



Optimizing Medical Treatment for Sepsis in Intensive Care: From Reinforcement Learning to Pretrial Evaluation

by Li, Albert-Smet, Faisal

Weiqing Wang^{1,2}
Cong Wei¹

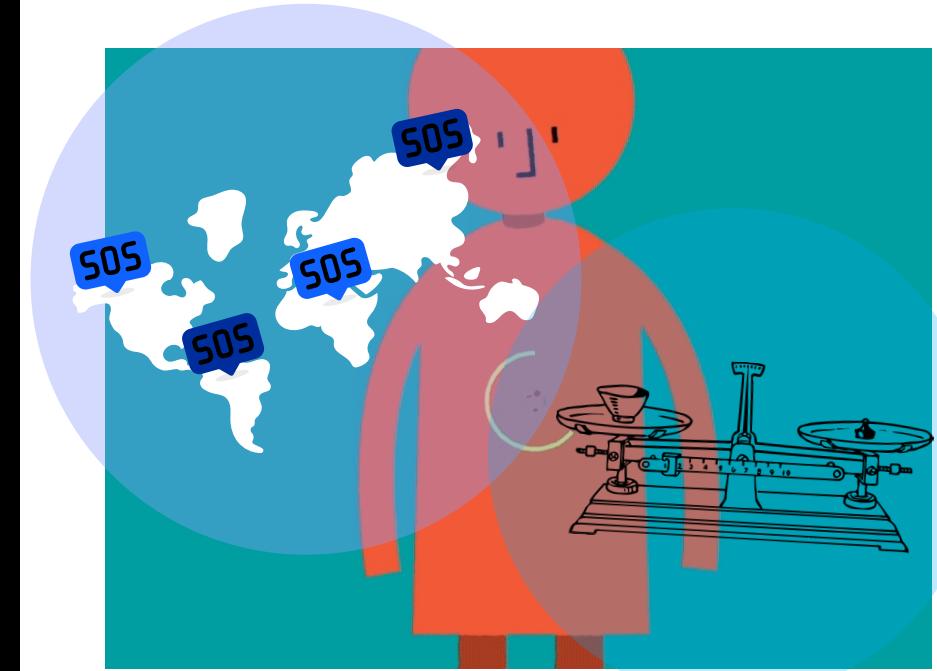


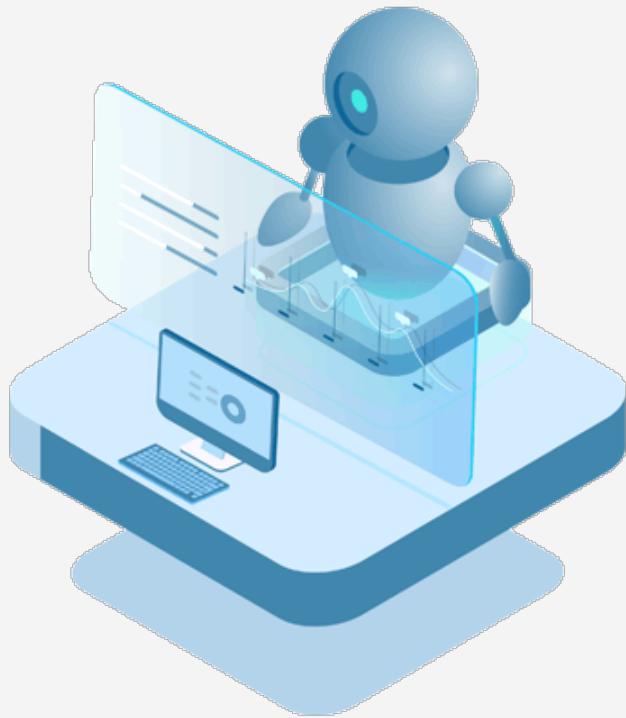
1. Department of Computer Science, University of Toronto.

2. Worldwide Ecosystem, International Business Machines (IBM)

The aim is to establish a policy for prospective clinical testing
in general, then

Why Sepsis?

