

# Training Agent using Upside-Down Reinforcement Learning (74)

Rupesh Kumar Srivastava<sup>1</sup>, Pranav Shyam<sup>2</sup>, Filipe Mutz<sup>3</sup>, Wojciech Jaśkowski<sup>4</sup>, Jürgen Schmidhuber<sup>1, 5</sup>

<sup>1</sup>NNAISENSE, <sup>2</sup>Google DeepMind

<sup>3</sup>Federal University of Espírito Santo (UFES), Brazil

<sup>4</sup>Snowflake, <sup>5</sup>The Swiss AI Lab IDSIA

NeurIPS Deep Reinforcement Learning Workshop 2019

2025.05.22.

김동민

# Lifetime credit assignment problem solver

- Jürgen Schmidhuber

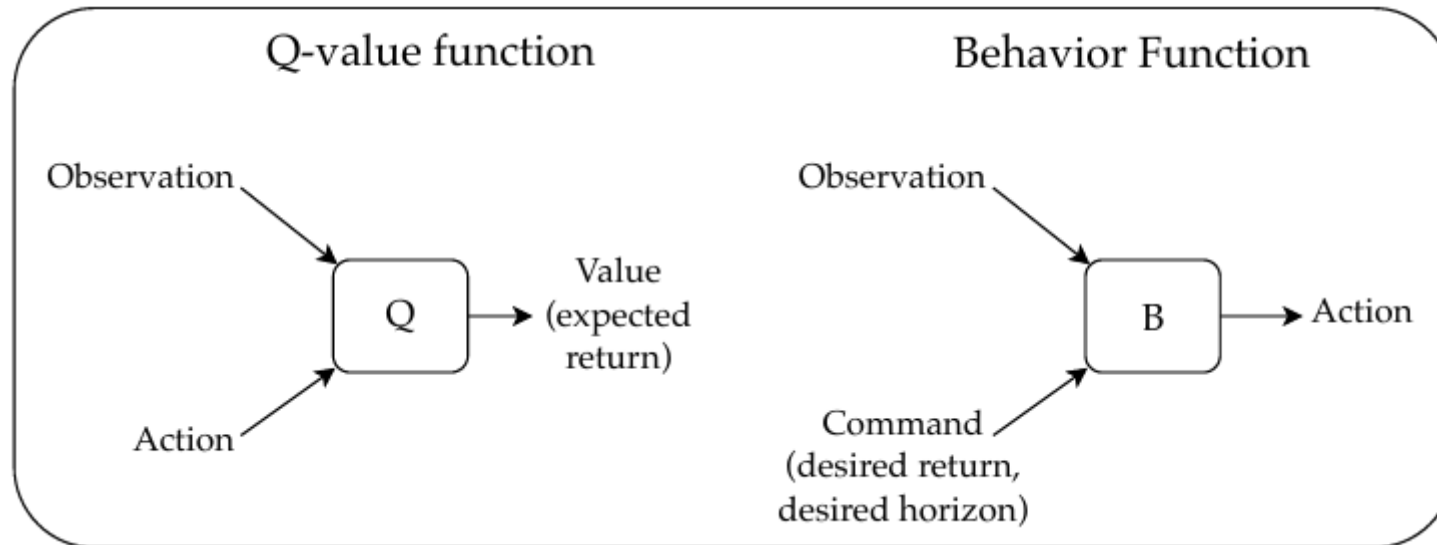
- LSTM (Long Short-Term Memory)
- Convolutional Neural Network(?)
- Generative Adversarial Networks (?)
- Transformer(?)



- 이 논문에서도 CNN을 사용하는데 LeCun을 인용하지 않았음 (NN 관련 Hinton의 연구도 1개만 인용)

# In Short

- 전통적인 RL문제들을 Supervised Learning (SL) 방법만 사용해서 풀어보자
  - To bring the simplicity, robustness and scalability of SL algorithms to traditional RL algorithms
- Key Idea: Use Upside-down function (Behavior function) instead of Q-value function



