# The Predictron: End-to-end Learning and Planning

Yoonho Lee

Department of Computer Science and Engineering
Pohang University of Science and Technology

December 27, 2016

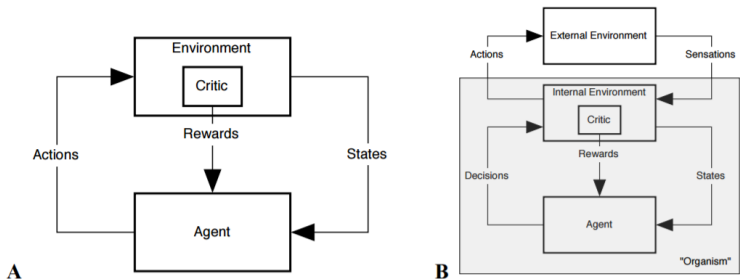# Hierarchical RL

motivation

# Hierarchical RL



Figure 1: *Agent-Environment Interaction in RL.* **A**: *The usual view.* **B**: *An elaboration.*
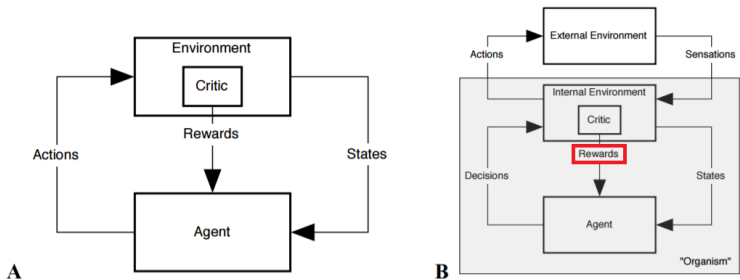
# Hierarchical RL



Figure 1: *Agent-Environment Interaction in RL.* **A**: *The usual view.* **B**: *An elaboration.*
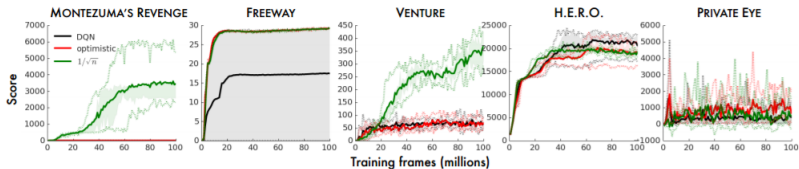
# Reward augmentation[1]



Figure 2: Average training score with and without exploration bonus or optimistic initialization in 5 Atari 2600 games. Shaded areas denote inter-quartile range, dotted lines show min/max scores.
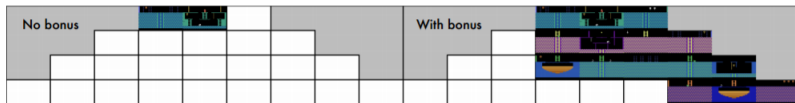


Figure 3: "Known world" of a DQN agent trained for 50 million frames with (**right**) and without (**left**) count-based exploration bonuses, in MONTEZUMA'S REVENGE.

[1]Bellemare et al. Unifying Count-Based Exploration and Intrinsic Motivation, NIPS 2016
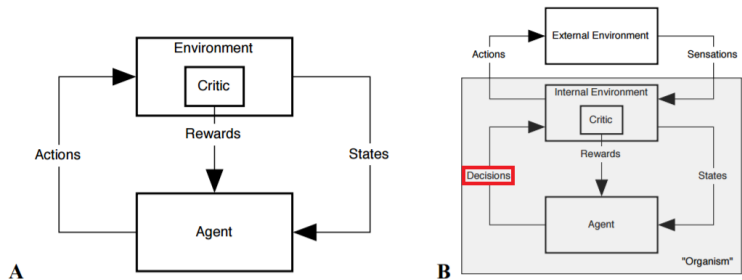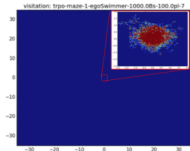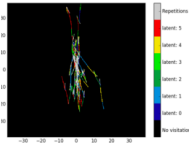
# Hierarchical RL



Figure 1: *Agent-Environment Interaction in RL.* **A**: *The usual view.* **B**: *An elaboration.*
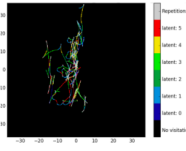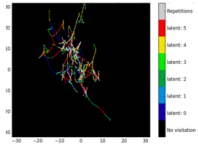
# Hierarchical actions[2]



