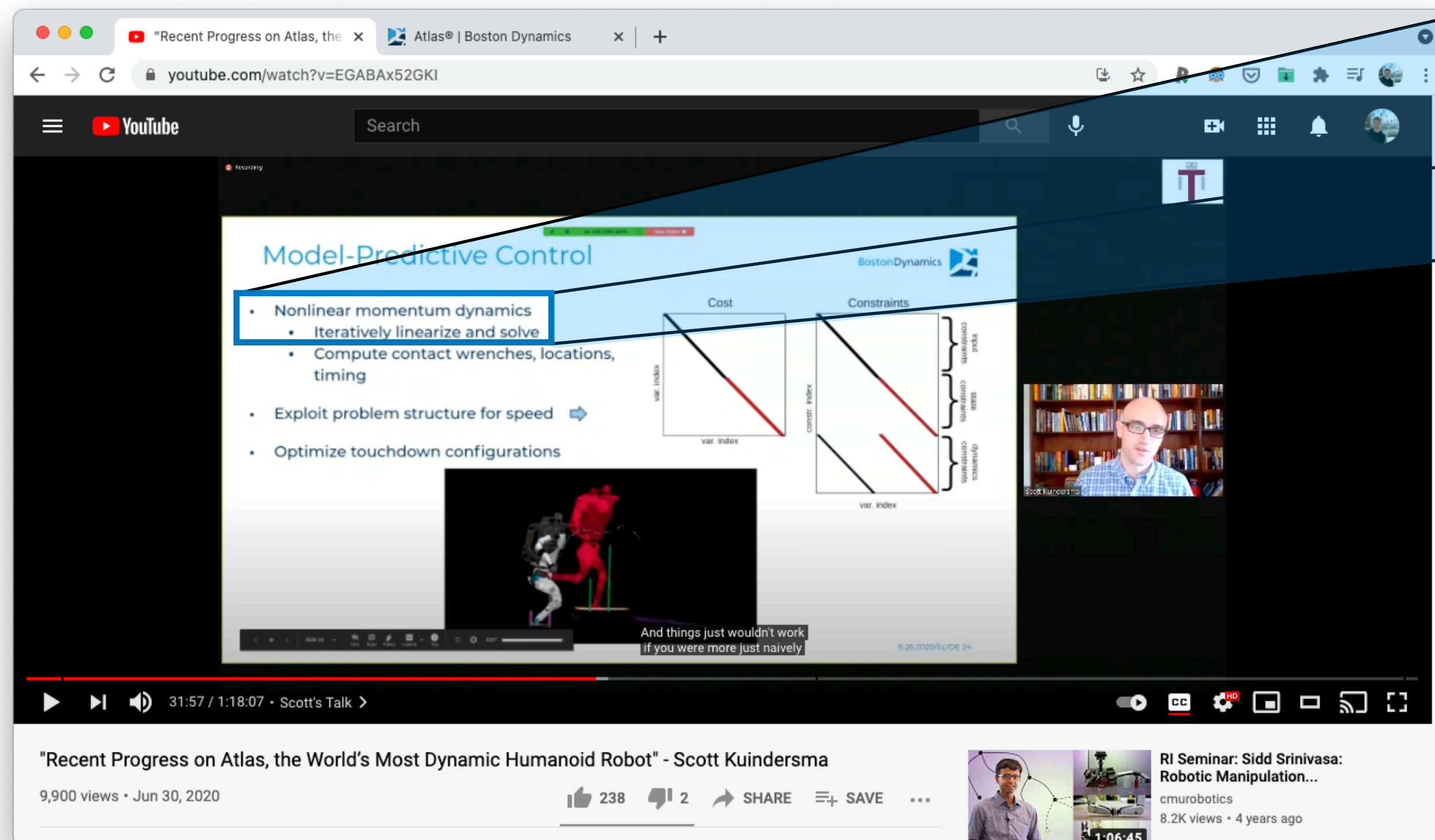
iLQR Overview

- Issues with iLQR:
 - Computationally expensive: requires calculating multiple Jacobians and solving multiple LQR problems every time we need to take an action.
 - Mostly suited to “regulator like” problems (i.e. not Atari playing) as ultimately is depending on smoothness in underlying domain to perform well (the “larger” the non-linearity in dynamics the poorer iLQR will perform, although still better than vanilla LQR)

Review

iLQR Overview

- Still, lots of robotics problems leverage ideas from iLQR
 - Boston Dynamics recently disclosed Atlas control details, likely builds off iLQR ideas: <https://www.youtube.com/watch?v=EGABAx52GKI>



The screenshot shows a YouTube video player with a presentation slide titled "Model-Predictive Control" by Boston Dynamics. The slide content includes:

- Nonlinear momentum dynamics
 - Iteratively linearize and solve
 - Compute contact wrenches, locations, timing
- Exploit problem structure for speed
- Optimize touchdown configurations

The slide also features two graphs: "Cost" and "Constraints", both showing a linear relationship between variables. A small inset video shows a speaker, and a larger video shows a humanoid robot (Atlas) performing a task. The YouTube interface shows the video is at 31:57 / 1:18:07, titled "Recent Progress on Atlas, the World's Most Dynamic Humanoid Robot" by Scott Kuindersma, with 9,900 views as of June 30, 2020.

