



Real-World RL

“Real world” *noun*

The existing state of things, as opposed to one that is imaginary, simulated, or theoretical. --

Oxford Dictionary

Gabriel Dulac-Arnold
Researcher
Google Research, Paris

Overall Goal

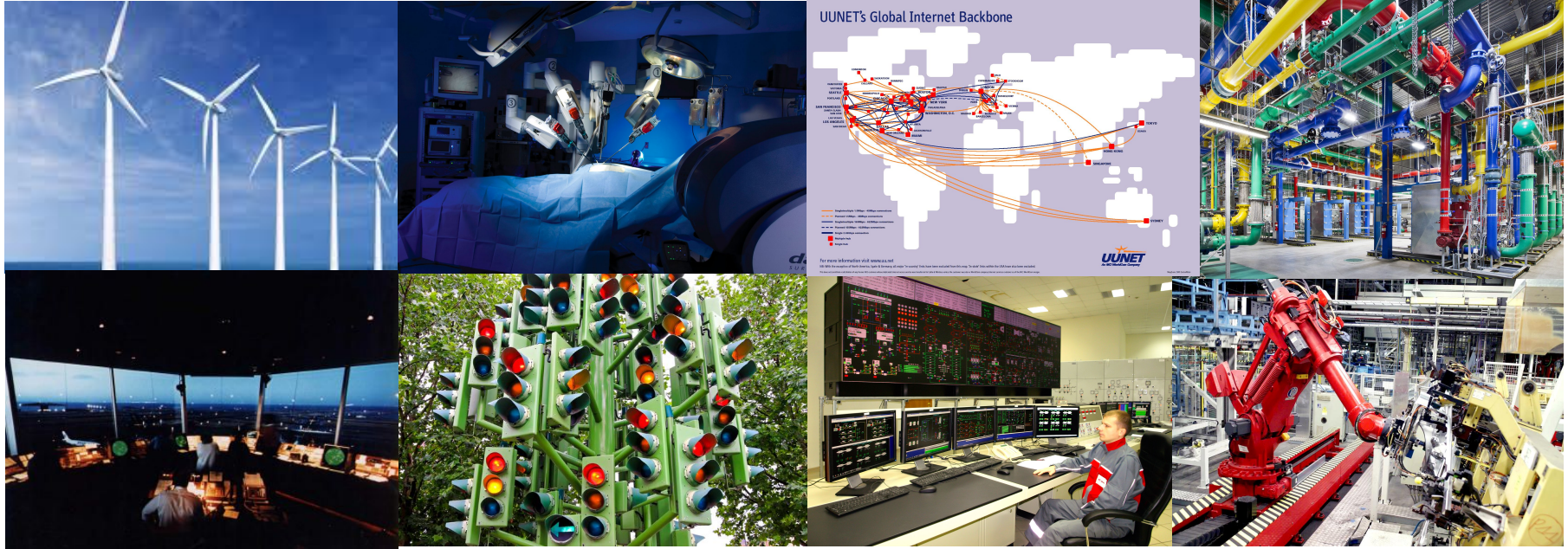
Create learning controllers to:

- Decrease energy and resource utilisation
- Allow for under-specified and imprecise systems
- Adapt to system degradations and increase longevity
- Make larger systems easier to instantiate and modify
- Democratize robotics
- Make programming obsolete

These things work



But what we really want is these



But what would we really like (until AGI)?











Can't we “just apply DQN / AlphaGo”?