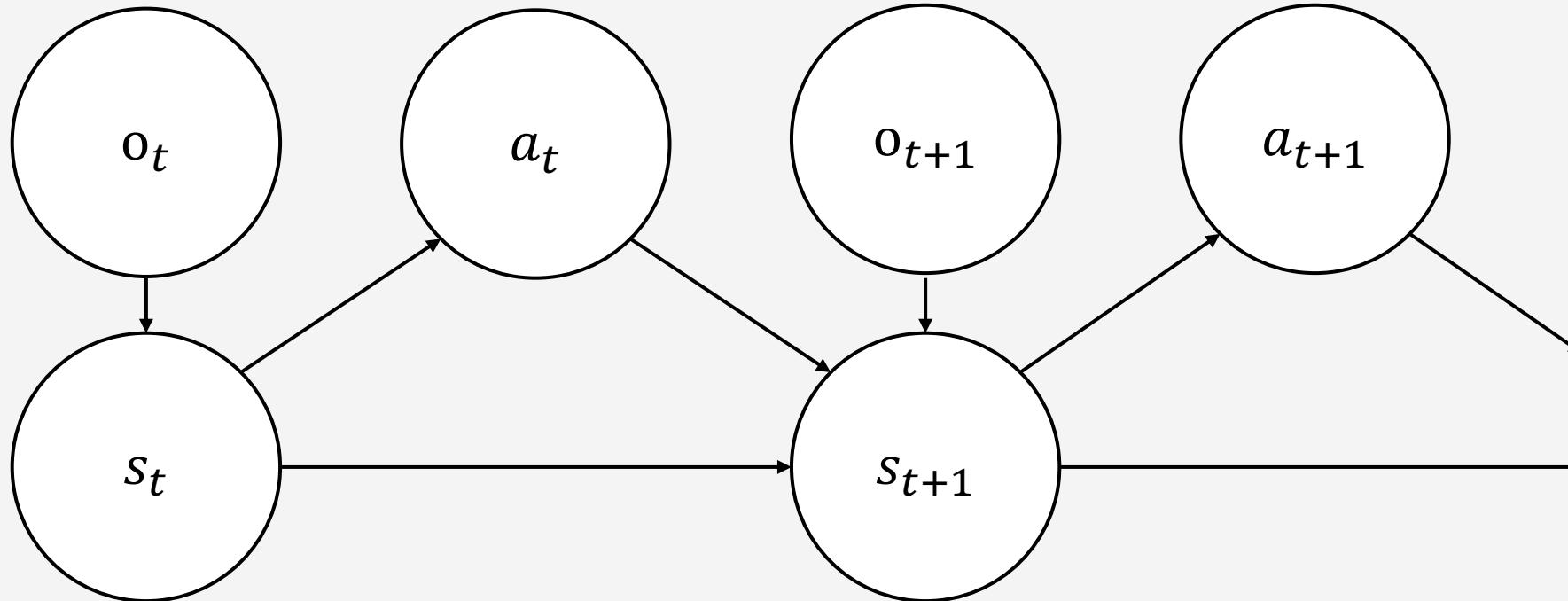




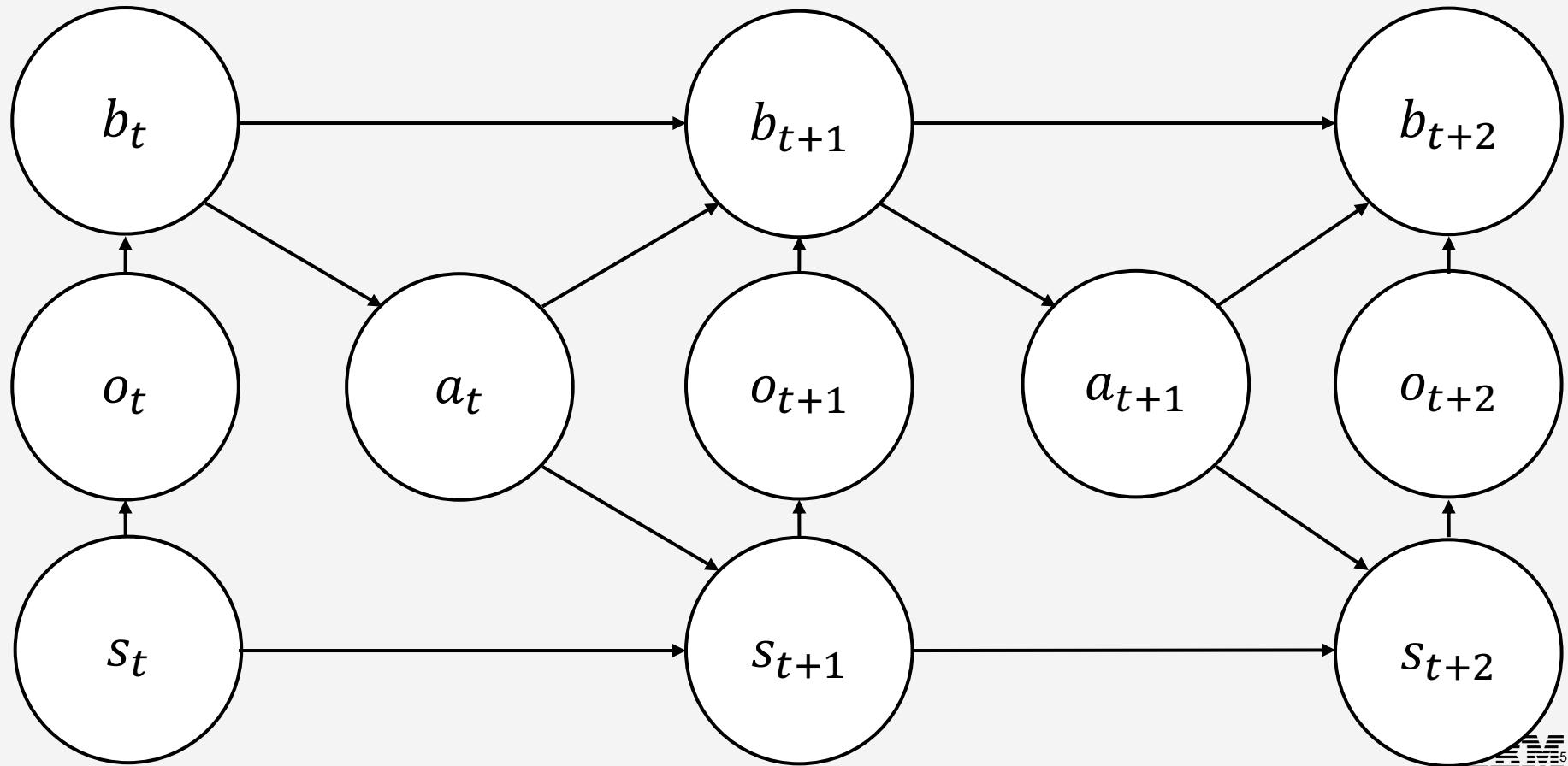
Why Reinforcement Learning?