- Nonlinear momentum dynamics
 - Iteratively linearize and solve



CURIOSITY DRIVEN EXPLORATION



TYPES OF MOTIVATION

- Extrinsic Motivation : pursuit of an external reward (such rewards can be sparse)
- Intrinsic Motivation : pursuit because the task is inherently enjoyable (curiosity, novelty, surprise)

Benefits of Intrinsic Motivation :

- Task independent
- Free of human supervision
- Does not require reward functions to incentivize agent

CURIOSITY

- Seek novelty / surprise
- Visit novel states
- Observe novel state transitions

CURIOSITY DRIVEN EXPLORATION

1. Ensembles of Q functions
2. State counting
3. Model prediction error
4. Episodic curiosity through reachability

ENSEMBLE OF Q-FUNCTIONS

- We are looking to model the uncertainty of Q values
 - Train multiple Q-functions each one using different subset of data (can also have a shared backbone)
 - Use posterior sampling to sample a Q-function ($Q \sim P(Q)$)
 - Choose actions according to the sampled Q-function
 - Update the Q distribution using the collected experience tuples
-
- Entropy of predictions of the network is high in areas where lesser data has been previously collected, thus encouraging exploration
 - When Q-functions agree on actions, entropy is lower and there we exploit instead of explore

EXPLORATION REWARD

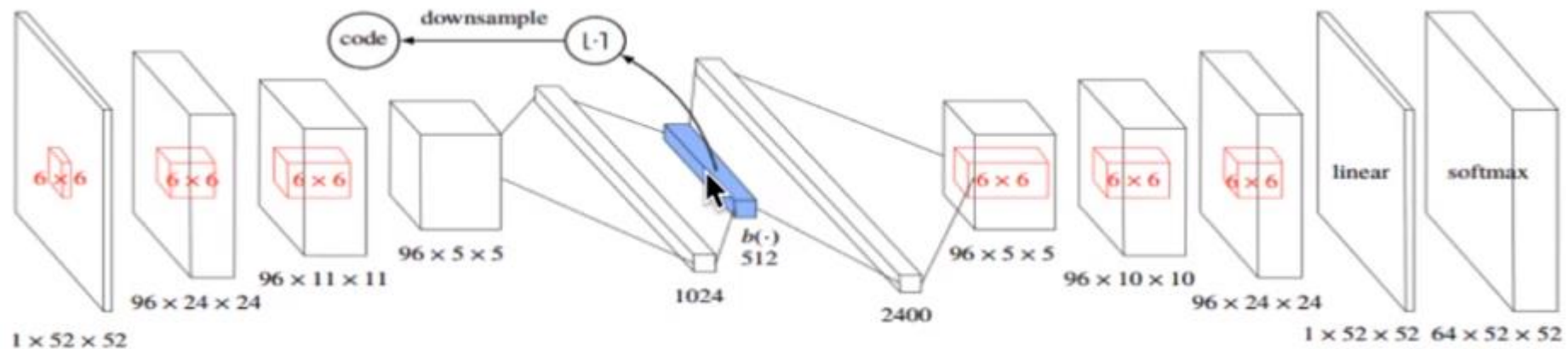
$$R^t(s, a, s') = \underbrace{r(s, a, s')}_{\text{extrinsic}} + \underbrace{\mathcal{B}^t(s, a, s')}_{\text{intrinsic}}$$

Independent of the task in hand!

- Add exploration reward bonuses to the extrinsic (task-related) rewards
- Exploration reward bonuses are non stationary

STATE COUNTING

- Count states in latent space (not in image space)
- Compression into latent space is important because the image space is very large
- Use autoencoders to get the image encoding (based on reconstruction)
- Try to visit states that have a lower count



MODEL PREDICTION ERROR

- We want a positive reinforcement whenever the system fails to correctly predict the environment
- Train a dynamics model alongside
- Exploration reward bonuses encourage policy to visit states that cause the trained dynamics model to fail

$$R^t(s, a, s') = \underbrace{r(s, a, s')}_{\text{extrinsic}} + \underbrace{\mathcal{B}^t(\|T(s, a; \theta) - s'\|)}_{\text{intrinsic}}$$

model error!

LIMITATIONS OF PREDICTION ERROR AS BONUS

- Agent can be rewarded even though the model can no longer improve
- Agent is attracted to the most noisy states. This can be random noise with unpredictable outcomes.
- Noisy TV problem

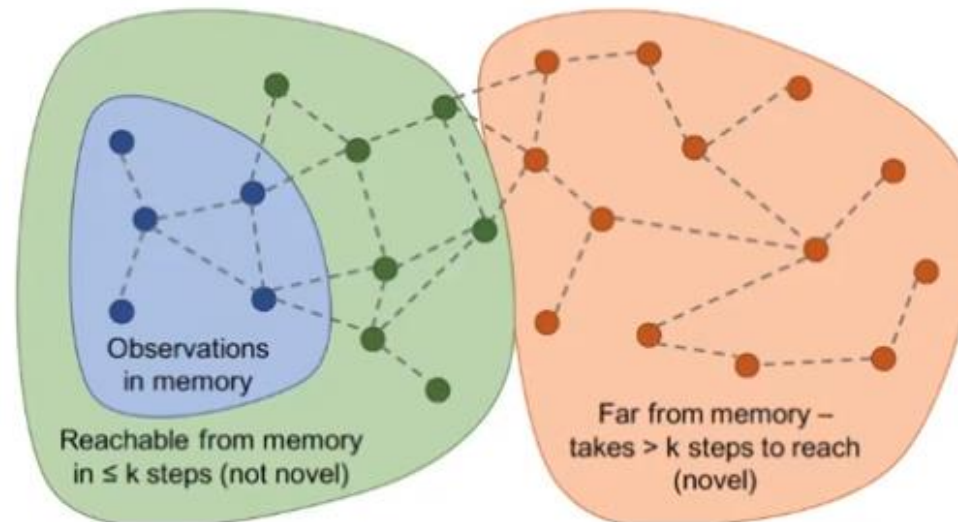
EPIODIC CURIOSITY THROUGH REACHABILITY

$$R^t(s, a, s') = \underbrace{r(s, a, s')}_{\text{extrinsic}} + \underbrace{\mathcal{B}^t(s, \mathcal{M})}_{\text{intrinsic}}$$