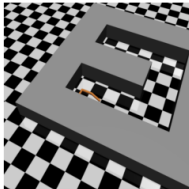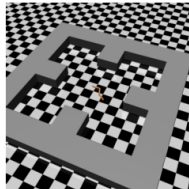(a) Gaussian noise with covariance $\Sigma = I$ (b) Hierarchichy with Multi-policy (c) Hierarchy with Bilinear-SNN,($\alpha_H = 0$) (d) Hierarchy with Bilinear-SNN, $\alpha_H = 0.01$
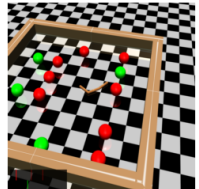


(a) Maze 0    (b) Maze 1    (c) Maze 2 or 3    (d) Food Gather

[2]Florensa et al. Stochastic Neural Networks for Hierarchical Reinforcement Learning, under review for ICLR 2017

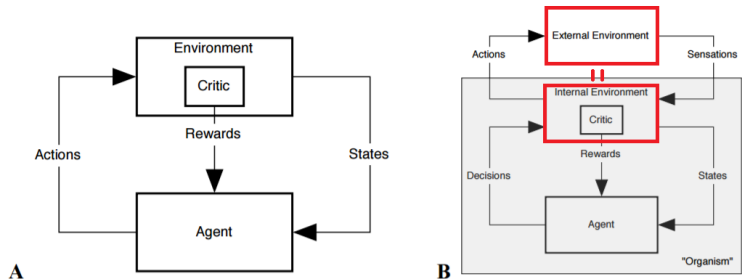# Hierarchical RL



Figure 1: *Agent-Environment Interaction in RL.* **A**: *The usual view.* **B**: *An elaboration.*

# Model-based video prediction[3]



[3]Oh et al. Action-Conditional Video Prediction using Deep Networks in Atari Games, NIPS 2015

# Value Iteration Networks

Aviv Tamar, Yi Wu, Garrett Thomas, Sergey Levine, Pieter Abbeel
Best paper at NIPS 2016

# Value Iteration Networks

# Value Iteration Networks

Network diagram

# Value Iteration Networks

## Bellman Operator

Every $V$ converges to $V^*$ under the Bellman Operator $T$ defined as:

$$(TV)(s) = \max_{a \in A}\{R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a)V(s')\} \qquad (1)$$

$$V^* = \lim_{n \to \infty} T^n V \quad \forall V \qquad (2)$$

The Bellman Operator can be viewed as a CNN forward pass:

$$(T^*V)(s) = \max_{a \in A}\{R(s, a) + \gamma conv_{P(a)}(V(s))\}$$

# Value Iteration Networks

VI Module

# Value Iteration Networks

Results



| Network | Train Error | Test Error |
|---------|-------------|------------|
| VIN     | 0.30        | 0.35       |
| CNN     | 0.19        | 0.73       |

Figure 4: Continuous control domain. Left: average distance to goal on training and test domains for VIN and CNN policies. Right: trajectories predicted by VI and CNN on random test domains.

# Value Iteration Networks

Summary

- Neural network architecture that plans using value iteration
- Assumes that the state is a sufficient statistic for the reward function
- The small MDP must have a finite state space
- Uses prior knowledge about the environment's structure

# The Predictron: End-to-End Learning and Planning

David Silver, Hado van Hasselt, Matteo Hessel, Tom Schaul,
Arthur Guez, Tim Harley, Gabriel Dulac-Arnold, David Reichert,
Neil Rabinowitz, Andre Barreto, Thomas Degris
Under review for ICLR 2017

# Predictron



$\vdots$      $\vdots$      $\vdots$

$$\mathbf{r}^2, \mathbf{s}^2, \gamma^2 \xrightarrow{v} \mathbf{v}^2 \qquad \mathbf{g}^2 = \mathbf{r}^1 + \gamma^1 \mathbf{r}^2 + \gamma^1 \gamma^2 \mathbf{v}^2 \qquad \mathbb{E}^m \left[\mathbf{g}^2\right]$$

$\uparrow m$      $\parallel$

$$\mathbf{r}^1, \mathbf{s}^1, \gamma^1 \xrightarrow{v} \mathbf{v}^1 \qquad \mathbf{g}^1 = \mathbf{r}^1 + \gamma^1 \mathbf{v}^1 \qquad \mathbb{E}^m \left[\mathbf{g}^1\right]$$