~~Hidden~~ Decision Process

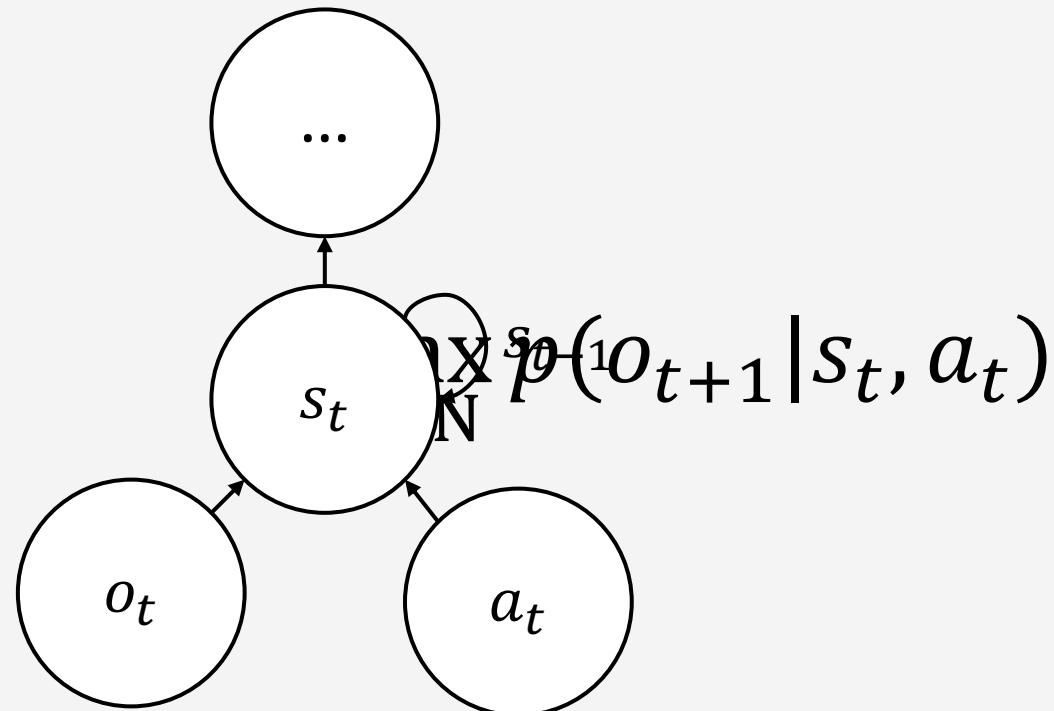
Markov



Partially Observable Markov Decision Process



$$\begin{aligned} o_{t+1} &= g(o_0, a_0, \dots, o_{t-1}, a_{t-1}, o_t, a_t) \\ &= h(s_{t-1}, a_t, o_t) \end{aligned}$$



Policy Optimization



Value-based

- Train $Q(S, a)$.
- Policy:
$$\operatorname{argmax}_{a \in \mathcal{A}} Q(S, a)$$

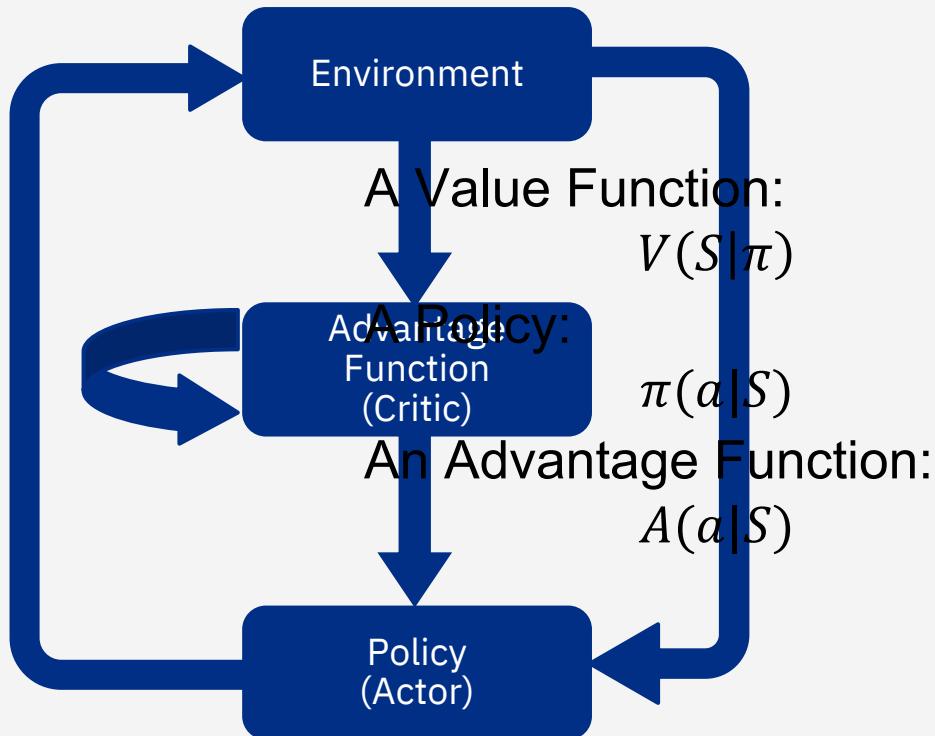
Policy Based

- π is stochastic.
- Choose π that maximises $V^\pi(S_0)$

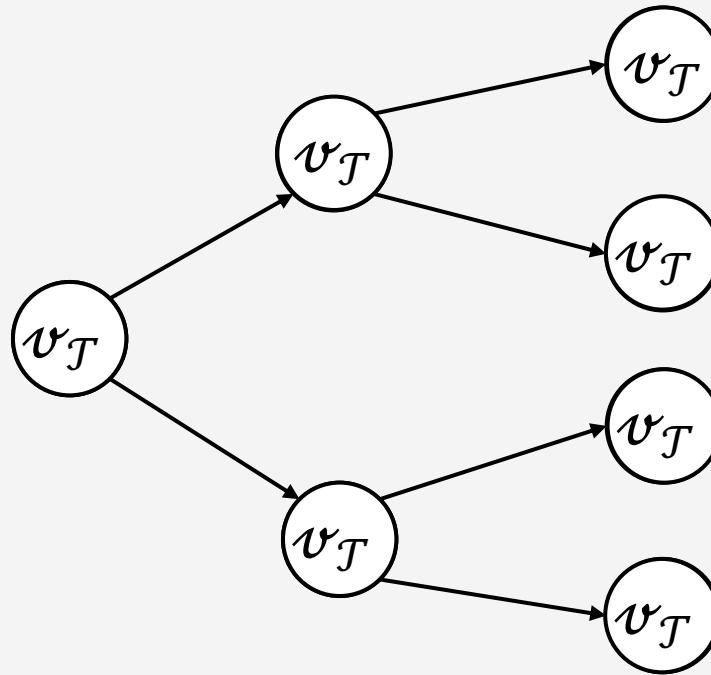
Model Based

- Agent Evaluation
- Agent Training
- World Model

Advantage Actor Critic (A2C)

