



南京航空航天大學

Nanjing University of Aeronautics and Astronautics



模式分析与机器智能
工业和信息化部重点实验室

MIIT Key Laboratory of
Pattern Analysis & Machine Intelligence

Fast Counterfactual Inference for History-Based Reinforcement Learning

AAAI | 2023



Background

- **History-based Decision Process (HDP)**
- $\langle O, H, A, T, P, R, \gamma \rangle$
 - O: the observation space
 - H: the history space
 - A: the action space
 - T: the horizon(time steps into the future when making decision)
 - P: the history-action transition probability**
 - R: reward function
 - γ : discount rate



Background

- **History-based RL**

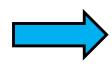
- current state + sequence of past states and actions = Decision Making

- **Pros**

- Partial Observability
- Handling Sequential Dependencies
- Adaptive Behavior

- **Cons**

- Data efficiency
- Computational Complexity
- **do not adjust for confounding correlations**



Counterfactual
Inference (CI)



Background

- **Experiment**

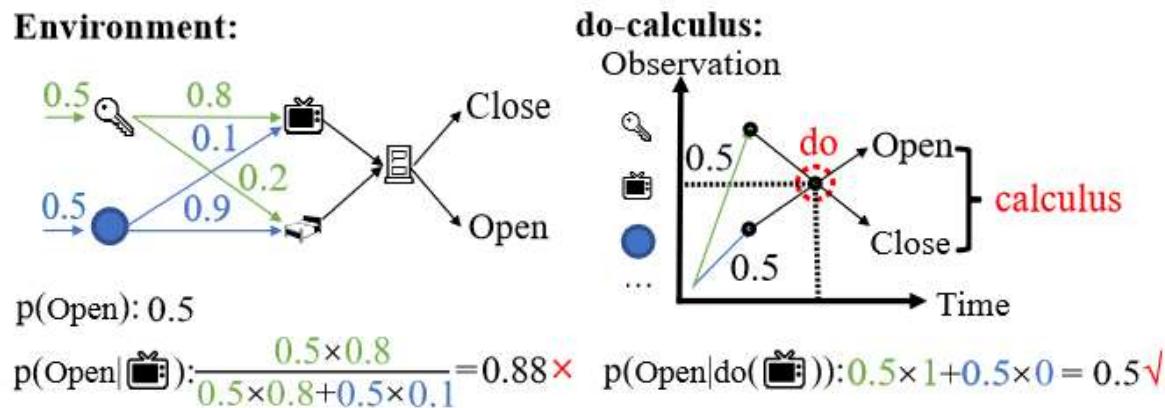


Figure 1: Key-to-door example. The high correlation on TV caused by sampling can be eliminated by do-calculus which separates confounders (key and ball).

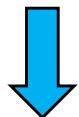


Background

- **Former Effort**

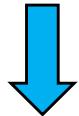
- Prior researchers perform CI on Markovian RL tasks for feature selection **where the states are Markovian**

单变量干预



- for history-based RL, we should intervene for **each non-Markovian observation-and-time combination** ot and estimate its causal effect

不适用



high computational complexity due to the large scale of historical observations

需要开发多变量干预





Preliminaries

- **Do-calculus and Backdoor Formula**

- **causal diagram G:**

- covariates $X := \{X_1, \dots, X_n\}$
 - a response variable Y

- **the intervention operation $\text{do}(X_j = x)$:**

adopted by imposing a certain value x on one of the covariates X_j

$$p(Y|\text{do}(X_j = x)) = \int p(Y|X_j = x, \mathbf{X}^{\text{ba}}) dp(\mathbf{X}^{\text{ba}}),$$

$X_{\text{ba}} \subseteq X$ are the backdoor variables relative to (X_j, Y)

- 1) X_{ba} contains no descendant of X_j
- 2) X_{ba} blocks each path pointed to X between X_j and Y