- No available simulator!
 - Can't learn on millions of episodes
- Not small discrete action spaces
- Not deterministic
- You can't do whatever you want when learning
- You need to be running at 99.999% something, mistakes are not an option

“Challenges of Real-World RL”

Gabriel Dulac-Arnold, Daniel Mankowitz, Todd Hester
2019 - Google Research, DeepMind

- **Training off-line from the fixed logs of an external behavior policy.**
 - Generally no direct learning on the system, existing logs sub-optimal or random.
- **Learning on the real system from limited samples.**
 - If allowed on the system, must be data-efficient.
- **High-dimensional continuous state and action spaces.**
 - Real problems have complex control interfaces.
- **Safety constraints that should never or at least rarely be violated.**
 - Real systems can easily destroy themselves and their surroundings.
- **Tasks that may be partially observable, non-stationary or stochastic.**
 - Need to be able robust in the face of these perturbations.
- **Reward functions that are unspecified, multi-objective, or risk-sensitive.**
 - Explaining a task is hard, especially to a computer.
- **System operators who desire explainable policies and actions.**
 - Otherwise no one will run your controller.
- **Inference that must happen in real-time at the control frequency of the system.**
 - Can't use slow, complex controllers.
- **Large and/or unknown delays in the system actuators, sensors, or rewards.**
 - Rewards can come much later, and systems take time to react to actions.

*non-exhaustive

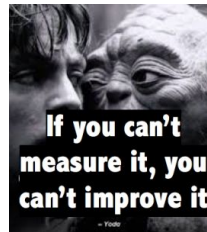
A formalism for RWRL experiments

Formally describing a goal, and measuring your distance from it, is necessary to attain it.

The goal of this work is to define a common language for describing aspects of RL problems that are not clearly enveloped by just an MDP + Reward function formalism. This allows both for the algorithm, but also the algorithm designer, to know whether they are moving in the right direction.

Three main points:

- 1. Training Regime**
- 2. Environmental Constraints**
- 3. Evaluation Metrics**



Example Safety Environments

Cart-Pole Variables: x, θ

Type	Constraint
Static	Limit range: $x_l < x < x_r$
Kinematic	Limit velocity near goal: $ \theta_c - \theta > \theta_L \vee \dot{\theta} < \dot{\theta}_V$
Dynamic	Limit cart acceleration: $\ddot{x} < A_{\max}$

Walker Variables: θ, u, F

Type	Constraint
Static	Limit joint angles: $\theta_L < \theta < \theta_U$ Enforce uprightiness: $0 < u_x$
Kinematic	Limit joint velocities: $\max_i \dot{\theta}_i < L_{\dot{\theta}}$
Dynamic	Limit foot contact forces: $F_{\text{foot}} < F_{\max}$

Manipulator Variables: θ, F, \mathcal{M}

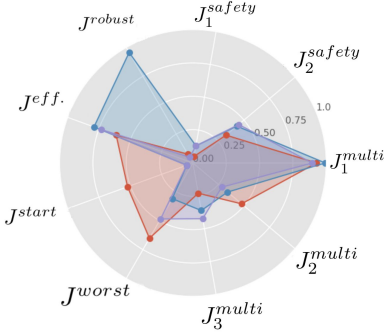
Type	Constraint
Static	Limit joint angles $\theta_L < \theta < \theta_U$ Avoid dynamic obstacles $\mathcal{M} \cap \mathcal{M}_{O,i} = \emptyset$ Avoid self-contact $\mathcal{M} \cap \mathcal{M} = \mathcal{M}$
Kinematic	Limit joint velocities: $\max_i \dot{\theta}_i < L_{\dot{\theta}}$
Dynamic	Acceleration Limits: $\max_i \ddot{\theta}_i < L_{\ddot{\theta}}$ Limit end effector forces: $F_{EE} < F_{\max}$

Humanoid Variables: θ, u, F

Type	Constraint
Static	Limit joint angles: $\theta_{L,i} < \theta_i < \theta_{U,i}$ Enforce uprightiness: $0 < u_x$
Kinematic	Limit joint velocities: $\max_i \dot{\theta}_i < L_{\dot{\theta}}$
Dynamic	Limit foot contact forces. $F_{\text{Foot}} < F_{\max}$ Encourage falls on posterior $F_i < F_{\max,1} \forall i \in \mathcal{C} \setminus i_{\text{post}}$ $F_{\text{post}} < F_{\max,2}$