**SL Problem for behavior function  $B_{\mathcal{T}}$**

$$B_{\mathcal{T}} = \arg \min_B \sum_{(\tau, t_1, t_2)} L(B(a_{t_1}, s_{t_1}, d^r, d^h), a_{t_1})$$

where  $0 \leq t_1 \leq t_2 \leq \text{len}(\tau)$  for  $\tau \in \mathcal{T}$ ,  $d^r = \sum_{t=t_1}^{t_2} r_t$  and  $d^h = t_2 - t_1$

# Upside-Down Reinforcement Learning (UDRL)

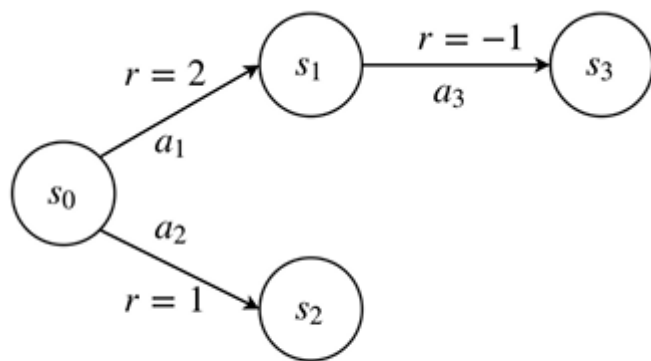
- 다음과 같은 질문에 답하기 위해 UDRL을 고안해냄
  - Is it possible to learn to act in high-dimensional environments efficiently using only SL, avoiding the issues arising from non-stationary learning objectives common in traditional RL algorithms?
- UDRL
  - a new approach to RL where instead of trying to predict rewards from states like value functions and Q-learning, the agents directly predict actions based on the **state** and the **command**
  - The goal of learning is no longer to maximize returns in expectation, but to learn to follow commands
- Command ([desired return, desired time horizon])
  - "achieve total reward R in next T time steps"
  - "reach state S in fewer than T time steps"