Preliminaries

- **Fine-Grained History Counterfactual Inference**

- **Covariates**

- the past observations, the current observation, and the action $\{h_{t-1}, o_t, a_t\}$

- **response variables**

- and the immediate reward and next observation $\{r_{t+1}, o_{t+1}\}$

Proposition 1 (Fine-grained CI). *Given o_t and a_t , h_{j-1} satisfies the backdoor formula relative to $(o_j, \{o_{t+1}, r_{t+1}\})$ for any historical observation $o_j (j < t)$, and the causal effect of $\text{do}(o_j = o)$ can be estimated with*

$$p(o_{t+1}, r_{t+1} | \text{do}(o_j = o), o_t, a_t) \quad (2)$$

$$= \int_{h_{j-1} \in \mathcal{H}_{j-1}} p(o_{t+1}, r_{t+1} | o_j = o, h_{j-1}, o_t, a_t) dp(h_{j-1}).$$





Preliminaries



$$p(o_{t+1}, r_{t+1} | \text{do}(o_j = \mathbf{0}), o_t, a_t) \quad (3)$$

$$= \int_{h_{j-1} \in \mathcal{H}_{j-1}} p(o_{t+1}, r_{t+1} | o_j = \mathbf{0}, h_{j-1}, o_t, a_t) dp(h_{j-1}).$$

- **Average Treatment Effect (ATE)**

The difference between the counterfactual and factual distributions

$$\text{ATE}(o_j^{\overline{o, o'}}) := \sum_{t=j+1}^T \mathbb{E}_{\substack{o_j \in \{o, o'\}, \\ o_t \in \mathcal{O}, a_t \in \mathcal{A}}} \quad (4)$$

$$\|p(\cdot | o_j, o_t, a_t) - p(\cdot | \text{do}(o_j = \mathbf{0}), o_t, a_t)\|_1.$$

ATE趋于0表示无因果

- History-base 需要遍历所有的历史观测，计算量大

Can we develop a method that can decide the causality among numbers of variables?



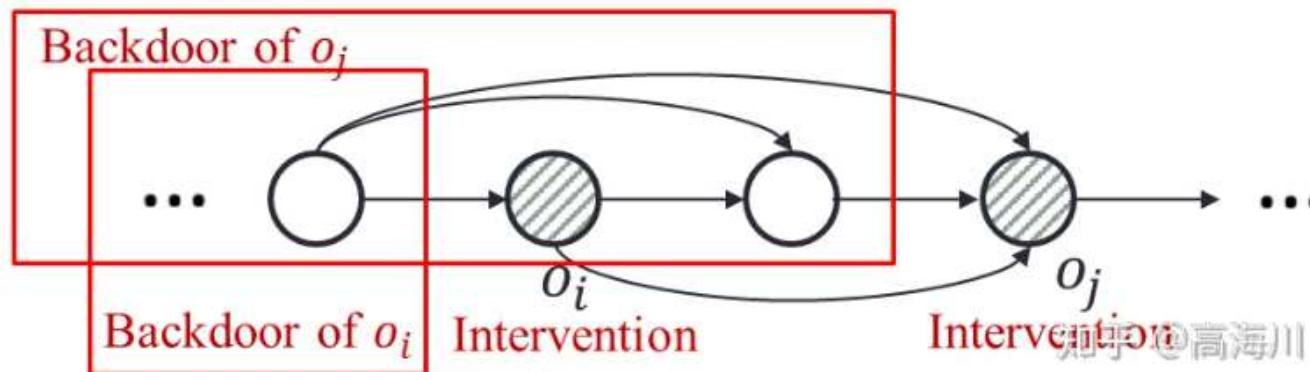


Coarse-to-Fine History Counterfactual Inference

- Steps

- estimate the causal effects of observations without timestamps. $\Omega(T \cdot |O|) \rightarrow \Omega(|O|)$. 空间推理
- Second, perform CI on observation sub-spaces. $\Omega(|O|) \rightarrow \Omega(\log|O|)$ 时间推理