Noise & Non-Stationarity

Env.	Noise	Non-Stationarity
Cart-Pole	Actuator and sensor delays	Track friction increasing with time
Walker	Noisy perception of terrain	Occasionally non-responsive leg actuator
Manipulator	Imprecise proprioception	Changes in gripper friction
Humanoid	Reduced torque on leg actuator	Varying payload CoGs

<h3>Training Regime</h3>	<ul style="list-style-type: none"> • Offline: Policies are learned on fixed datasets of observations from the environment. • This data can be either random, safe & mediocre, or “optimal”. • Train as long as you want on the data, and then deploy on the system. 												
<h3>Env. Constraints</h3>	<ul style="list-style-type: none"> • Safety: <ul style="list-style-type: none"> ○ Static, Kinematic & Dynamic constraints - non-linear constraints. ○ Inescapable attractors - this is hard, likely described by static box constraints. ○ Stay near demonstrator. ○ Can have a fallback controller, want to minimize falling back on it. • Non-stationarity & environmental shift: <ul style="list-style-type: none"> ○ Physics constants changing over time (friction, torques, available force). ○ Same, but between logs and reality (‘sim2real’ shift). 												
<h3>Evaluation Metrics</h3>	<table border="1" data-bbox="428 594 1136 952"> <thead> <tr> <th>Challenge</th> <th>Evaluator</th> </tr> </thead> <tbody> <tr> <td>Off-line</td> <td>$J^{start} = R(\text{Train}(D_{\pi_B}))$</td> </tr> <tr> <td>Efficient</td> <td>$J^{eff.} = \min D_i \text{ s.t. } R(\text{Train}(D_i)) > R_{\min}$</td> </tr> <tr> <td>Safe</td> <td>$J^{safety}(\pi) = \left(\sum_{i=1}^T c_j(s_i, a_i) \right)_{1 \leq j \leq K} \in \mathbb{R}^K$</td> </tr> <tr> <td>Robust</td> <td>$J^{robust}(\pi) = \frac{1}{K} \sum_{p \in \mathbf{P}} \mathbf{E}^p \left[\sum_{i=1}^T r(s_i, a_i) \right]$</td> </tr> <tr> <td>Discerning</td> <td>$J^{multi}(\pi) = \left(\sum_{i=1}^{T_n} r_j(s_i, a_i) \right)_{1 \leq j \leq K} \in \mathbb{R}^K$</td> </tr> </tbody> </table>  <p>The radar chart displays eight performance metrics on a circular scale from 0.00 to 1.00. The metrics are: J_{robust}, J_1^{safety}, J_2^{safety}, J_1^{multi}, J_2^{multi}, J_3^{multi}, J_{worst}, and J_{start}. The chart shows two overlapping data series: a blue one and a red one. The red series generally shows higher values for J_{robust}, J_1^{safety}, and J_2^{safety}, while the blue series shows higher values for J_1^{multi}, J_2^{multi}, and J_3^{multi}. The J_{start} and $J_{eff.}$ metrics are also present but not explicitly labeled with values on the chart.</p>	Challenge	Evaluator	Off-line	$J^{start} = R(\text{Train}(D_{\pi_B}))$	Efficient	$J^{eff.} = \min D_i \text{ s.t. } R(\text{Train}(D_i)) > R_{\min}$	Safe	$J^{safety}(\pi) = \left(\sum_{i=1}^T c_j(s_i, a_i) \right)_{1 \leq j \leq K} \in \mathbb{R}^K$	Robust	$J^{robust}(\pi) = \frac{1}{K} \sum_{p \in \mathbf{P}} \mathbf{E}^p \left[\sum_{i=1}^T r(s_i, a_i) \right]$	Discerning	$J^{multi}(\pi) = \left(\sum_{i=1}^{T_n} r_j(s_i, a_i) \right)_{1 \leq j \leq K} \in \mathbb{R}^K$
Challenge	Evaluator												
Off-line	$J^{start} = R(\text{Train}(D_{\pi_B}))$												
Efficient	$J^{eff.} = \min D_i \text{ s.t. } R(\text{Train}(D_i)) > R_{\min}$												
Safe	$J^{safety}(\pi) = \left(\sum_{i=1}^T c_j(s_i, a_i) \right)_{1 \leq j \leq K} \in \mathbb{R}^K$												
Robust	$J^{robust}(\pi) = \frac{1}{K} \sum_{p \in \mathbf{P}} \mathbf{E}^p \left[\sum_{i=1}^T r(s_i, a_i) \right]$												
Discerning	$J^{multi}(\pi) = \left(\sum_{i=1}^{T_n} r_j(s_i, a_i) \right)_{1 \leq j \leq K} \in \mathbb{R}^K$												