- desired return  $d^r = \sum_{t=t_1}^{t_2} r_t$

- desired horizon  $d^h = t_2 - t_1$



# Toy Example

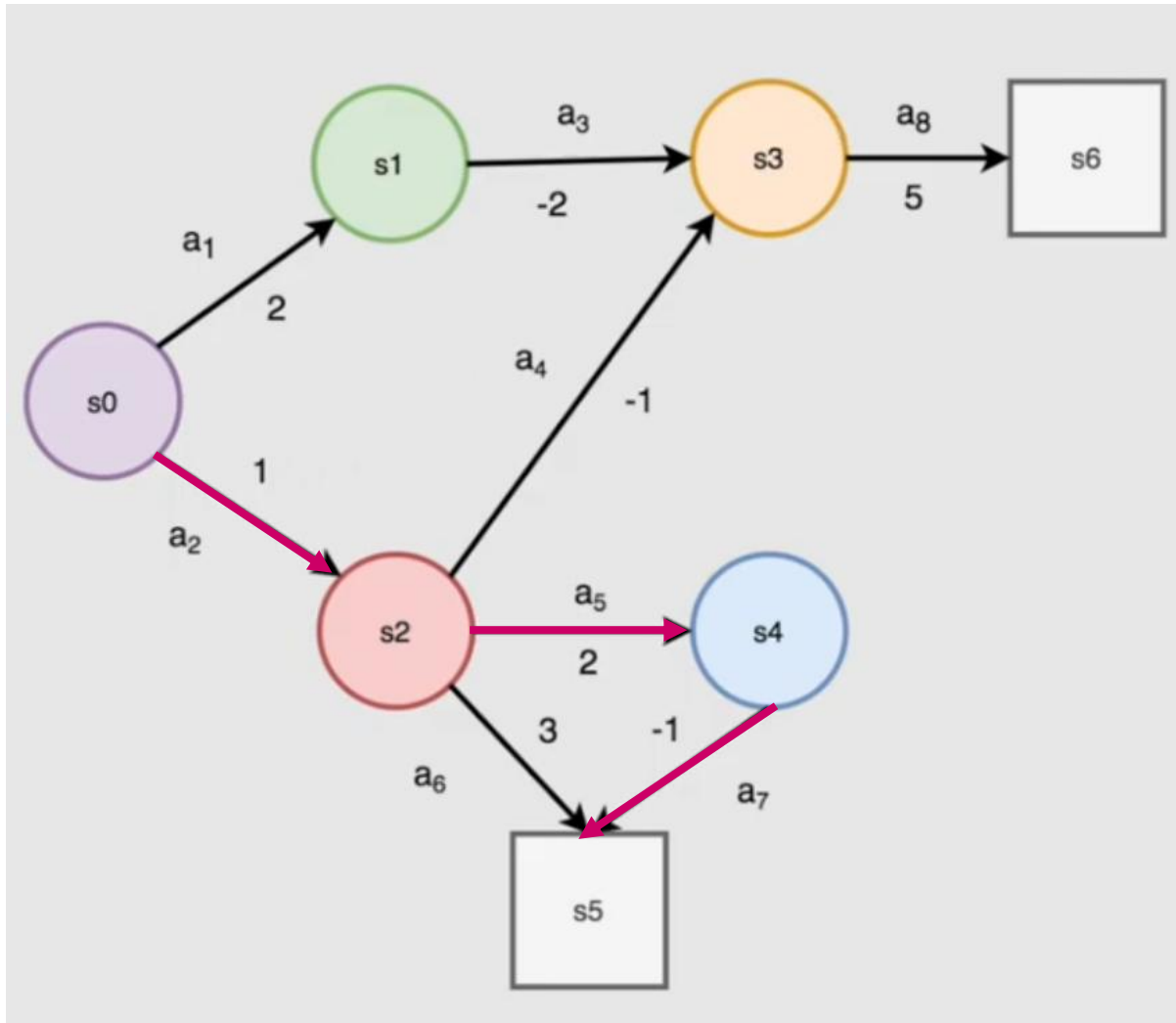


**Figure 2:** A toy environment with four states.

**Table 1.** A behavior function based on all unique trajectories for the toy environment (left image).

| State | Desired Return | Desired Horizon | Action |
|-------|----------------|-----------------|--------|
| $s_0$ | 2              | 1               | $a_1$  |
| $s_0$ | 1              | 1               | $a_2$  |
| $s_0$ | 1              | 2               | $a_1$  |
| $s_1$ | -1             | 1               | $a_3$  |