- Problem



multiple intervened historical observations have no common backdoor variables



Coarse-to-Fine History Counterfactual Inference

- Counterfactual Inference on Observations

Definition 1 (Step-backdoor adjustment formula). $X_j^{s-\text{ba}} = X_j^{\text{ba}} \setminus (X_i^{\text{ba}} \cup \{X_i\})$ is step-backdoor relative to (X_j, Y) if 1) $X_j^{s-\text{ba}}$ has no descendant of X_j , 2) X_i^{ba} , X_i , and $X_j^{s-\text{ba}}$ block each path between X_j and Y , and 3) conditioned on X_i and X_i^{ba} , the distribution of $X_j^{s-\text{ba}}$ is identifiable.

Theorem 1. Given a set of intervened variables with different timestamps, if every two temporally adjacent variables meet the step-backdoor adjustment formula, then the overall causal effect can be estimated with



Coarse-to-Fine History Counterfactual Inference

$$\begin{aligned}
 & p(Y|do(X_i = x, X_j = x', X_k = x'', \dots)) \\
 &= \int \dots \int p(Y|\mathbf{X}_i^{ba}, \mathbf{X}_j^{s-ba}, \mathbf{X}_k^{s-ba}, \dots, \\
 & \quad X_i = x, X_j = x', X_k = x'', \dots) \tag{6}
 \end{aligned}$$

计算 X_i 后门 (后门准则)
计算 X_j 和 X_i 的 step-backdoor (SBAF)

$$\begin{aligned}
 & dp(\mathbf{X}_i^{ba}) \\
 & dp(\mathbf{X}_j^{s-ba}|X_i = x, \mathbf{X}_i^{ba}) \\
 & dp(\mathbf{X}_k^{s-ba}|X_i = x, \mathbf{X}_i^{ba}, X_j = x', \mathbf{X}_j^{s-ba}) \dots
 \end{aligned}$$

Theorem 2 (CI on observations). *Given o_t and a_t , the causal effect of $Do(o)$ can be estimated by*

$$p(o_{t+1}, r_{t+1} | Do(o), o_t, a_t) \tag{7}$$

$$= \int_{h_{t-1} \in \mathcal{H}_{t-1}^o} p(o_{t+1}, r_{t+1} | h_{t-1}, o_t, a_t) dp(h_{t-1} | \mathcal{H}_{t-1}^o),$$

$\Omega(T \cdot |\mathcal{O}|) \rightarrow \Omega(|\mathcal{O}|)$. 空间推理

where \mathcal{H}_{t-1}^o denotes the history sub-space where each history $h_{t-1} \in \mathcal{H}_{t-1}^o$ contains at least one observation with value o .



Coarse-to-Fine History Counterfactual Inference

- Counterfactual Inference on Sub-Spaces

Proposition 2 (CI on observation sub-spaces). *Given o_t and a_t , the causal effect of $\text{Do}(\mathcal{O}^i)$ can be estimated by*

$$\begin{aligned} & p(o_{t+1}, r_{t+1} | \text{Do}(\mathcal{O}^i), o_t, a_t) \\ &= \int_{h_{t-1} \in \mathcal{H}_{t-1}^{\mathcal{O}^i}} p(o_{t+1}, r_{t+1} | h_{t-1}, o_t, a_t) dp(h_{t-1} | \mathcal{H}_{t-1}^{\mathcal{O}^i}), \end{aligned} \quad (8)$$

where $\mathcal{H}_{t-1}^{\mathcal{O}^i}$ represents the history subspace where each history $h_{t-1} \in \mathcal{H}_{t-1}^{\mathcal{O}^i}$ contains at least one observation belonging to the observation subspace \mathcal{O}^i .

\mathcal{O}^{Ca} : observation sub-space with historical causality.