- Augmented rewards as before
- \mathcal{M} is a non-parametric memory populated with past image observations
- Curiosity reward will use a comparator neural net that takes in two images and predicts whether they are close (few actions apart) or far
- Comparisons are done between current observation and observations in memory

EPIODIC CURIOSITY THROUGH REACHABILITY

- At each timestep, agent compare current observation with the ones in memory
- If current observation is novel, agent gets rewarded and novel observation is added to the memory



SIM2REAL TRANSFER



CHOICES

1. Use a physics simulator to learn and then transfer to real world
2. Learn policies directly in the real world
3. Combine simulators with deeply learned residuals

PROS AND CONS OF SIMULATION

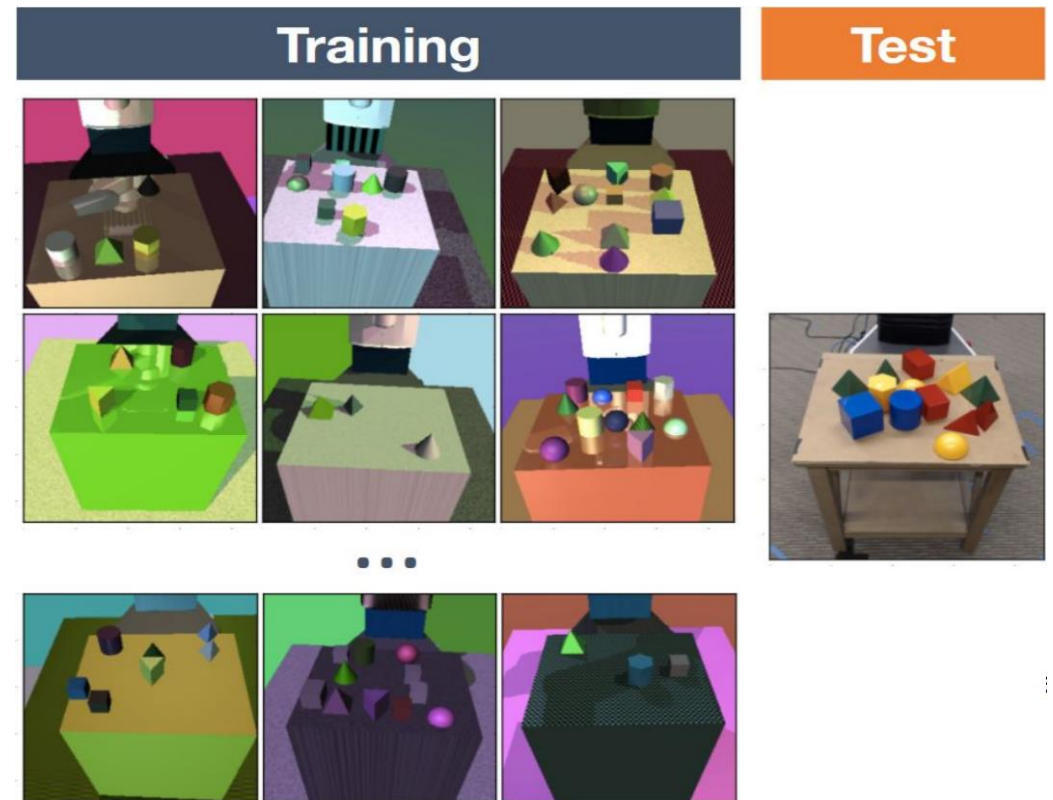
Pros	Cons
Can afford many samples	Under modeling
Safe	Large engineering effort
Allow more exploration	Needs accurate parameters

STRATEGIES FOR SIM2REAL TRANSFER

1. Domain Randomization
2. Adaptive Domain Randomization
3. Residual Physics (combine analytical models with deeply learned residuals)
4. Abstraction

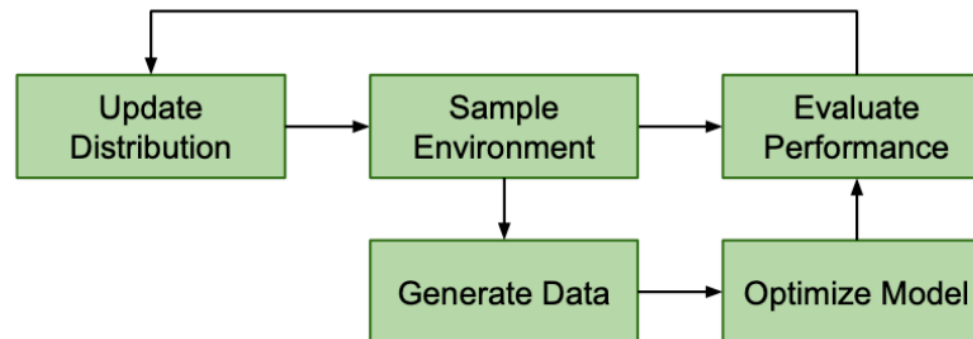
DOMAIN RANDOMIZATION

- Create many version of simulation environments by randomizing parameters
- Parameters of the environment : camera position and orientation, textures, lighting, contrast
- Train the detector on the diverse simulation environments