# Example 2



trajectory of traditional RL

$(s_0, a_2, 1, s_2)$

$(s_2, a_5, 2, s_4)$

$(s_4, a_7, -1, s_5)$

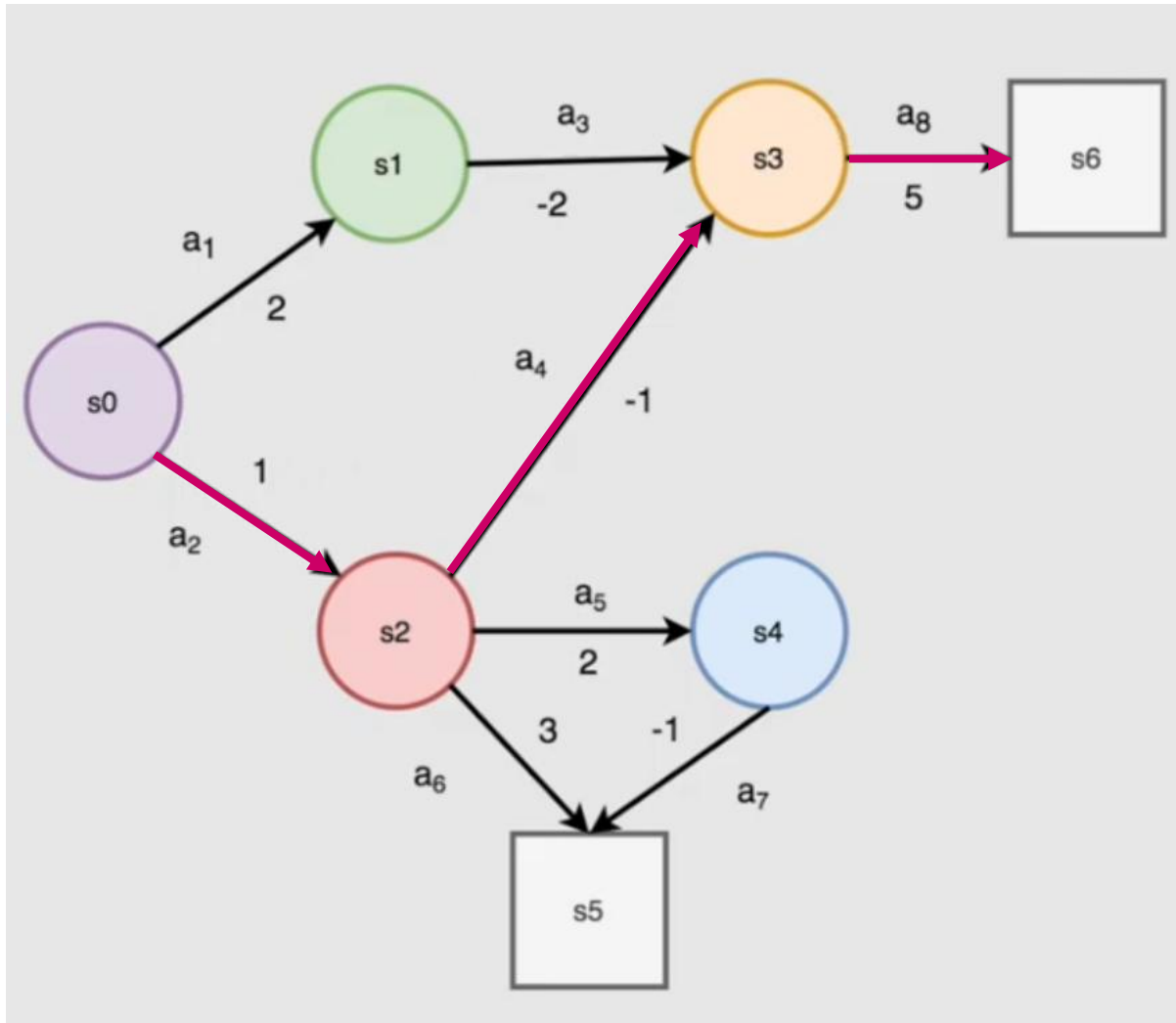
trajectory and SL of UDRL

$L(B(s_0, 2, 3), a_2)$

$L(B(s_2, 1, 2), a_5)$

$L(B(s_4, -1, 1), a_7)$

# Example 2



trajectory of traditional RL

$(s_0, a_2, 1, s_2)$

$(s_2, a_4, -1, s_3)$

$(s_3, a_8, 5, s_6)$

trajectory and SL of UDRL

$L(B(s_0, 5, 3), a_2)$

$L(B(s_2, 4, 2), a_4)$

$L(B(s_3, 5, 1), a_8)$

# UDRL Algorithm 1

- Behavior function은 NN을 function approximator로 활용

---

**Algorithm 1** Upside-Down Reinforcement Learning: High-level Description.

---

- 1: Initialize replay buffer with warm-up episodes using random actions // Section 2.2.1
  - 2: Initialize a behavior function // Section 2.2.2
  - 3: **while** stopping criteria is not reached **do**
  - 4:   Improve the behavior function by training on replay buffer // Section 2.2.3
  - 5:   Sample exploratory initial commands based on replay buffer // Section 2.2.4
  - 6:   Explore using sampled commands in **Algorithm 2** and add to replay buffer // Section 2.2.5
  - 7:   If evaluation is required, evaluate using evaluation commands in **Algorithm 2** // Section 2.2.6
  - 8: **end while**
- 

- Sampling exploratory commands

1. A number of episodes with the highest returns are selected from the replay buffer.
2. The exploratory desired horizon  $d_0^h$  is set to the mean of the lengths of the selected episodes.
3. The exploratory desired returns  $d_0^r$  are sampled from the uniform distribution  $\mathcal{U}[M, M + S]$  where  $M$  is the mean and  $S$  is the standard deviation of the selected episodic returns.