Suppose:

an observation space is divided into $Z \geq |\mathcal{O}^{Ca}|$ sub-spaces
at least $Z - |\mathcal{O}^{Ca}|$ parts containing no causal observation

Proposition 3 (Coarse-to-fine CI). *If $Z \geq |\mathcal{O}^{Ca}|$, the number of interventions for coarse-to-fine CI is $\Omega(\log|\mathcal{O}|)$.*

$$\Omega(|\mathcal{O}|) \rightarrow \Omega(\log|\mathcal{O}|) \quad \text{时间推理}$$



Coarse-to-Fine History Counterfactual Inference

- Combining RL and Coarse-to-Fine Counterfactual Inference

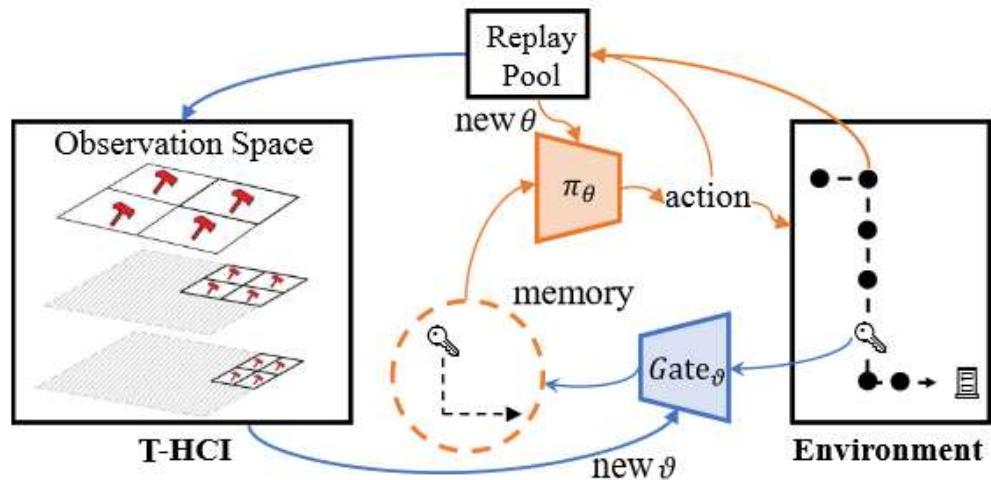


Figure 3: T-HCI Algorithm framework. The blue and orange lines respectively mark the CI loop and RL loop.

算法包含两个loops

- 1) T-HCI loop,
- 2) 策略学习loop

两者交换进行：在策略学习loop里，agent被采样学习一定回合数量，并将样本存在replay pool中；在T-HCI loop中，利用存储的样本进行上述的因果推理过程



Coarse-to-Fine History Counterfactual Inference

- Combining RL and Coarse-to-Fine Counterfactual Inference

$\hat{\mathcal{O}}$: the constructed discrete observation space

\mathcal{O}^i : the current intervened observation subspace

\mathcal{O}^w : the observations that have been eliminated by previous inference

$$\text{ATE}(\overline{\mathcal{O}^i}, \overline{\mathcal{O}^w}) \quad (9)$$

$$= \min_{\zeta} \sum_{t=1}^T \mathbb{E}_{h_t \in \mathcal{H}_t} \left| \left| p_{\zeta}(\cdot | h_{t-1}^{\overline{\mathcal{O}^i}, \overline{\mathcal{O}^w}}, o_t, a_t) - \hat{p}(\cdot | h_t, a_t) \right| \right|_1. \quad \mathbf{I}(\cdot): \text{the indicator function}$$

$$\mathcal{L}_{\text{Gate}}(\vartheta) = \frac{1}{|\hat{\mathcal{O}}|} \sum_{o \in \hat{\mathcal{O}}} |\text{Gate}_{\vartheta}(o) - I(o \in \mathcal{O}^w)|, \quad (10)$$

Gate $\theta(o_j)$ taking value 1 if and 0 otherwise

$$\text{Gate}_{\vartheta}(o_j) = \begin{cases} 1 & o_j \in \mathcal{O}^w \\ 0 & \text{otherwise} \end{cases}$$