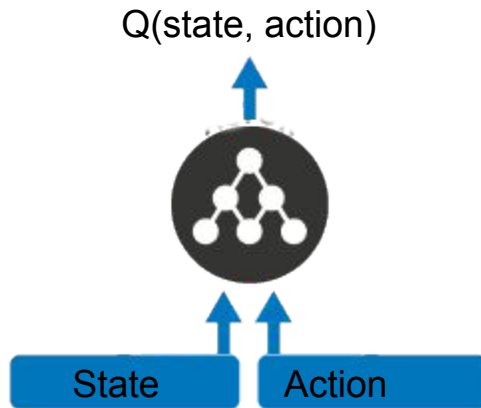
Deep Reinforcement Learning in Large Action Spaces

Gabriel Dulac-Arnold, Richard Evans, Hado van Hasselt,
Peter Sunehag, Jon Hunt, Timothy Lillicrap, Timothy
Mann, Thomas Degris, Theophane Weber



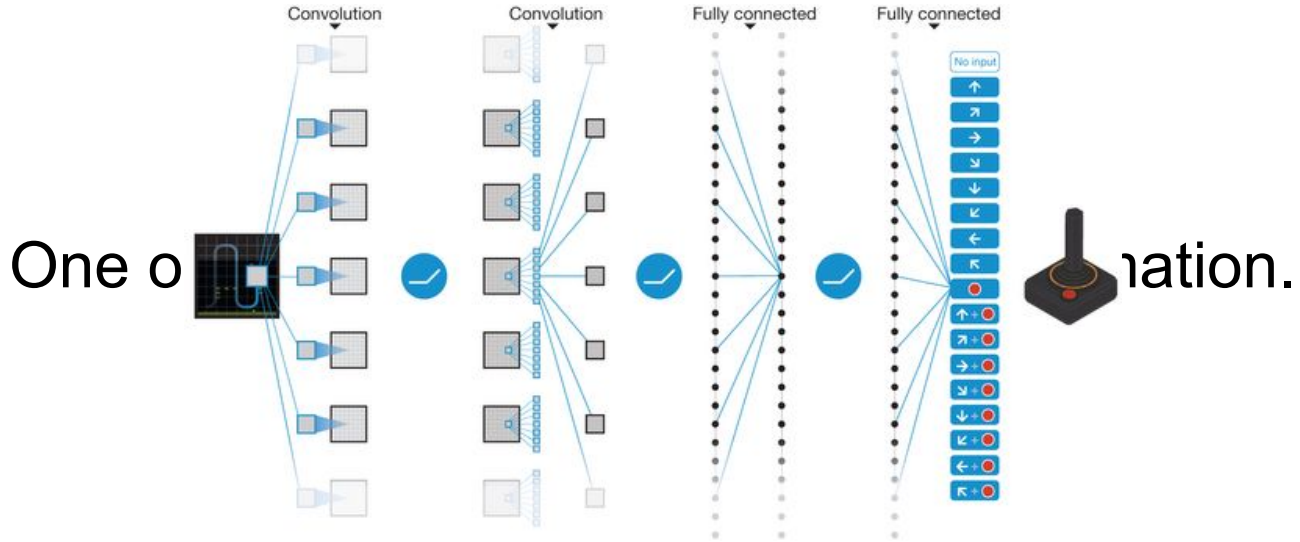
Q-networks don't scale

Need to be evaluated for every action.