ADAPTIVE DOMAIN RANDOMIZATION

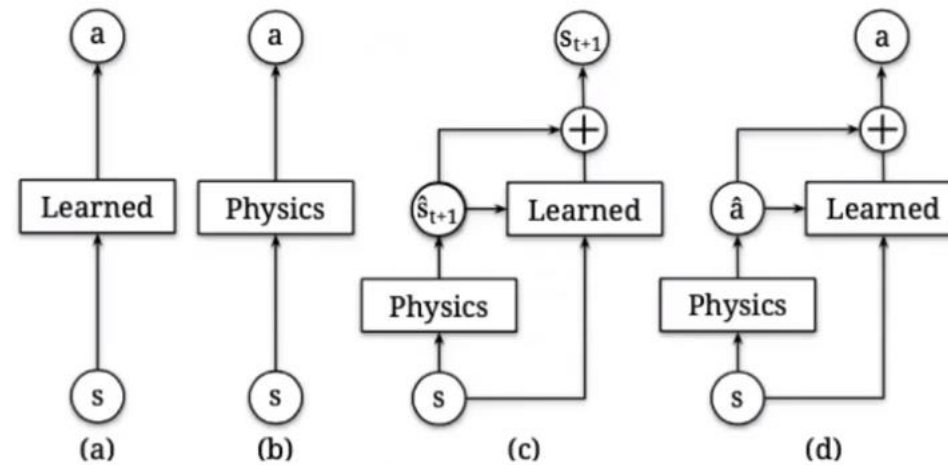
- Policy might fail if the distribution of environments is too wide
- Widen the distribution as long as the policy succeeds
- Increase the range of values of the parameters (that characterize the environment) if the policy performs well



RESIDUAL PHYSICS

- Use neural networks to capture the residual between analytical physics models and real world physics
- Paper to read : “TossingBot : Learning to Throw Arbitrary Objects with Residual Physics”

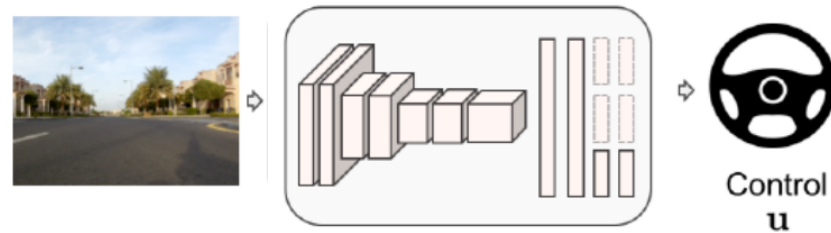
RESIDUAL PHYSICS



- a) Pure learned model for predicting actions
- b) Pure analytical physics model for predicting actions
- c) Combination of analytical and learned model for predicting next state
- d) Combination of analytical and learned model for predicting actions (TossingBot)

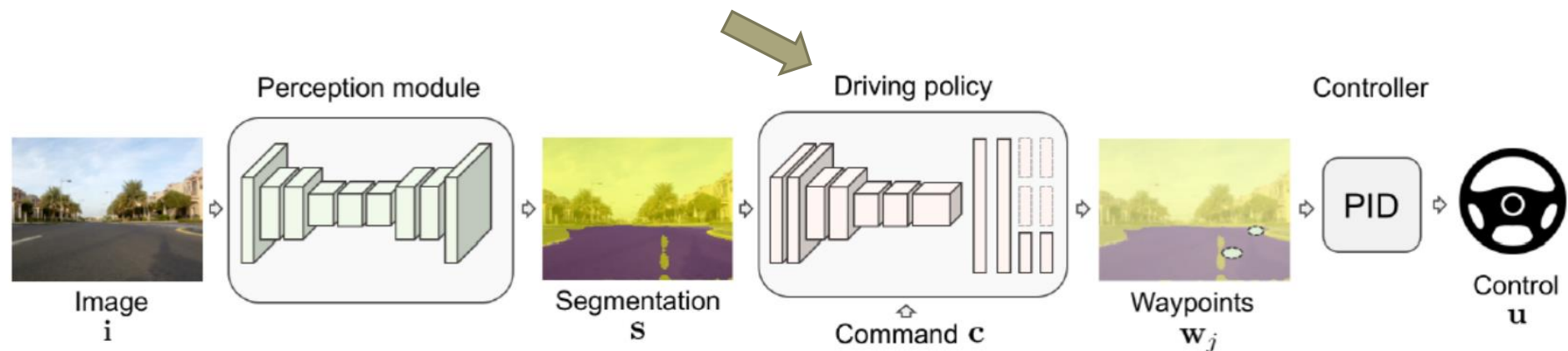
ABSTRACTION

- Learn higher level policies instead of low level controllers
- Paper to read : “Driving Policy Transfer via Modularity and Abstraction”
- Image pixels to steering wheel mapping is not sim2real transferable
 - Image quality mismatch between simulation and real world
 - Dynamics mismatch (physics of the car and steering) between simulation and real world



ABSTRACTION

- Mapping of labelled map to waypoints is better sim2real transferable
- Low level controller can take the car from waypoint to waypoint in the real world

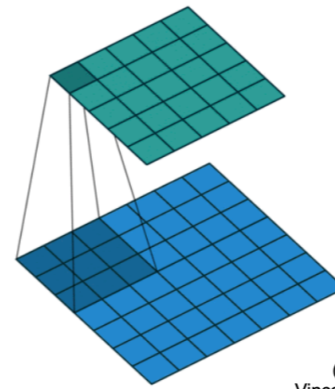


Graph Neural Network / Learning from Demonstrations and Task Rewards / Adversarial Imitation Learning

