SDRL: Symbolic Deep Reinforcement Learning

Liu

Introduction

Backgroun

Method

Experiment

Conclusion and Future Work SDRL: Interpretable and Data-efficient Deep Reinforcement Learning Leveraging Symbolic Planning

Bo Liu

Auburn University, Auburn, AL, USA

#### Collaborators

SDRL: Symbolic Deep Reinforcement Learning

Liu

Introduction Background Method

Experiment

Conclusion and Future Work









Fangkai Yang NVIDIA Corporation Redmond, WA, USA

Steven Gustafson Maana Inc. Bellevue, WA, USA

## Sequential Decision-Making

SDRL: Symbolic Deep Reinforcement Learning

Liu

#### Introduction Background Method

Experiment





- Sequential decision-making (SDM) concerns an agent making a sequence of actions based on its behavior in the environment.
- Deep reinforcement learning (DRL) achieves tremendous success on sequential decision-making problems using deep neural networks (Mnih et al., 2015).

## Challenge: Montezuma's Revenge

SDRL: Symbolic Deep Reinforcement Learning

Introduction Background

Method

Experiment



- The avatar: climbs down the ladder, jumps over a rotating skull, picks up the key (+100), goes back and uses the key to open the right door (+300).
- Vanilla DQN achieves 0 score (Mnih et al., 2015).

## Challenge: Montezuma's Revenge



Introduction

Background

Method

Experiment

Conclusion and Future Work



Problem: long horizon sequential actions, sparse and delayed reward.

- poor data efficiency.
- lack of interpretability.

#### Our Solution

SDRL: Symbolic Deep Reinforcement Learning

Liu

#### Introduction

Background

Method

Experiment

Conclusion and Future Work

#### Solution: task decomposition

- Symbolic planning: subtasks scheduling (high-level plan).
- DRL: subtask learning (low-level control).
- Meta-learner: subtask evaluation.

#### Goal

- Symbolic planning drives learning, improving task-level interpretablility.
- DRL learns feasible subtasks, improving data-efficiency.

# Background: Symbolic Planning with Action Language

SDRL: Symbolic Deep Reinforcement Learning

Liu

Introduction

Background

Method

Experiment

Conclusion and Future Work **Action language** (*Gelfond & Lifschitz, 1998*): a formal, declarative, logic-based language that describes dynamic domains.

Dynamic domains can be represented as a transition system.



## Action Language $\mathcal{BC}$

SDRL: Symbolic Deep Reinforcement Learning

Liu

Introduction

Background

Method

Experiment

Conclusion and Future Work Action Language *BC* (*Lee et al., 2013*) is a language that describes the transition system using a set of *causal laws*. ■ *dynamic laws* describe transition of states

 $move(x, y_1, y_2)$  causes  $on(x, y_2)$  if  $on(x, y_1)$ .

static laws describe value of fluents inside a state

 $intower(x, y_2)$  if  $intower(x, y_1)$ ,  $on(y_1, y_2)$ .



## Background: Reinforcement Learning



- Reinforcement learning is defined on a Markov Decision Process (S, A, P<sup>a</sup><sub>ss'</sub>, r, γ). To achieve optimal behavior, a policy π : S × A → [0, 1] is learned.
  - An option is defined on the tuple (*I*, π, β), which enables the decision-making to have a hierarchical structure:
    - the initiation set  $I \subseteq S$ ,
    - policy  $\pi: \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ ,
    - **probabilistic termination condition**  $\beta: S \mapsto [0, 1].$

# SDRL: Symbolic Deep Reinforcement Learning

SDRL: Symbolic Deep Reinforcement Learning Liu

Introductior

Background

Method

Experiment



- **Symbolic Planner**: orchestrates sequence of *subtasks* using high-level symbolic plan.
- **Controller**: uses DRL approaches to learn the subpolicy for each subtask with *intrinsic rewards*.
- Meta-Controller: measures learning performance of subtasks, updates intrinsic goal to enable reward-driven plan improvement.

## Symbolic Planner



# Symbolic Planner: Planning with Intrinsic Goal

SDRL: Symbolic Deep Reinforcement Learning

Liu

Introduction

Backgroun

Method

Experiment

Conclusion and Future Work Intrinsic goal: a linear constraint on plan quality quality ≥ quality(Πt) where Πt is the plan at episode t.
 Plan quality: a utility function

$$\mathit{quality}(\mathsf{\Pi}_t) = \sum_{\langle s_{i-1}, g_{i-1}, s_i 
angle \in \mathsf{\Pi}_t} 
ho^{g_{i-1}}(s_{i-1})$$

where  $\rho^{g_i}$  is the gain reward for subtask  $g_i$ .