Network is likely a deep net (expensive).

DQN Nets scale a bit more, but generalize less



Issues with current approaches

Value-based policy architectures:

- Require explicit argmax of costly function: **Bad**
- Generalize over actions : **Good**

Actor-based policy architectures:

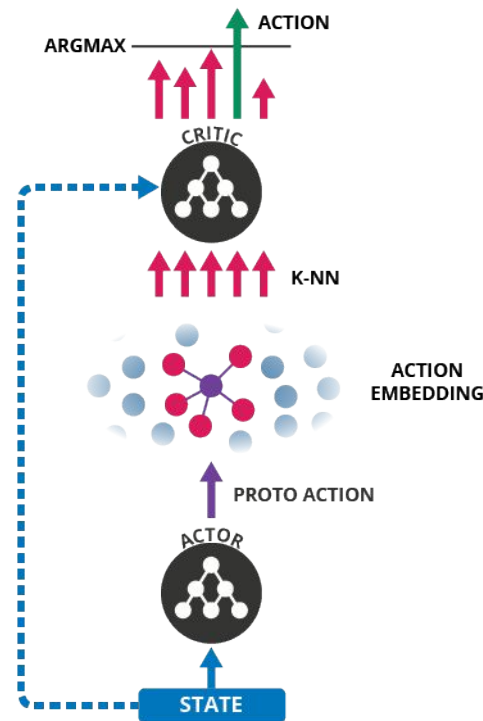
- Avoids explicit argmax on costly function: **Good**
- Most don't generalize over *discrete* actions : **Bad**

Can we have both?

Solution: Wolpertinger Architecture



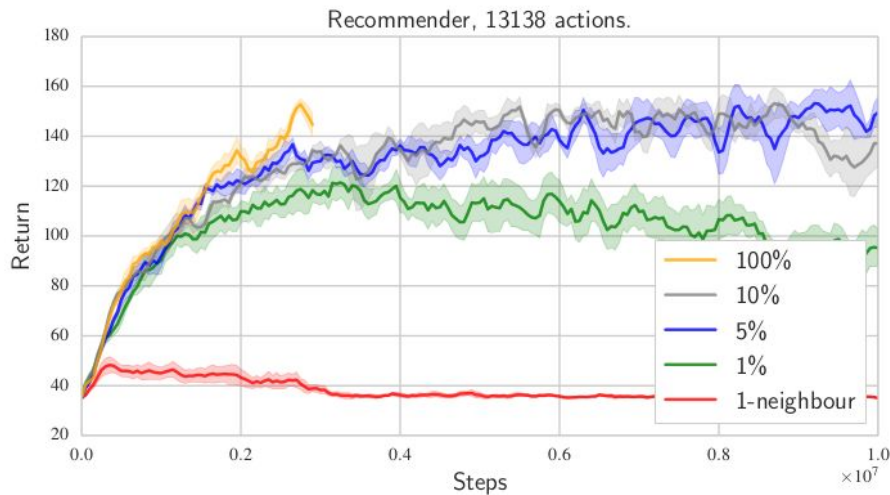
1. **Actions in latent space.**
2. **Learn a continuous-control actor that maps a given state to a point in this space.**
3. **Approximate nearest-neighbor finds K closest valid actions.**
4. **Take $\text{argmax}[Q(s,a)]$ over this K-set.**



Results



13,138-action recommender task w/ 200-dimensional action features.