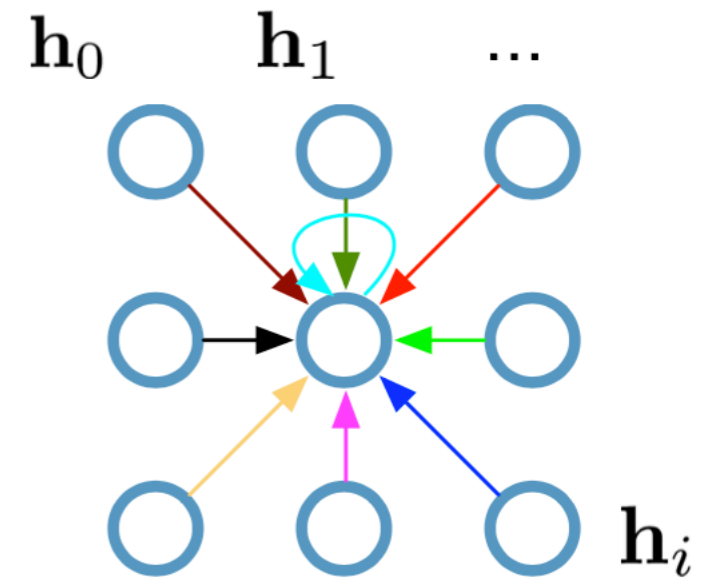
**-
A review**

Graph Neural Network

Single CNN layer
with 3x3 filter:



(Animation by
Vincent Dumoulin)



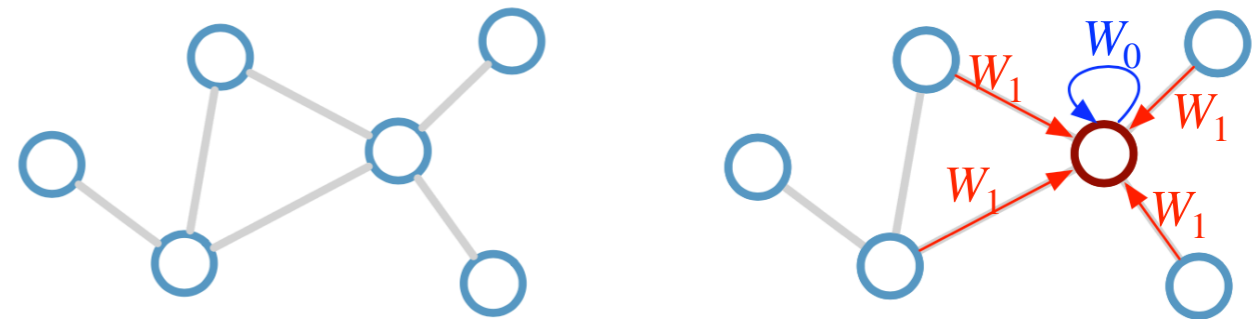
From CNN to GNN

Hidden layer activations of pixels/nodes $h_i \in \mathbb{R}^F$

$$h_4^{(l+1)} = \sigma \left(\underbrace{W_0^{(l)}} h_0^{(l)} + \underbrace{W_1^{(l)}} h_1^{(l)} + \dots + \underbrace{W_8^{(l)}} h_8^{(l)} \right)$$

• Different weight!

Undirected graph



From CNN to GNN

Hidden layer activations of nodes $h_i \in \mathbb{R}^F$

Neighbor indices \mathcal{N}_i

Fixed, trainable norm c_{ij}

$$h_i^{(l+1)} = \sigma \left(W_0^{(l)} h_i^{(l)} + \sum_{j \in \mathcal{N}_i} \frac{1}{c_{ij}} W_1^{(l)} h_j^{(l)} \right)$$

- Sharing weight over neighboring locations
- Permutation invariant
- Linear complexity

Paper I

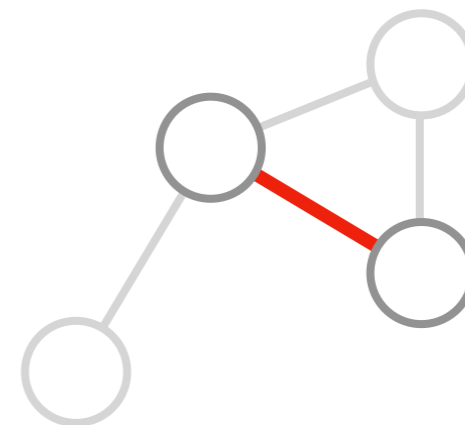
Interaction Networks for Learning about Objects, Relations and Physics, Battaglia et al., 2016

- Object / object part → graph's node
- Input
 - Object state
 - Dynamic: position, velocity
 - Static: mass, size, shape
 - Relation attribute
 - Coefficients of connectivity, restitution, spring constant, etc.
- Output: future velocity

Paper I

Interaction Networks for Learning about Objects, Relations and Physics, Battaglia et al., 2016

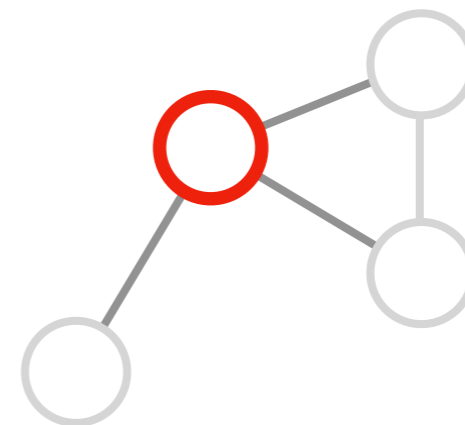
- Relational / **edge** NN
 - Input
 - Two object (two sides of an edge) states
 - Relational attributes
 - Output
 - Feature vector



Paper I

Interaction Networks for Learning about Objects, Relations and Physics, Battaglia et al., 2016

- Object / **node** NN
 - Input
 - One object state
 - Summation of all incoming edge message
 - Output
 - Future object velocity