# UDRL Algorithm 2

---

**Algorithm 2** Generates an Episode using the Behavior Function.

---

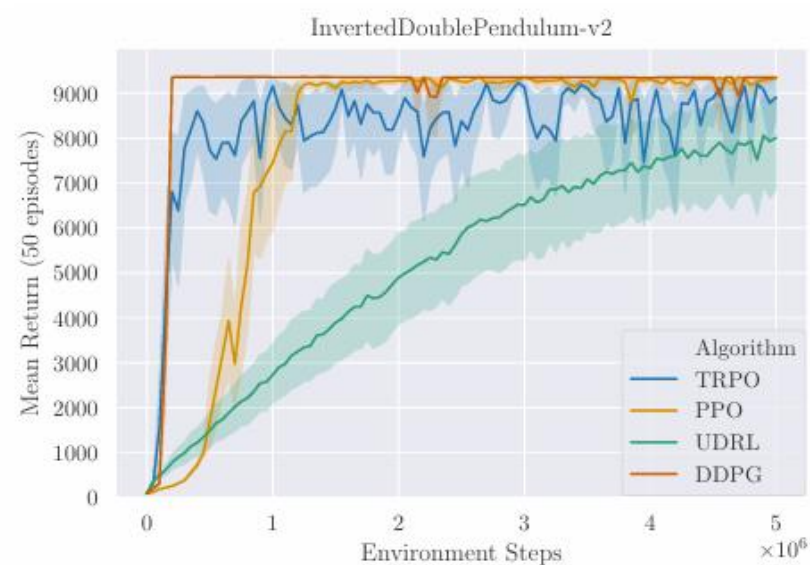
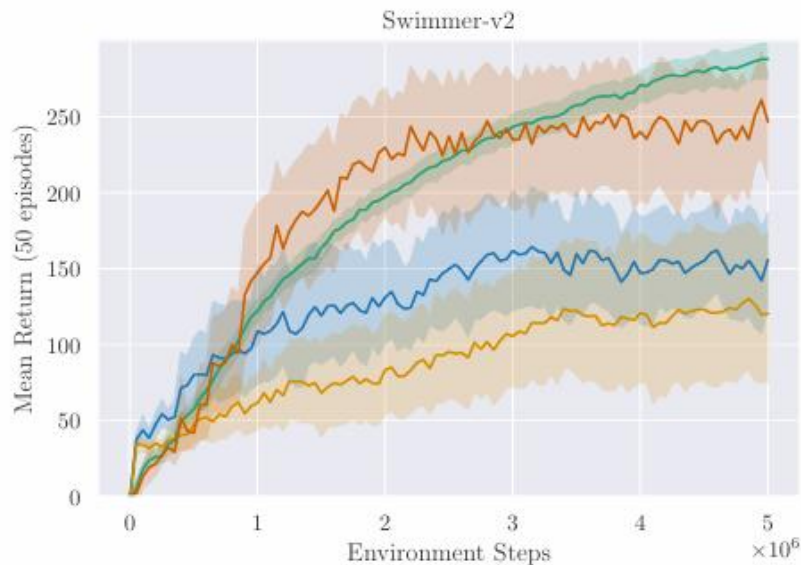
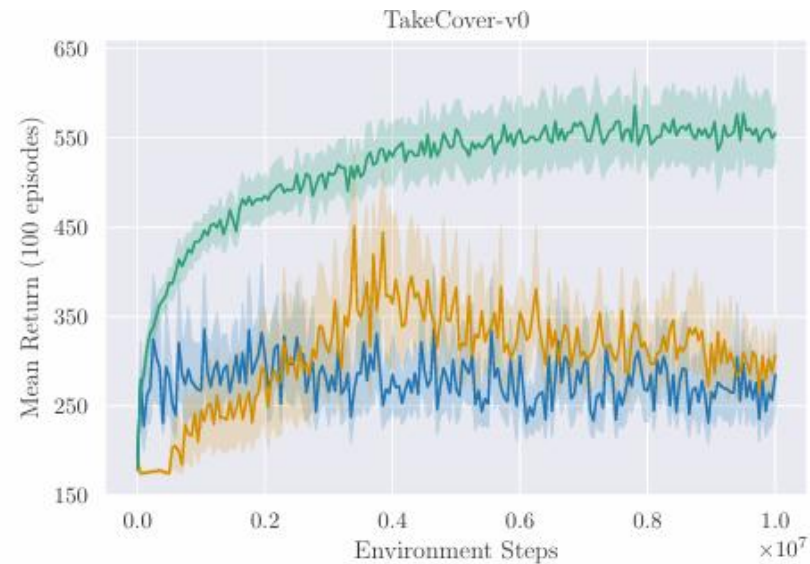
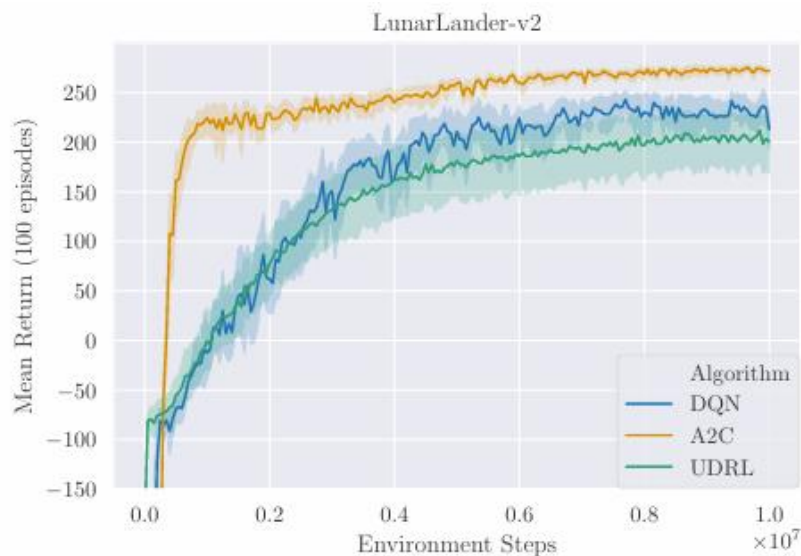
**Input:** Initial command  $c_0 = (d_0^r, d_0^h)$ , Initial state  $s_0$ , Behavior function  $B$  parameterized by  $\theta$