Coarse-to-Fine History Counterfactual Inference

- **Policy learning**
 - Built on the memories of causal historical observations,
mapping $\phi : \{o_0, \dots, o_{t-1}, o_t\} \rightarrow \{\text{Gate}(o_0) \ominus o_0, \dots, \text{Gate}(o_{t-1}) \ominus o_{t-1}, o_t\}$
 - name ψ as causal memory. a policy $\pi_\theta(a_t | \phi_\theta(h_t))$
 - Train with Deep Monte Carlo (DMC) , A2C and PPO



Coarse-to-Fine History Counterfactual Inference

Algorithm 1 PPO, Actor-Critic Style

```
for iteration=1,2,... do
    for actor=1,2,...,N do
        Run policy  $\pi_{\theta_{\text{old}}}$  in environment for  $T$  timesteps
        Compute advantage estimates  $\hat{A}_1, \dots, \hat{A}_T$ 
    end for
    Optimize surrogate  $L$  wrt  $\theta$ , with  $K$  epochs and minibatch size  $M \leq NT$ 
     $\theta_{\text{old}} \leftarrow \theta$ 
end for

for iteration 1,2... do,
    for iteration 1,2,... do,
        PPO algorithm (update mapping  $\phi$ )
    end for
    store  $\theta$  into replay pool
    for iteration 1,2..., do
        T-HCL algorithm (update  $\vartheta$ )
    end for
end for
```