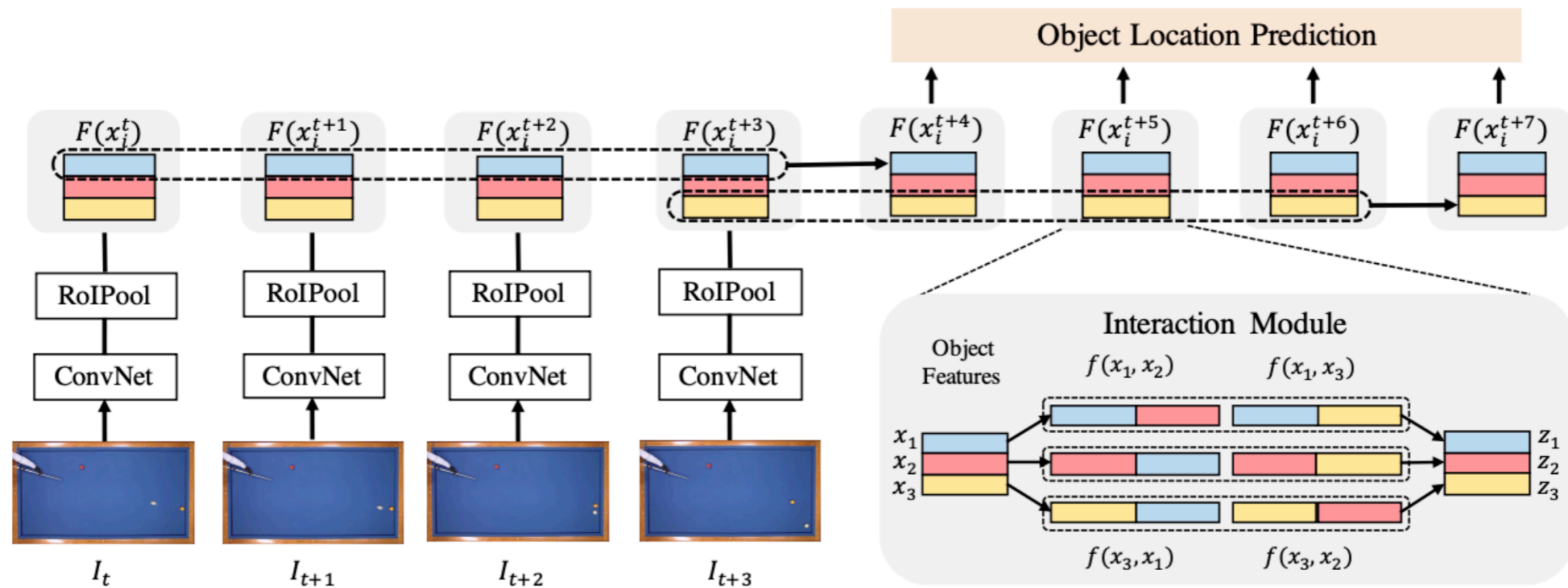
Paper II

Long-term Visual Dynamics with Region Proposal Interaction Networks, Qi et al., 2021

- Learning object interactions from visual input
- Input
 - Detected object appearance (ConvNet, pooling)
 - Motion history
- Output
 - Future motion trajectory

Paper II

Long-term Visual Dynamics with Region Proposal Interaction Networks, Qi et al., 2021



Paper III

Graph Networks as Learnable Physics Engines for Inference and Control, Gonzalez et al., 2018

- Physical system's bodies and joints → graph's nodes and edges
- Input
 - Global feature
 - None / Gravity, wind, density, etc.
 - Node feature
 - Dynamic: position, 4D quaternion orientation, linear and angular velocity
 - Static: mass, inertia, etc.
 - Edge feature
 - Magnitude of action at joint, etc.
- Output: dynamic features, temporal difference.

Paper III

Graph Networks as Learnable Physics Engines for Inference and Control, Gonzalez et al., 2018

G.1. Dynamic graph

We retrieved the the absolute position, orientation, linear and angular velocities for each body:

- **Global: None**
- **Nodes:** (for each body)
 - Absolute body position (3 vars): `mjData.xpos`
 - Absolute body quaternion orientation position (4 vars): `mjData.xquat`
 - Absolute linear and angular velocity (6 vars): `mj_objectVelocity (mjOBJ_XBODY, flg_local=False)`
- **Edges:** (for each joint) Magnitude of action at joint: `mjData.ctrl` (0, if not applicable).

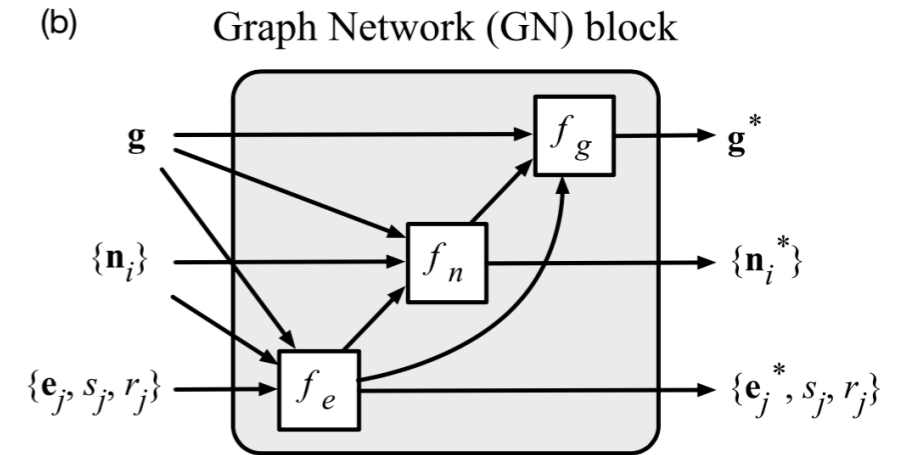
Paper III

Graph Networks as Learnable Physics Engines for Inference and Control, Gonzalez et al., 2018

G.2. Static graph

We performed an exhaustive selection of global, body, and joint static properties from mjModel:

- **Global:** mjModel.opt.{timestep, gravity, wind, magnetic, density, viscosity, impratio, o_margin, o_solref, o_solimp, collision_type (one-hot), enableflags (bit array), disableflags (bit array)}.
- **Nodes:** (for each body) mjModel.body_{mass, pos, quat, inertia, ipos, iquat}.
- **Edges:** (for each joint)
 - Direction of edge (1: parent-to-child, -1: child-to-parent).
 - Motorized flag (1: if motorized, 0 otherwise).
 - Joint properties: mjModel.jnt_{type (one-hot), axis, pos, solimp, solref, stiffness, limited, range, margin}.
 - Actuator properties: mjModel.opt.actuator_{biastype (one-hot), biasprm, cranklength, ctrllimited, ctrlrange, dyntype (one-hot), dynprm, forcelimited, forcerange, gaintype (one-hot), gainprm, gear, invweight0, length0, lengthrange}.