But can be slow...

Deep Q-Learning from Demonstrations

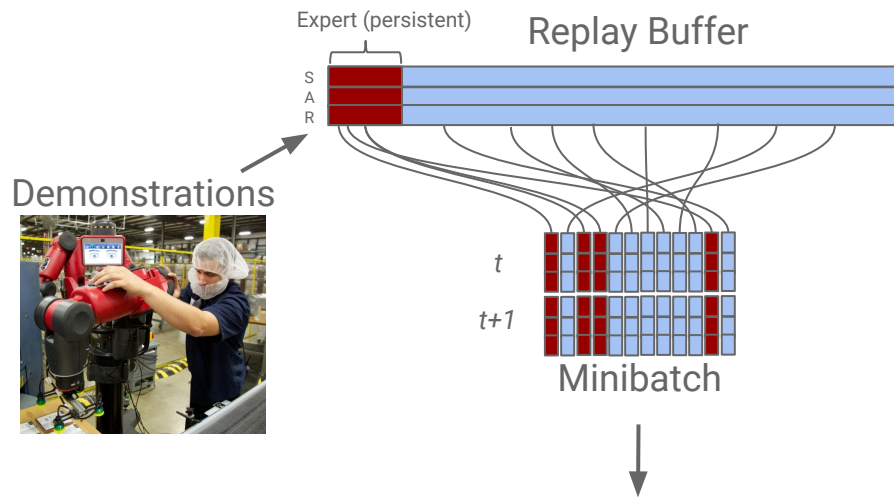
Todd Hester, Matej Vecerik, Olivier Pietquin, Marc Lanctot, Tom Schaul, Bilal Piot, Dan Horgan,
John Quan, Andrew Sendonaris, Ian Osband, Gabriel Dulac-Arnold, John Agapiou, Joel Z. Leibo,
Audrunas Gruslys



Goals

- Start out with good performance, rather than random action selection
- Learn to out-perform the demonstrations using the task rewards
- Use demonstrations to solve the exploration problem in hard exploration tasks, out-performing baseline deep RL approaches

Deep Q-Learning from Demonstrations



- Prioritized replay sampling
- Priority based on TD error + bonus for demonstration samples

$$J(Q) = J_{DQ}(Q) + \lambda_1 J_n(Q) + \lambda_2 J_E(Q) + \lambda_3 J_{L2}(Q)$$

Losses

- We train the network with a combination of four losses:

$$J(Q) = J_{DQ}(Q) + \lambda_1 J_n(Q) + \lambda_2 J_E(Q) + \lambda_3 J_{L2}(Q)$$

- Double Q-Learning 1-step: $r_{t+1} + \gamma Q_{\theta_t^-}(s_{t+1}, \operatorname{argmax}_{a \in \mathcal{A}} Q_{\theta_t}(s_{t+1}, a))$
- Double Q-Learning n-step (n=10): $r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \max_{a \in \mathcal{A}} \gamma^n Q(s_{t+n}, a)$
- Supervised Classification of Demonstrator Actions
- L2 Weight Regularization: $|\theta|_2$

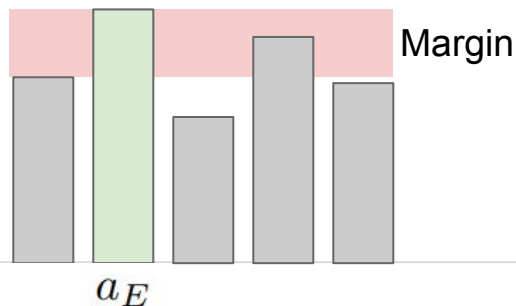
Supervised Loss

Add supervised loss using large margin classification

$$J_E(Q) = \max_{a \in A} [Q(s, a) + l(a_E, a)] - Q(s, a_E)$$

Where $l(a_E, a)$ is 0 if $a = a_E$, and positive otherwise

Forces value of expert action to be some margin higher than all other values.