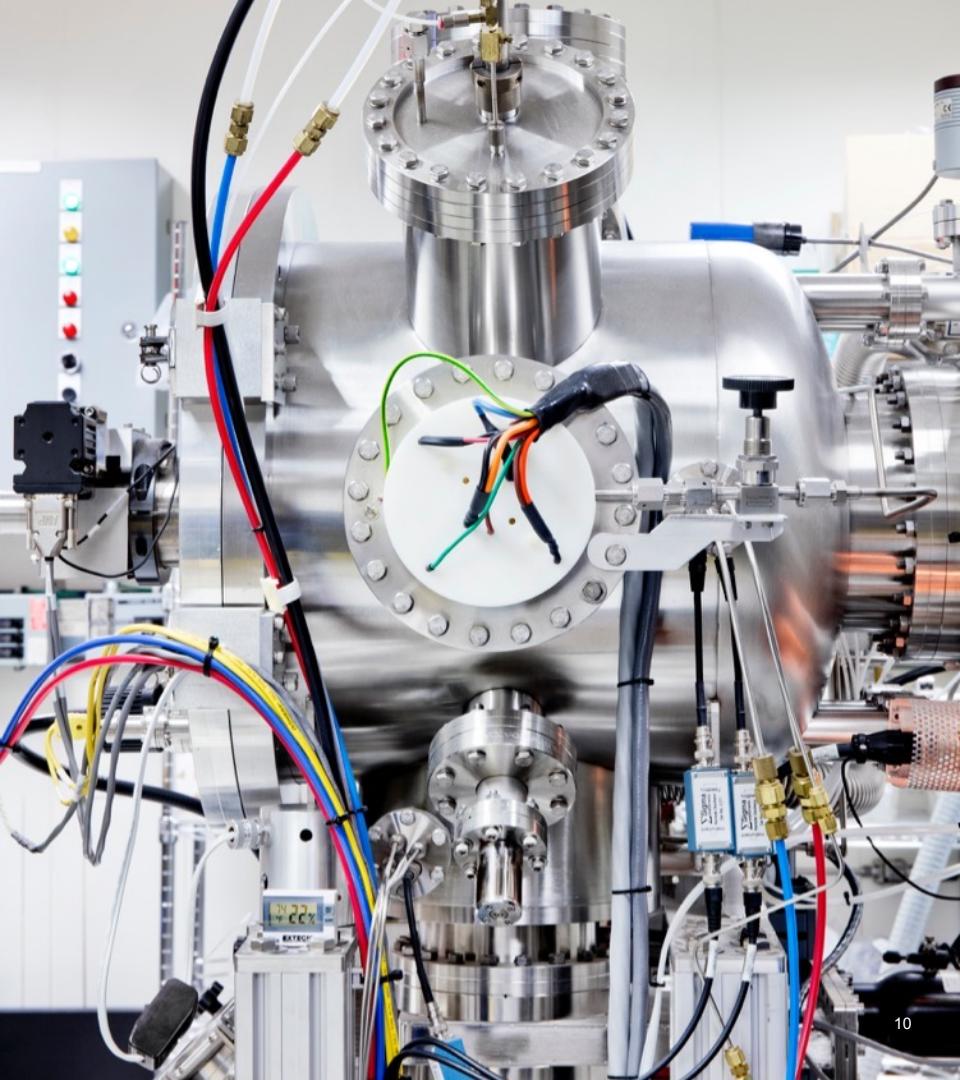
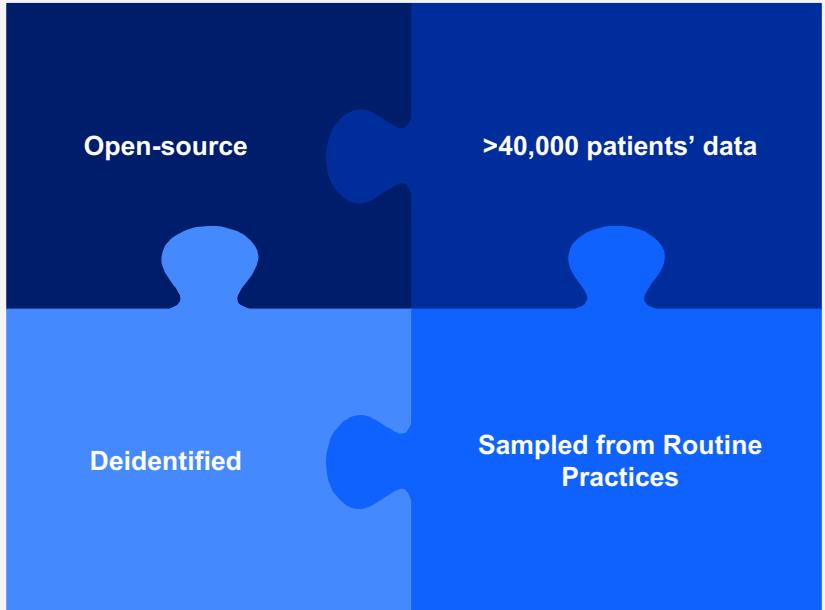


Heuristic Tree Search

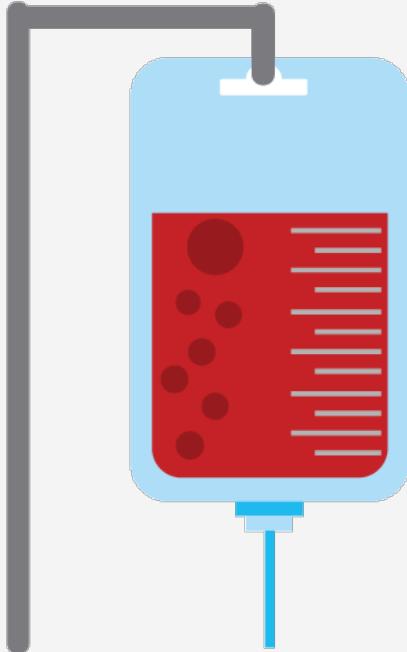


Preclinical Training Dataset

Medical Information Mart for Intensive Care (MIMIC)



The Action Space



Maximum dose of
vasopressors administered.