Paper III

Graph Networks as Learnable Physics Engines for Inference and Control, Gonzalez et al., 2018

Algorithm 1 Graph network, GN

Input: Graph, $G = (\mathbf{g}, \{\mathbf{n}_i\}, \{\mathbf{e}_j, s_j, r_j\})$

for each edge $\{\mathbf{e}_j, s_j, r_j\}$ **do**

Gather sender and receiver nodes $\mathbf{n}_{s_j}, \mathbf{n}_{r_j}$

Compute output edges, $\mathbf{e}_j^* = f_e(\mathbf{g}, \mathbf{n}_{s_j}, \mathbf{n}_{r_j}, \mathbf{e}_j)$

end for

for each node $\{\mathbf{n}_i\}$ **do**

Aggregate \mathbf{e}_j^* per receiver, $\hat{\mathbf{e}}_i = \sum_{j/r_j=i} \mathbf{e}_j^*$

Compute node-wise features, $\mathbf{n}_i^* = f_n(\mathbf{g}, \mathbf{n}_i, \hat{\mathbf{e}}_i)$

end for

Aggregate all edges and nodes $\hat{\mathbf{e}} = \sum_j \mathbf{e}_j^*, \hat{\mathbf{n}} = \sum_i \mathbf{n}_i^*$

Compute global features, $\mathbf{g}^* = f_g(\mathbf{g}, \hat{\mathbf{n}}, \hat{\mathbf{e}})$

Output: Graph, $G^* = (\mathbf{g}^*, \{\mathbf{n}_i^*\}, \{\mathbf{e}_j^*, s_j, r_j\})$

-
- Relation between node, edge, global features
 - Sequence of update
 - Architecture of GN block

Paper IV

Learning to Simulate Complex Physics with Graph Networks, Gonzalez et al., 2020

- Object → graph of particles, scene → graph of all particles from all objects
- Input
 - Particle velocity of last 5 time steps
- Output: particle acceleration

Paper IV

Learning to Simulate Complex Physics with Graph Networks, Gonzalez et al., 2020

- Particles are different from nodes. Only consider displacements.
- Inject noise to particle velocities during training time to tackle accumulating error.
- Use edges to encode two particles' relative position. No longer use absolute position compared to previous papers' node features. This is for translation invariance.
- Use multiple rounds of message passing, each round has different node, edge weights.

—

Learning from Demonstrations and Task Rewards

Learning from demonstrations

- Accelerate trial-and-error learning by suggesting good actions to try
- Train initial safe policies, to deploy in the real world
- Time consuming
- Sub-optimality, noise, diversity in performing task
- Can not surpass expert

Learning from task rewards

- Cheap supervision
- Right end task encoded via task rewards
- Sample inefficient
- Unsafe initial policy, unsafe to deploy in real world