- Symbolic planner: generates a new plan that
  - explores new subtasks,
  - exploits more rewarding subtasks.

#### From Symbolic Transition to Subtask

SDRL: Symbolic Deep Reinforcement Learning

Liu

Introduction

Background

Method

Experiment

- Assumption: given the set S of symbolic states and S of sensory input, we assumed there is an Oracle for symbol grounding: 𝔅 : 𝔅 × 𝔅 → {t, f}.
- Given  $\mathbb{F}$  and a pair of symbolic states  $s, s' \in \mathcal{S}$ :
  - initiation set  $I = \{ \tilde{s} \in \tilde{S} : \mathbb{F}(s, \tilde{s}) = \mathbf{t} \},\$
  - $\pi:\widetilde{\mathcal{S}}\mapsto\widetilde{\mathcal{A}}$  is the subpolicy for the corresponding subtask,
  - $\beta$  is the termination condition such that

$$eta(\widetilde{s'}) = \left\{ egin{array}{cc} 1 & \mathbb{F}(s',\widetilde{s'}) = \mathbf{t}, ext{for } \widetilde{s'} \in \widetilde{\mathcal{S}}, \\ 0 & ext{otherwise.} \end{array} 
ight.$$





#### Controllers: DRL with Intrinsic Reward

SDRL: Symbolic Deep Reinforcement Learning

Liu

Introductior

Method

Experiment

Conclusion and Future Work Intrinsic reward: pseudo-reward crafted by the human.
 Given a subtask defined on (*I*, π, β), intrinsic reward

$$r_i(\tilde{s'}) = \begin{cases} \phi & \beta(\tilde{s'}) = 1 \\ r & \text{otherwise} \end{cases}$$

where  $\phi$  is a positive constant encouraging achieving subtasks and r is the reward from the environment at state  $\tilde{s'}$ .

#### Meta-Controller



#### Meta-Controller: Evaluation with Extrinsic Reward

SDRL: Symbolic Deep Reinforcement Learning

Liu

Introduction

Method

Experiment

Conclusion and Future Work Extrinsic rewards: r<sub>e</sub>(s, g) = f(ε) where ε can measure the competence of the learned subpolicy for each subtask.
 For example, let ε be the success ratio, f can be defined as

$$F(\epsilon) = \left\{egin{array}{cc} -\psi & \epsilon < ext{threshold} \ r(m{s},m{g}) & \epsilon \geq ext{threshold} \end{array}
ight.$$

f

- $\psi$  is a positive constant to punish selecting unlearnable subtasks,
- r(s,g) is the cumulative environmental reward by following the subtask g.

#### Experimental Results I.

SDRL: Symbolic Deep Reinforcement Learning

Liu

Introduction Background Method

#### Experiment

% object declaration			
location(mp;rd;ls;lll;lrl;key).			
% dynamic causal law declaration			
move(L) causes loc=L if location(L).			
move(L) causes cost=L+Z if rho((at(L1)),move(L))=Z			
loc=L1, picked(key)=false.			
move(L) causes cost=L+Z if rho((at(L1),picked(key)			
move(L))=Z,loc=L1,picked(key)=true.			
inertial loc. inertial quality.			
% static causal law declaration			
picked(key)=true if loc=key.			
nonexecutable move(key) if picked(key).			
default rho((at(L1)),move(L))=10.			
<pre>default rho((at(L1),picked(key)),move(L))=10.</pre>			

No.	subtask	policy learned	in optimal plan
1	MP to LRL, no key	· ·	· · · ·
2	LRL to LLL, no key	~	1
3	LLL to key, no key	~	✓
4	key to LLL, with key	1	✓
5	LLL to LRL, with or without key	~	✓
6	LRL to MP, with or without key	~	✓
7	MP to RD, with key	<ul> <li>✓</li> </ul>	✓
8	LRL to LS, with or without key	√	
9	LS to key, with or without key	✓	
10	MP to RD, no key	1	
11	LRL to key, with or without		
12	key to LRL, with key		
13	LRL to RD, with key		







#### Experimental Results II.









Conclusion and Future Work



Baseline: Kulkarni et. al, Hierarchical Deep Reinforcement Learning: Integrating Temporal Abstraction and Intrinsic Motivation, NIPS'2016. 19 / 20

#### Conclusion

SDRL: Symbolic Deep Reinforcement Learning

Liu

Introductior

Background

Method

Experiment

- We present a **SDRL** framework features:
  - High-level symbolic planning based on intrinsic goal
  - **Low-level policy control** with DRL.
  - **Subtask learning evaluation** by a meta-learner.
- This is the first work on integrating symbolic planning with DRL that achieves both task-level interpretability and data-efficiency for decision-making.
- Future work.