The total volume of
intravenous fluids injected
over each hourly period.

Reward Function

Survival



-10

Death



+10

Discharge

Pre-clinical Testing



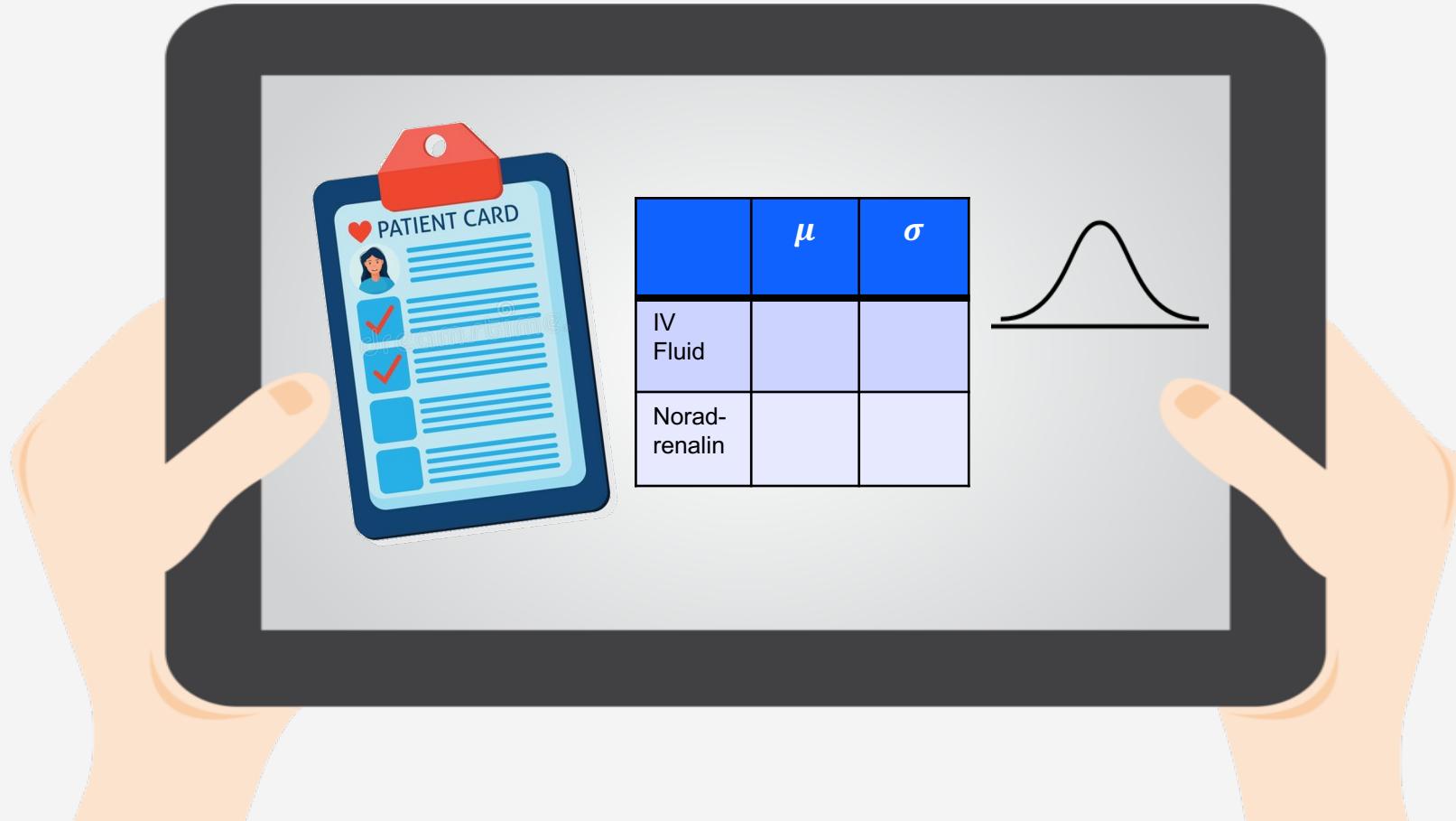
Consistency?

Whether AI is consistent with
human clinicians?