$\uparrow m$      $\parallel$

$$\mathbf{s}^0 \xrightarrow{v} \mathbf{v}^0 \qquad \mathbf{g}^0 = \mathbf{v}^0 \qquad \mathbb{E}^m \left[\mathbf{g}^0\right]$$

$\uparrow m$      $\parallel$

$\uparrow f$      $\mathbb{E}_p \left[\mathbf{g}_t\right]$

$s$

$(a)$      $(b)$      $(c)$

# Predictron



$(a)$ $(d)$

# Predictron



Figure 2: **Left:** Two sample mazes from the random-maze domain. Light blue cells are empty, darker blue cells contain a wall. One maze is connected from top-left to bottom-right (indicated in black), the other is not. **Right:** An example trajectory in the pool domain (before downsampling). It was selected by maximising the prediction of pocketing balls, using the predictron.

# Predictron



Figure 4: **Comparing predictron to baselines.** Aggregated prediction errors on random mazes (top) and pool (bottom) over all predictions for the eight architectures corresponding to the cube on the left. Each line is the median of RMSE over five seeds; shaded regions encompass all seeds. The full $(r, \gamma, \lambda)$-predictron (**red**), consistently outperformed conventional deep network architectures (**black**), with and without skips or and with and without weight sharing.

# Predictron



$$g^2 = r^1 + \gamma^1 r^2 + \gamma^1 \gamma^2 v^2$$

$$g^1 = r^1 + \gamma^1 v^1$$

$$g^0 = v^0$$

$(d)$

$(b)$

# Predictron



Figure 5: **Semi-supervised learning.** Prediction errors of the $(r, \gamma, \lambda)$-predictrons (shared core, no skips) using 0, 1, or 9 consistency updates for every update with labelled data, plotted as function of the amount of labelled data consumed. Early learning performance improved with more consistency updates.
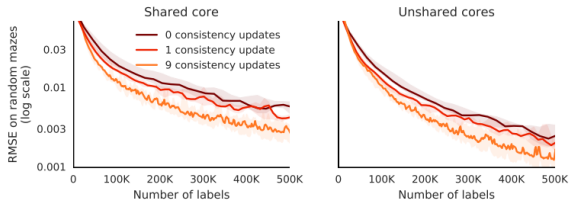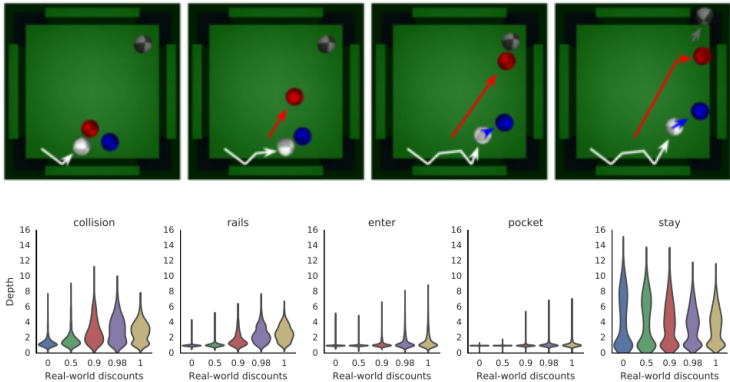
# Predictron



Figure 6: **Thinking depth.** Distributions of thinking depth on pool for different types of predictions and for different real-world discounts. Depth is defined as $\boldsymbol{\lambda}^0(1 + \boldsymbol{\gamma}^1 \boldsymbol{\lambda}^1(1 + \boldsymbol{\gamma}^2 \boldsymbol{\lambda}^2(1 + \ldots)))$.

# Predictron

Summary

- End-to-end simulation of an MRP
- Works with arbitrary state space for small MRP
- Small MRP has a non-interpretable state space
- Designed for RL, but does not take actions into account

# Summary

- Hierarchical Reinforcement Learning attempts to identify crucial decisions
- Agents can now use NN-based planning for better decision making

- Research directions
  - Theoretical bounds for optimal policy of a smaller MDP
  - Learning a smaller MDP with abstract actions
  - End-to-end planning based Q network or policy network

# Thank You

Value Iteration Networks
The Predictron: End-to-End Learning and Planning