Experiments

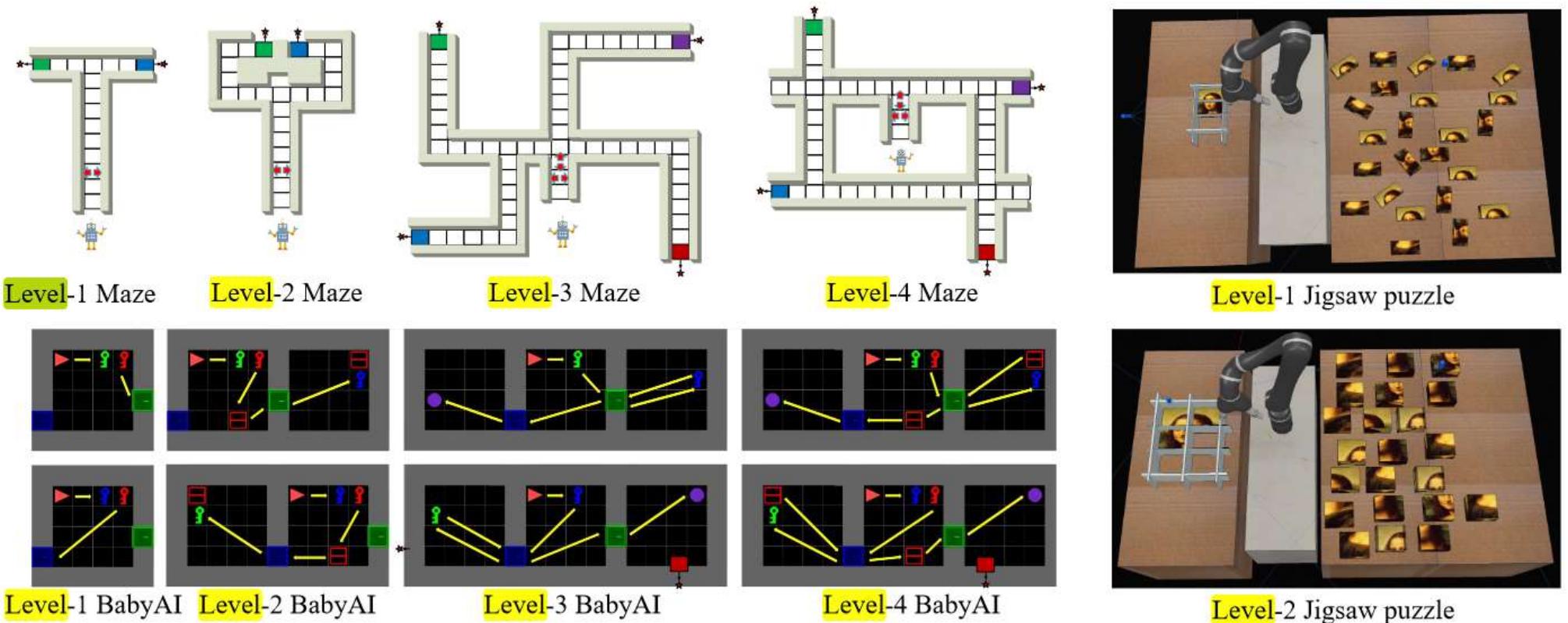


Figure 4: Environments of Maze, BabyAI, and Jigsaw Puzzle tasks.