Model-agnostic

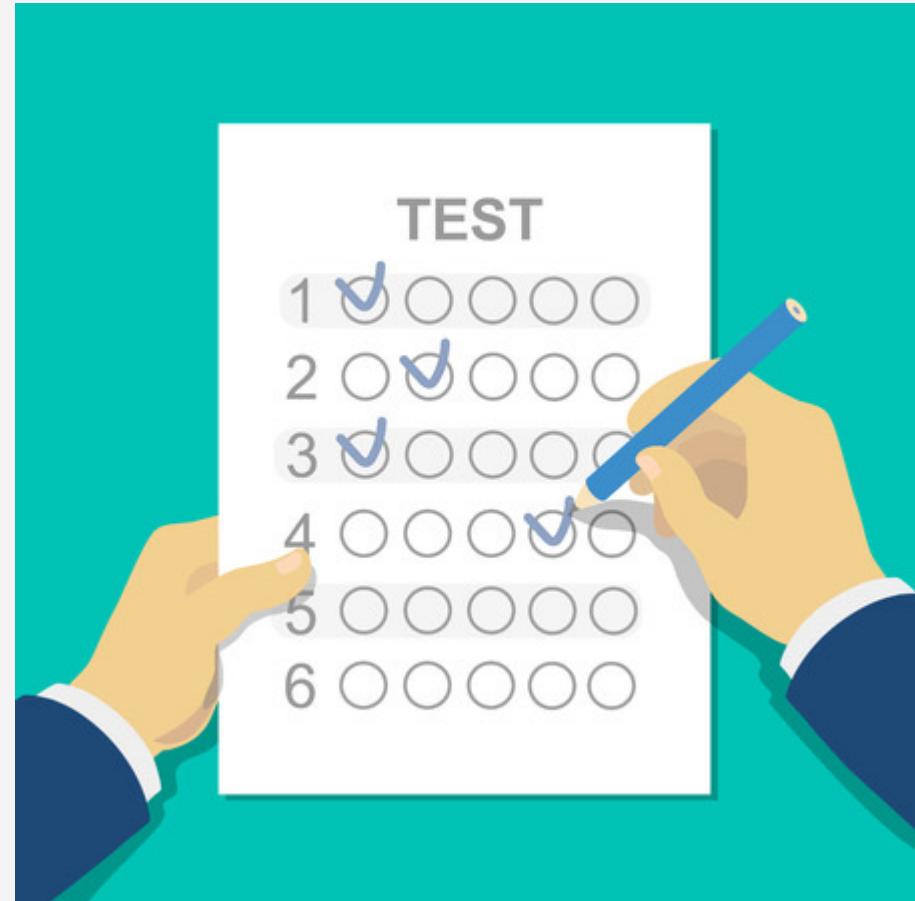
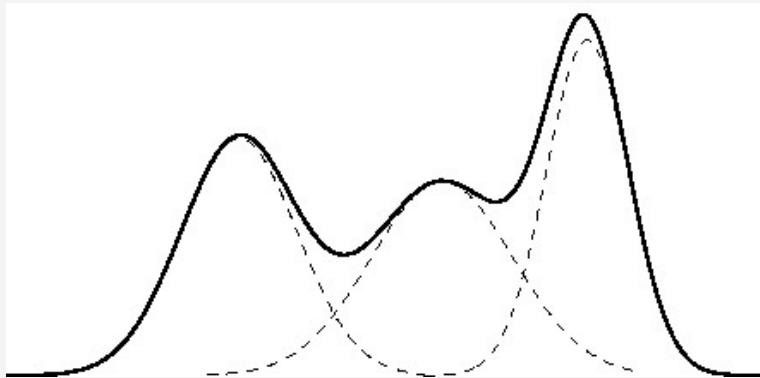
A novel framework is proposed.