**Output:** Episode data  $E$

- 1:  $E \leftarrow \emptyset$
  - 2:  $t \leftarrow 0$
  - 3: **while** episode is not done **do**
  - 4:   Compute  $P(a_t | s_t, c_t) = B(s_t, c_t; \theta)$
  - 5:   Execute  $a_t \sim P(a_t | s_t, c_t)$  to obtain reward  $r_t$  and next state  $s_{t+1}$  from the environment
  - 6:   Append  $(s_t, a_t, r_t)$  to  $E$
  - 7:    $t \leftarrow t + 1$
  - 8:    $d_t^r \leftarrow d_{t-1}^r - r_{t-1}$  // Update desired return
  - 9:    $d_t^h \leftarrow d_{t-1}^h - 1$  // Update desired horizon
  - 10:    $c_t \leftarrow (d_t^r, d_t^h)$
  - 11: **end while**
-

# Experiments (1)

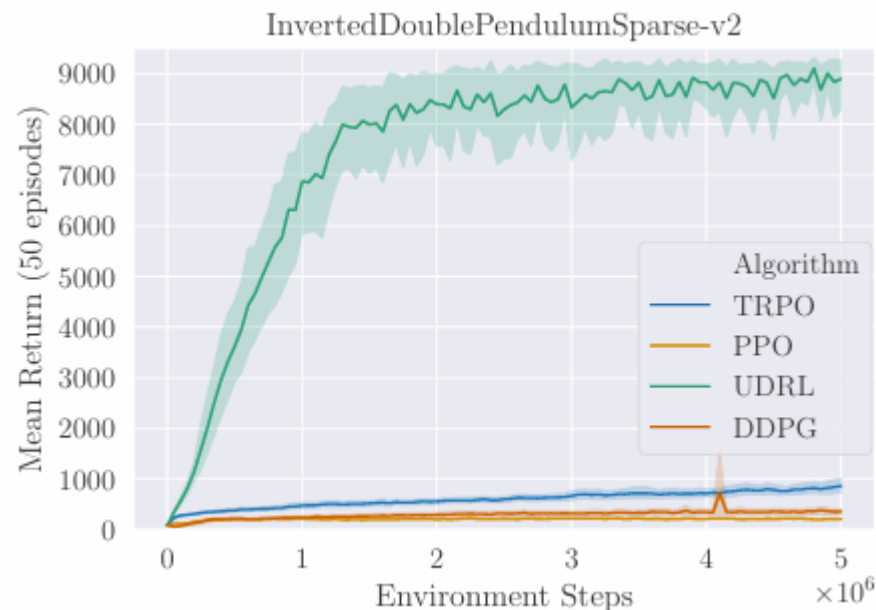
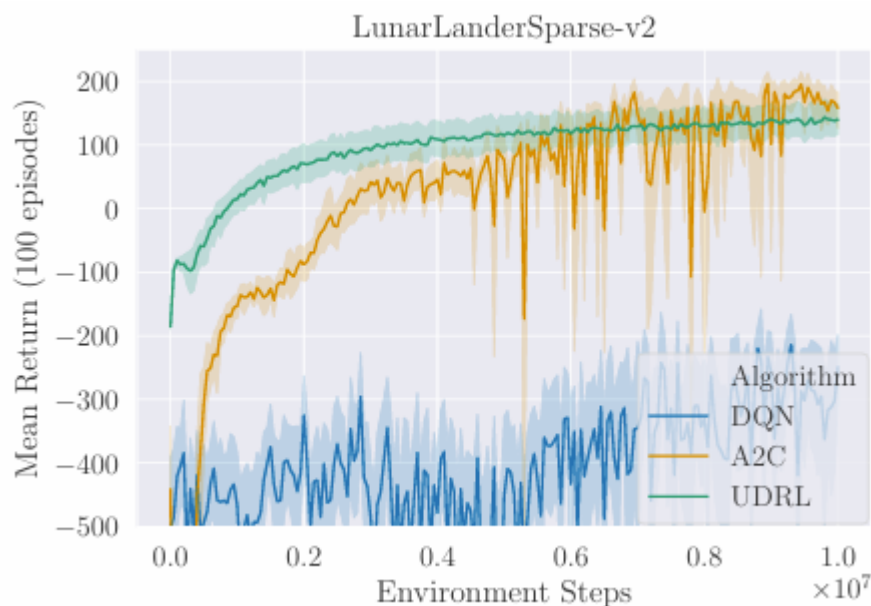
- Discrete action env
  - UDRL vs DQN vs A2C
- Continuous action env
  - UDRL vs TRPO vs PPO vs DDPG



# Experiments (2)

- Results on sparse delayed reward versions

- Sparse delayed reward 환경에서도 UDRL은 별다른 수정 없이 안정적인 성능을 유지, 반면 기존 RL 알고리즘들은 불안정하거나 실패