Pre-Training Phase

- Agent is trained solely from demonstration data without any environment interaction
- Enable agent to start out acting in environment much better than random
- All four losses are critical for pre-training
 - Q-Learning only: many actions were never taken and will have random values
 - Supervised only: Will not learn a real value function to learn from

Comparison Algorithms

- Replay Buffer Spiking (RBS) (Lipton et al. 2016)
 - Pre-fill replay buffer with demonstrations, eventually overwritten
- Human Experience Replay (HER) (Hosu and Rebedea 2016)
 - Keep demonstration data, sample from mixed replay buffer
- Accelerated DQN with Expert Trajectories (ADET) (Lakshminarayanan, Ozair, and Bengio 2016)
 - Use cross-entropy supervised loss. No pre-training.

Goals

- Start out with good performance, rather than random action selection
- Learn to out-perform the demonstrations using the task rewards
- Use demonstrations to solve the exploration problem in hard exploration tasks, out-performing baseline deep RL approaches

General Results

- Start out with good performance
 - Starts better than PDD DQN on 41/42 games
 - Starts better than imitation on 31/42 games
 - TD loss helps even without additional interactions
- Learn to out-perform the demonstrations
 - Higher score than worst demonstration on 29/42 games
 - Higher score than best demonstration on 14/42 games

State of the Art for Deep RL on 11 games

Game	DQfD	Prev. Best	Algorithm
Alien	4745.9	4461.4	Dueling DQN (Wang et al. 2016)
Asteroids	3796.4	2869.3	PopArt (van Hasselt et al. 2016)
Atlantis	920213.9	395762.0	Prior. Dueling DQN (Wang et al. 2016)
Battle Zone	41971.7	37150.0	Dueling DQN (Wang et al. 2016)
Gravitar	1693.2	859.1	DQN+PixelCNN (Ostrovski et al. 2017)
Hero	105929.4	23037.7	Prioritized DQN (Schaul et al. 2016)
Montezuma Revenge	4739.6	3705.5	DQN+CTS (Ostrovski et al. 2017)
Pitfall	50.8	0.0	Prior. Dueling DQN (Wang et al. 2016)
Private Eye	40908.2	15806.5	DQN+PixelCNN (Ostrovski et al. 2017)
Q-Bert	21792.7	19220.3	Dueling DQN (Wang et al. 2016)
Up N Down	82555.0	44939.6	Dueling DQN (Wang et al. 2016)

DQfD Results

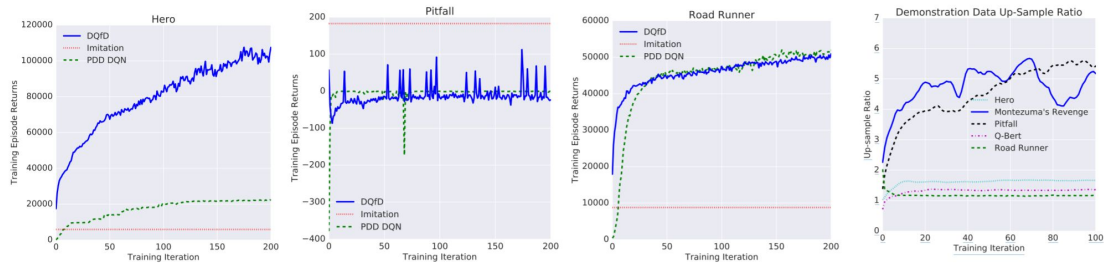


Figure 1: On-line scores of the algorithms on the games of Hero, Pitfall, and Road Runner. On Hero and Pitfall, DQfD leverages the human demonstrations to achieve a higher score than any previously published result. The last plot shows how much more frequently the demonstration data was sampled than if data were sampled uniformly, for five different games.