Pre-clinical Testing

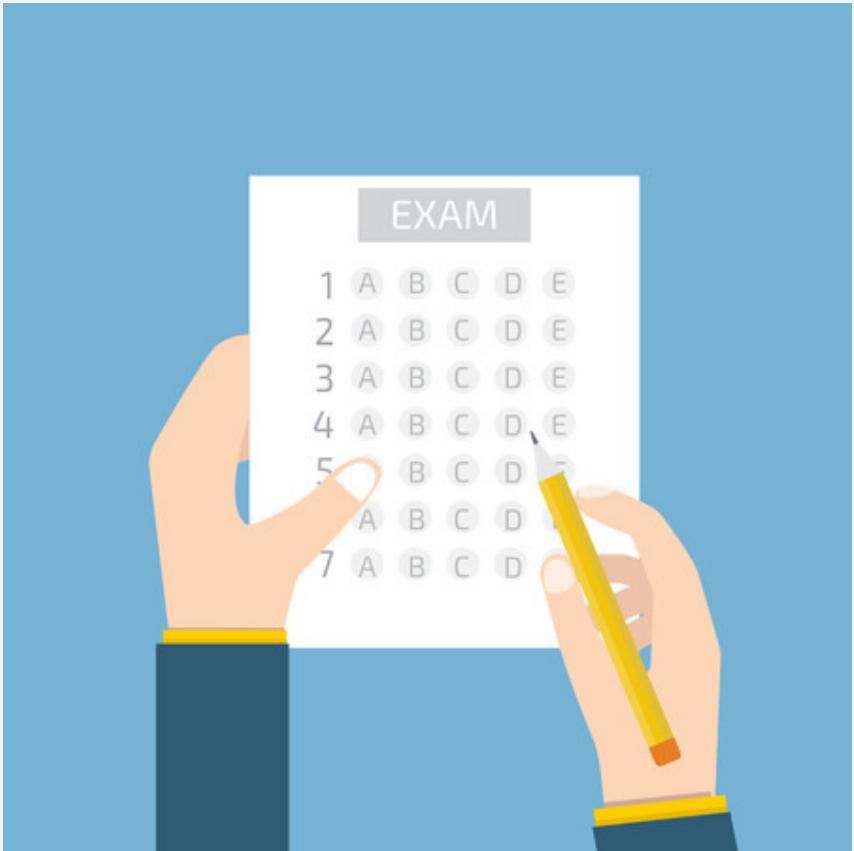


Evaluation – P-score

$$P_{\text{score}}(a) = \frac{1}{N} \sum_{i=1}^N f_{\mu_i, \sigma_i}(a)$$



Evaluation – C-score

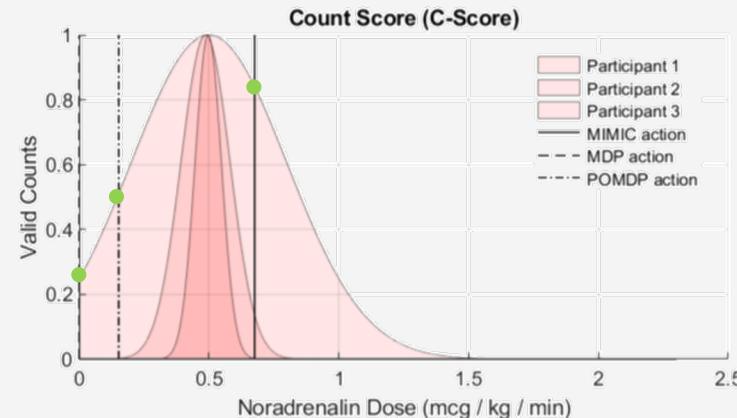
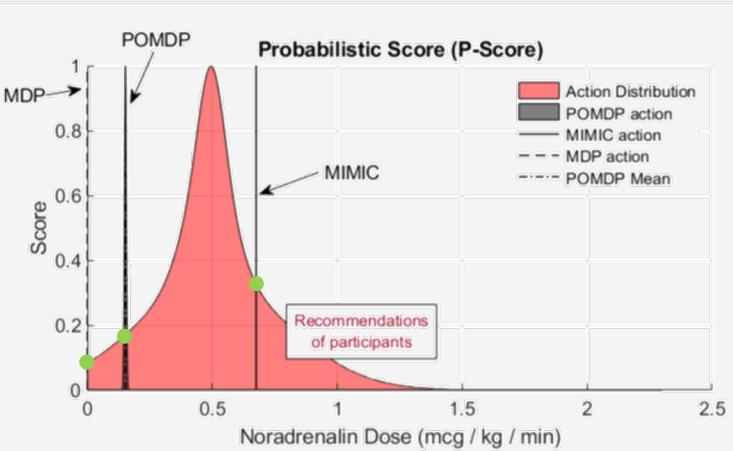
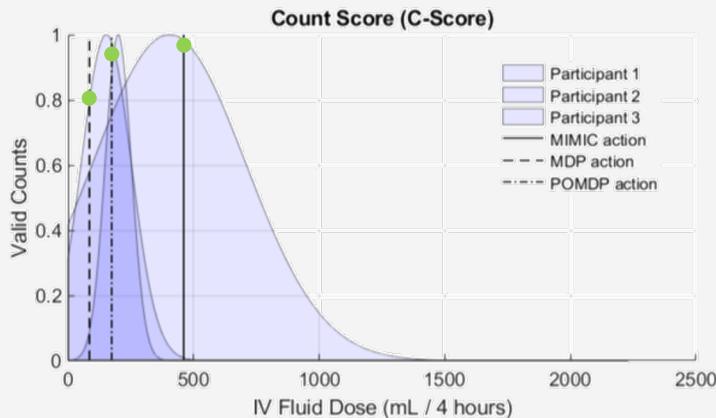
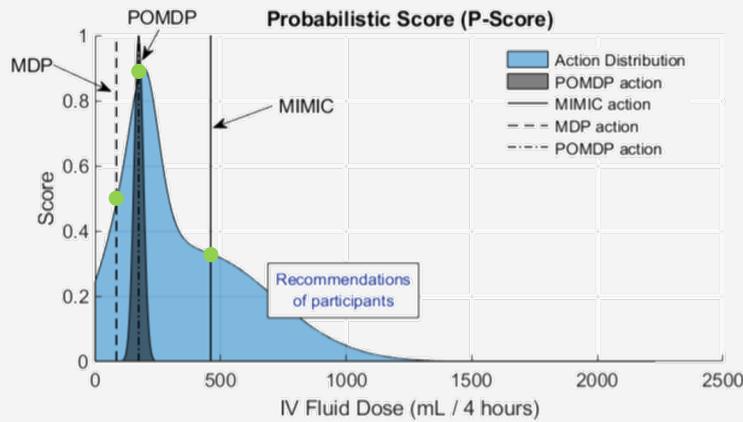


$$C_{\text{score}}(a) = \hat{p}(\mu_* - 3\sigma_* \leq a \leq \mu_* + 3\sigma_*)$$

$$= \frac{1}{N} \sum_{i=1}^N I(\mu_i - 3\sigma_i \leq a \leq \mu_i + 3\sigma_i)$$

$$\approx p(\mu_* - 3\sigma_* \leq a \leq \mu_* + 3\sigma_*)$$

Result



Result

Table 1: This table contains the average validation scores observed when three separate intensive care consultants evaluated 10 patient trajectories from the MIMIC dataset. The best score for each row is shown in bold. The column labels stand for: the original action recorded (MIMIC), a discrete and non-time dependant model (MDP) by Komorowski et al. (2018), and the model we propose in this work (POMDP).

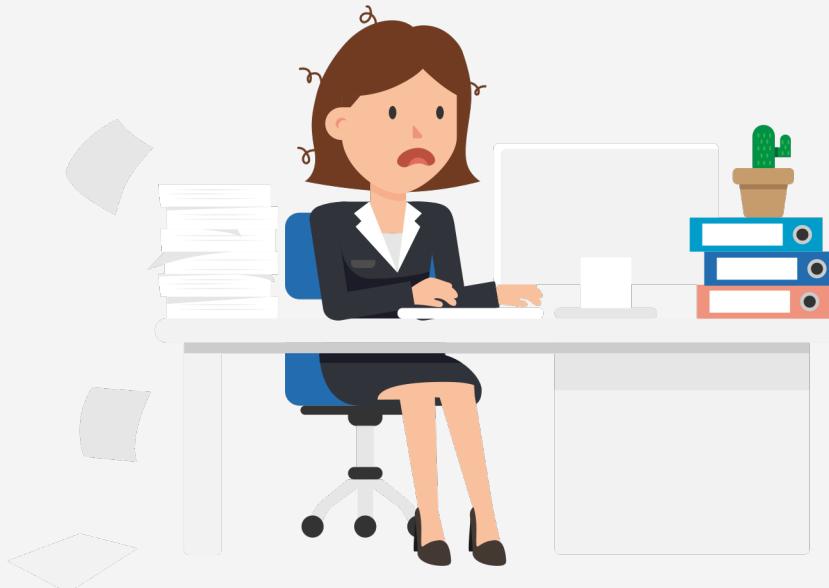
| ACTION | SCORE | MIMIC | MDP | POMDP |
|--------------|------------|--------------|-------|--------------|
| IV Fluids | P-Score | 0.454 | 0.434 | 0.448 |
| | C-Score | 0.578 | 0.622 | 0.644 |
| | Zero Count | 0.133 | 0.100 | 0.033 |
| Vasopressors | P-Score | 0.584 | 0.286 | 0.488 |
| | C-Score | 0.700 | 0.467 | 0.644 |
| | Zero Count | 0.033 | 0.233 | 0.033 |