- UDRL은 설계상 장기적인 보상 할당(credit assignment)에 강건하여, POMDP 상황에서도 Dense reward에서의 PPO 수준의 성능에 근접



# 후일담

- RL문제를 SL을 활용해 기존 RL의 문제를 풀 수 있는 가능성을 보여줌
  - 물론 관련된 더 오래된 연구들이 있으며 논문에서도 자세하게 언급함
  - 그러나 SL을 보조적 수단으로 사용하는 연구들이 많은 것에 비해 이 연구는 SL을 본격적으로 RL문제에 적용하고자 하였음
- Decision transformer, Q-transformer, 더 나아가 VLA 연구에 영감을 줬다고 할 수 있음
- 강화학습연구는 앞으로 어떤 방향으로 나아갈까?
  - 벨만 연구들이 더 많이 나오지 않을까?
    - 벨만 기반 RL은 reward signal 없이 학습이 어렵고, sample efficiency가 낮음
  - SL을 활용하면 대규모 데이터 활용에 더 적합
    - logged dataset 기반으로 pretraining → downstream fine-tuning 가능
  - 학습 안정성과 재현성
    - 벨만 기반은 non-stationarity, off-policy instability 등의 문제가 있을 수 있음
    - supervised-style 학습은 훈련 안정성과 재현성이 높음 (특히 산업계에 중요)
  - 멀티모달 입력 통합
    - VLA(Vision-Language-Action) 모델이나 world model 기반 학습은 복잡한 high-dimensional 입력에서 벨만 방정식 사용이 어렵거나 비효율적
    - transformer 또는 imitation learning이 더 실용적
- 편견이 좀 있었지만 Schmidhuber는 훌륭한 연구자라는 것을 다시 한 번 느낌