Learning from demonstrations and task rewards

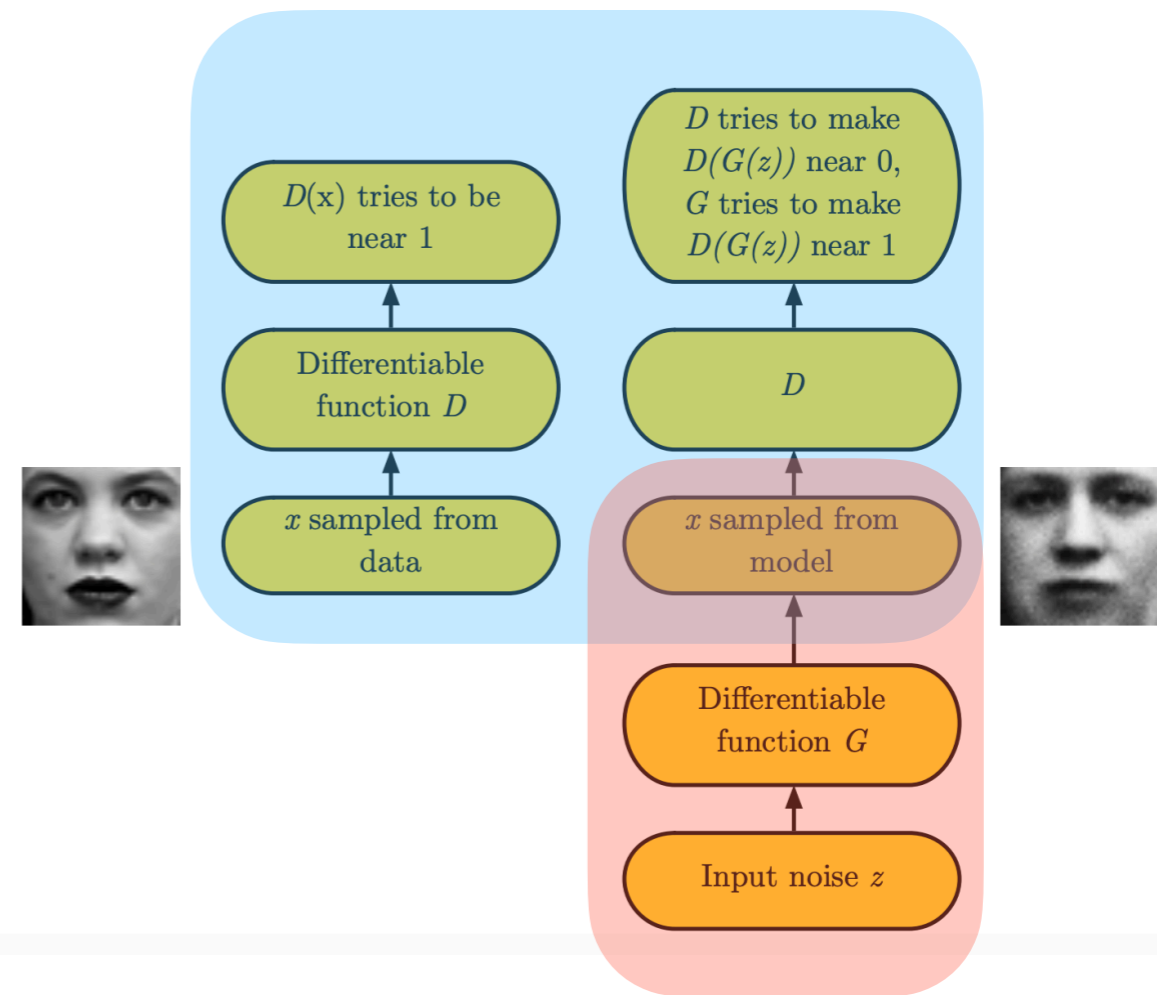
Major considerations

- **Sample efficiency**
- **Outperforming expert**
- **Safe to deploy in real world with good/safe initial performance**

Learning from demonstrations and task rewards

- Initialize the replay buffer with demonstrations
- Pre-train the model-free RL method with a demonstration only buffer, then fine-tune it.
- Combine imitation and task rewards (**see in later discussion**)
- Exploit the temporal structure, and solve longer horizon tasks progressively, rather than solving at once.

Adversarial Imitation Learning



(Goodfellow 2016)

Generative Adversarial Nets (GAN)

Discriminator D

Generator G

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(x)))]$$

- Please refer to lecture slides for explanation of min-max game, which is very clear and specific.
- We usually use Adam (adaptive optimizer) to train discriminator and generator simultaneously.
- It is important to maintain a *balance* between the discriminator and generator, i.e. avoiding one being much stronger than the other.

Generative Adversarial Imitation Learning (GAIL)

Discriminator D_ϕ

- Distinguish between state-action pairs from expert demonstrations and ones generated/visited by agent policy π_θ

Generator π_θ

- Policy network that conditions on state $\pi_\theta(s)$

Reward $r(s, a) = \log D_\phi(s, a), (s, a) \sim \pi_\theta$

- “The **reward** for the policy optimization is how well I matched the demonstrator’s trajectory distribution, else, how well I confused the discriminator.”

Generative Adversarial Imitation Learning (GAIL)

How to update the discriminator D_ϕ ?

$$\mathbb{E}_{(s,a) \sim Demo} \left[\nabla_\phi \log(1 - D_\phi(s, a)) \right] + \mathbb{E}_{(s,a) \in \tau_i} \left[\nabla_\phi \log D_\phi(s, a) \right]$$

How to update the generator / policy π_θ ? Policy gradient (e.g. REINFORCE)

$$\mathbb{E}_{(s,a) \in \tau_i} \left[\nabla_\theta \log \pi_\theta \log D_{\phi_{i+1}}(s, a) \right]$$

GAIL vs. BC (Behavior Cloning)

GAIL outperforms BC, why?

- **GAIL has more interaction with environment**
- **BC reduces imitation learning to supervised learning w.r.t. single action.**

GAIL vs. BC (Behavior Cloning)

7 Discussion and outlook

As we demonstrated, our method is generally quite sample efficient in terms of expert data. However, it is not particularly sample efficient in terms of environment interaction during training. The number of such samples required to estimate the imitation objective gradient (18) was comparable to the number needed for TRPO to train the expert policies from reinforcement signals. We believe that we could significantly improve learning speed for our algorithm by initializing policy parameters with behavioral cloning, which requires no environment interaction at all.

Fundamentally, our method is model free, so it will generally need more environment interaction than model-based methods. Guided cost learning [7], for instance, builds upon guided policy search [13] and inherits its sample efficiency, but also inherits its requirement that the model is well-approximated by iteratively fitted time-varying linear dynamics. Interestingly, both our Algorithm 1 and guided cost learning alternate between policy optimization steps and cost fitting (which we called discriminator fitting), even though the two algorithms are derived completely differently.

Our approach builds upon a vast line of work on IRL [31, 1, 29, 28], and hence, just like IRL, our approach does not interact with the expert during training. Our method explores randomly to determine which actions bring a policy’s occupancy measure closer to the expert’s, whereas methods that do interact with the expert, like DAgger [24], can simply ask the expert for such actions. Ultimately, we believe that a method that combines well-chosen environment models with expert interaction will win in terms of sample complexity of both expert data and environment interaction.

Combining imitation and task rewards

Original reward

$$r(s, a) = \log D_{\phi}(s, a), (s, a) \sim \pi_{\theta}$$

Combined reward

$$r(s, a) = \lambda r_{GAIL}(s, a) + (1 - \lambda) r_{task}(s, a), \lambda \in [0, 1]$$

$$r_{GAIL}(s, a) = -\log(1 - D(s, a))$$

Paper

Reinforcement and Imitation Learning for Diverse Visuomotor Skills, Zhu et al., 2018

- Combines imitation and task rewards
- Start episodes by setting the world in states of the demonstration trajectories.
- Please review the paper.

Thank you

—
A review

Yuning Wu, 5/7/2021