Strength

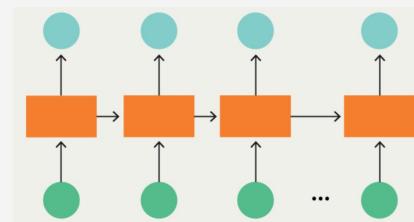
- ✓ POMDP fits Medical Setting
- ✓ An accurate state representation
- ✓ Preclinical Validation Framework



Limitations



- ✓ Large Dataset Required
- ✓ Small Number of Participants
- ✓ Pitfall of RNN



Thank you

Weiqing Wang

weiqing.wang@utoronto.ca | Weiqing.Wang@ibm.com

Cong Wei

cong.wei@mail.utoronto.ca



Appendix: Actor Critic

$$A(S, a) = \begin{cases} \delta_t & \text{last nonterminal step} \\ \gamma c_t \max(A(s_{t+1}, a_{t+1}), 0) + \delta_t & \text{otherwise} \end{cases}$$

- δ_t is the off-policy temporal difference error for V at time t .
- c_t is the contracted Importance Sampling Ratio.
- γ is the discount factor.

Appendix: Actor Critic

$$s^* = \operatorname{argmax}_{s \in \mathcal{F}} \gamma^{D(s)-1} \prod_{d=0}^{D(s)-1} \hat{p}(s^{d+1} | s^d, a_{\mathcal{T}}(s^d)) \delta_{\mathcal{T}(s^0, s)}(s^d)$$

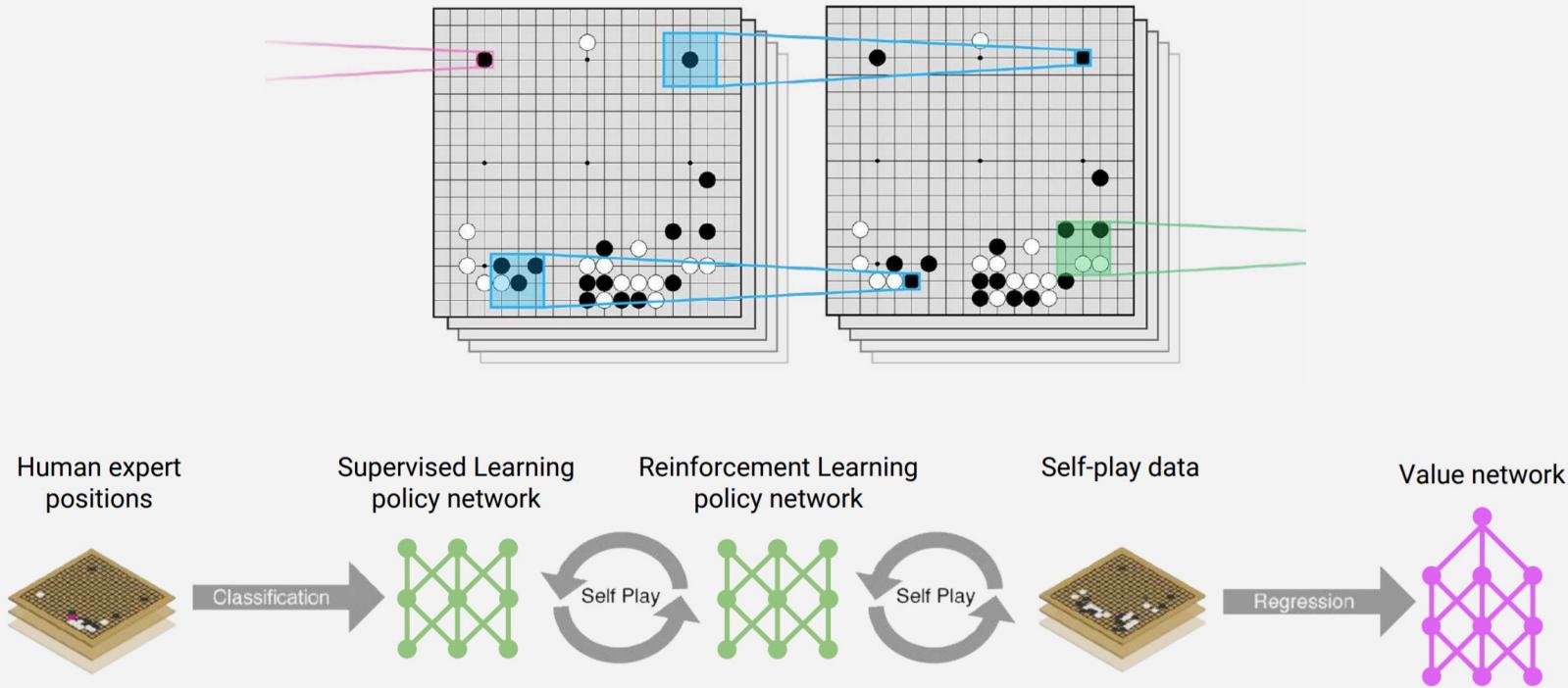
- $D(s)$ is the depth of the node.
- δ_x is the Dirac delta function concentrated at x .
- $a_{\mathcal{T}}(s^d)$ deterministic tree policy in s^d .

Appendix: Actor Critic

$$a_{\mathcal{T}}(s^d) = \operatorname{argmax}_{a \in \mathcal{A}} r(s^d, a) + \gamma \sum_{s^{d+1} \in C(s^d, a)} \hat{p}(s^{d+1} | s^d, a) V(s^{d+1})$$

$$v_{\mathcal{T}}(s^d) = \begin{cases} V(s^d) & s^d \text{ is a leaf} \\ r(s^d, a_{\mathcal{T}}(s^d)) + \gamma \sum_{s^{d+1} \in C(s^d, a_{\mathcal{T}}(s^d))} \hat{p}(s^{d+1} | s^d, a_{\mathcal{T}}(s^d)) v_{\mathcal{T}}(s^{d+1}) & \text{otherwise} \end{cases}$$

Appendix: Limitation



Appendix: Pitfall of RNN