Figure 2: The left plots show on-line rewards of DQfD with some losses removed on the games of Montezuma's Revenge and Q-Bert. Removing either loss degrades the performance of the algorithm. The right plots compare DQfD with three algorithms from the related work section. The other approaches do not perform as well as DQfD, particularly on Montezuma's Revenge.

Montezuma's Revenge



Extension to DDPG for continuous control



DDPGfD (Scholz et al.)



Model-Based RL for Real-World RL



Some solutions to these problems...

Model-Based RL:

Deep Reinforcement Learning in a Handful of Trials using Probabilistic Dynamics Models

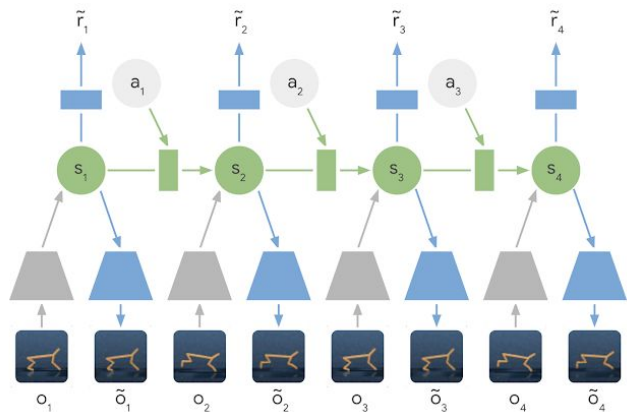
Learning Latent Dynamics for Planning from Pixels (PlaNet paper, already pinged by Todd and what I'm currently trying to train on Mujoco tasks)

Environmental robustness:

LEARNING TO ADAPT IN DYNAMIC, REAL-WORLD ENVIRONMENTS THROUGH META-REINFORCEMENT LEARNING

Safety:

Safe Exploration in Continuous Action Spaces

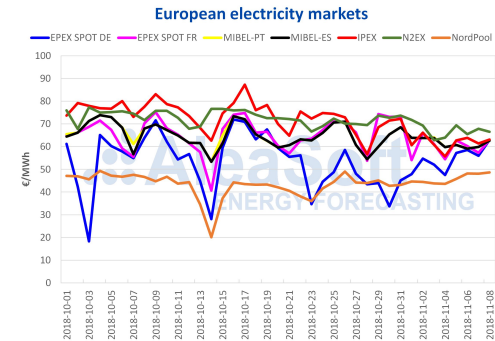


MBRL: Off-line, Efficient, Stable, Explainable

- Can learn from logged data w/out environmental interactions
 - Data-efficient learning of environment dynamics
 - Underlying learning problem is supervised and generally more stable
 - Allow for something closer to symbolic planning
 - Explainable
 - Easier to express safety constraints
 - Can integrate discrete & changing constraints & reward
 - Can be more easily controlled by a hierarchy
 - Meta-learning can allow for better domain adaptation/sys-id
- VS. Model-free approach

Example: Datacenters

- Integration of electricity futures
- Integration of scheduled maintenance
- Long-term goal of MPC plan is visible
 - Can potentially see which beliefs were mistaken
- Easier to enforce safety constraints



Example: Mobile Robots & Manipulation

- More data efficient
- Easier to create control abstractions
- Can reason about safety without experience
- Can deal with infinite number of policy specifications
- Can reason generic cost (energy & time minimization) and only goal state.
- Can deal with low-quality sensors (filtering)
- Can deal with low-quality actuators (visual closed loop)



What needs to be tried & what's missing

- Current situation
 - Good model learning from pixels (PlaNet, etc.)
 - Sub-optimal inefficient planning (CEM)
 - Aleatoric uncertainty
- Next steps
 - Better planning
 - Gradients?
 - Epistemic uncertainty - know when to trust yourself
 - Combining model-free & model-based (policy caching)
 - Trying these out on real systems
 - Adding dynamic constraints on plan
 - Goal-defined policies
 - ????