Experiments

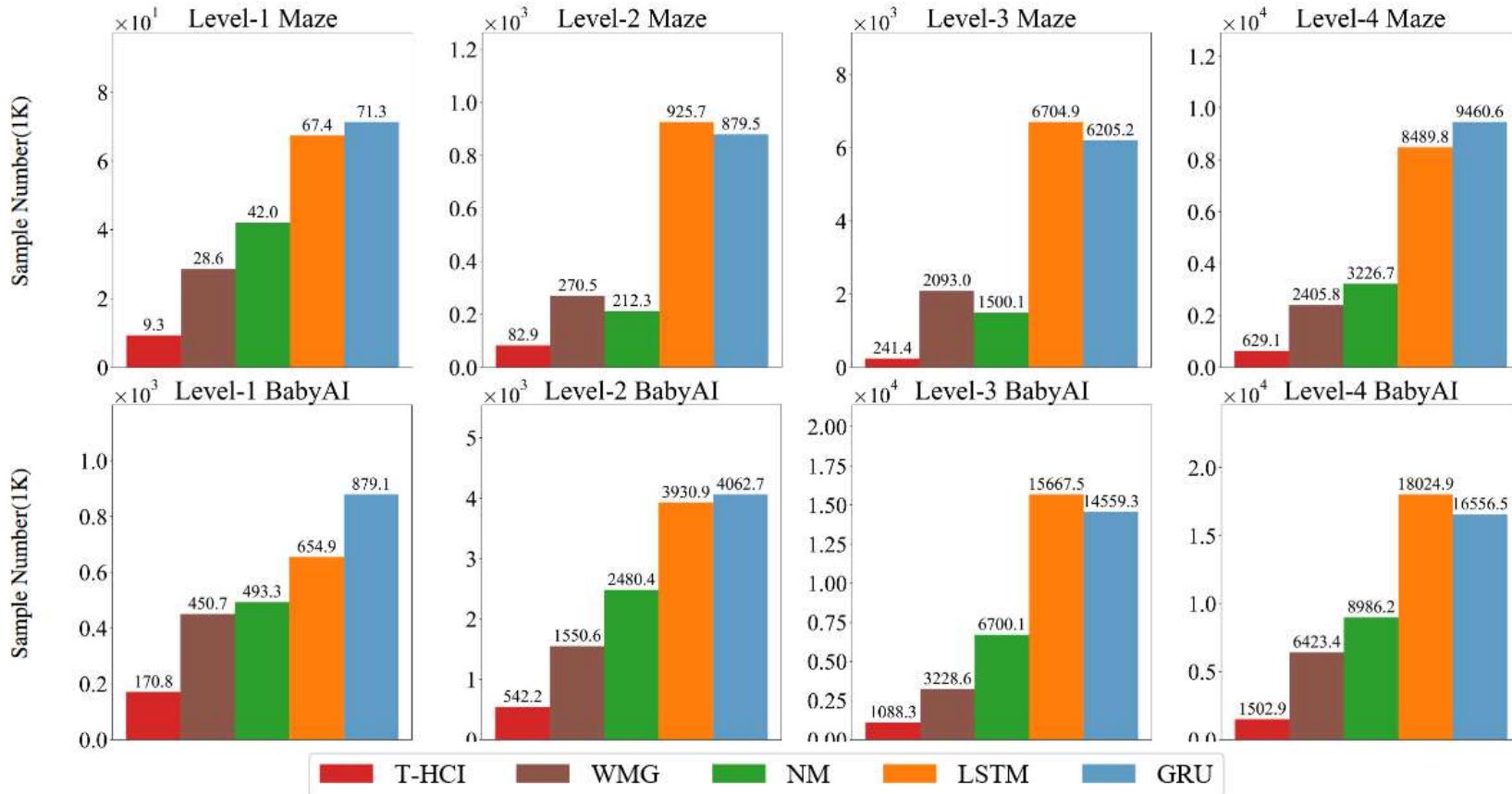


Figure 5: Comparison of average numbers of samples (thousand) in 10 trials of the Maze and BabyAI tasks.



Experiments

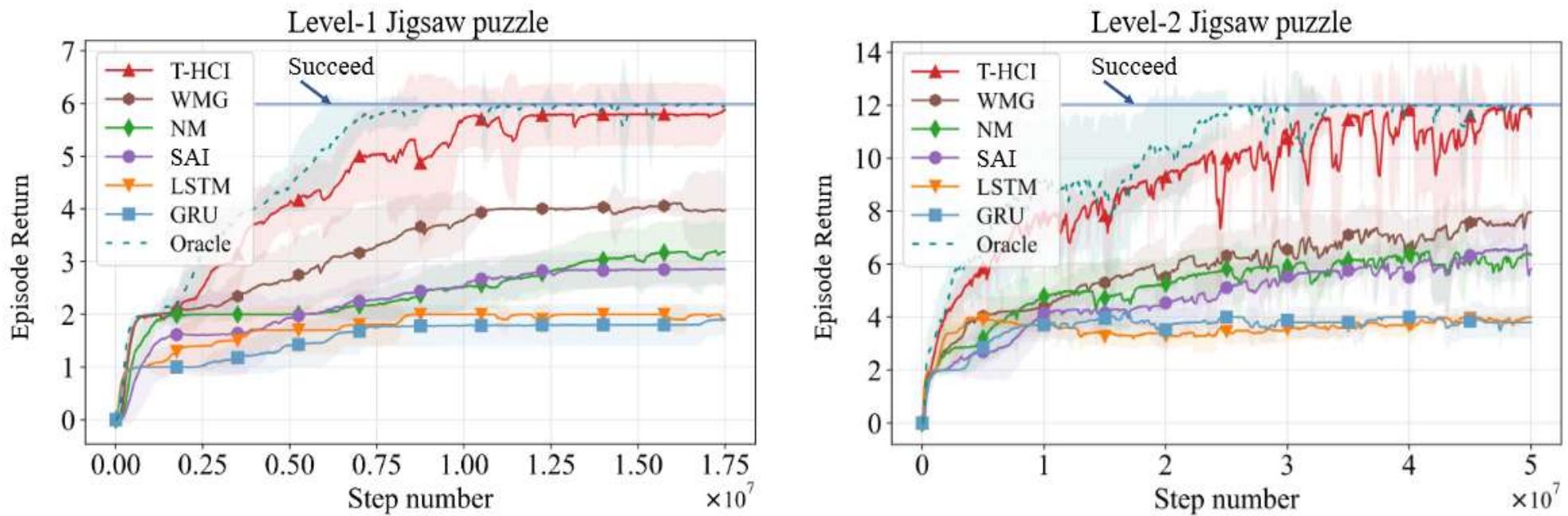


Figure 6: Learning curves of the Jigsaw puzzle tasks.