

# Interactive RL for Robotics

برديا طور انداز

علي رضا عباسي

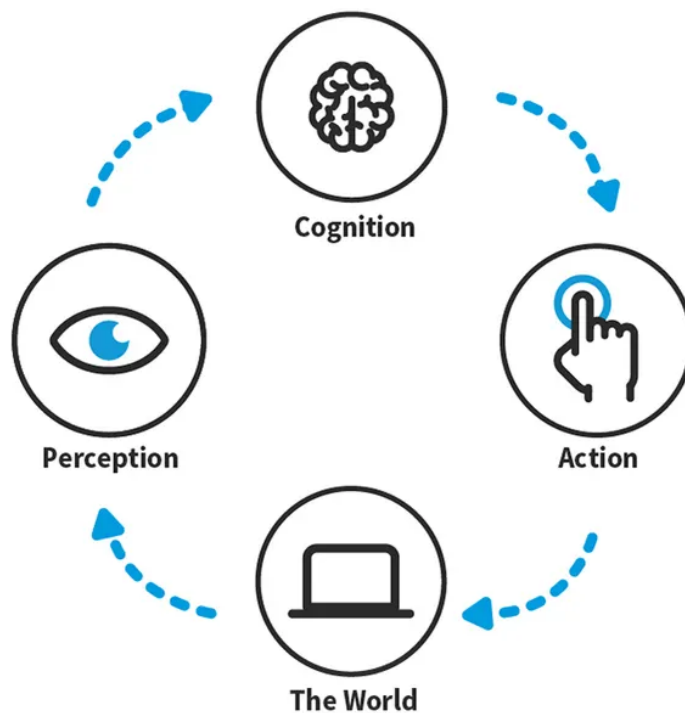




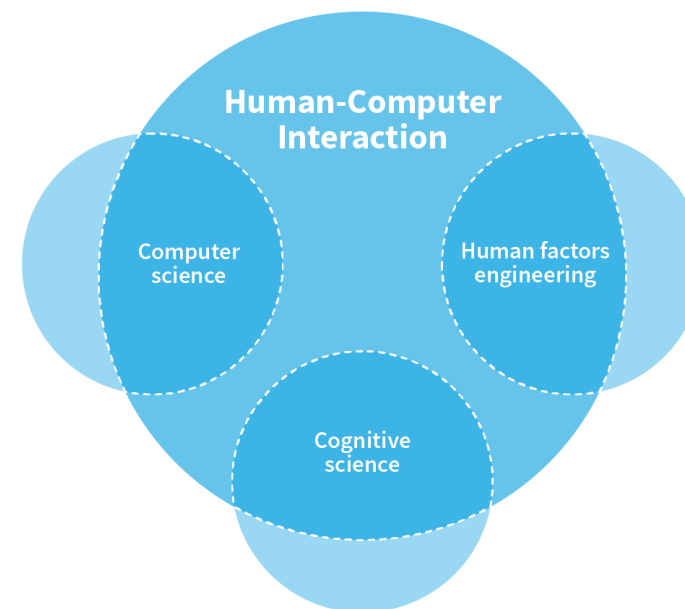


# یادگیری تقویتی تعاملی برای رباتیک

- انسان در حلقه (Human-in-the-loop)
- تعامل انسان و کامپیوتر



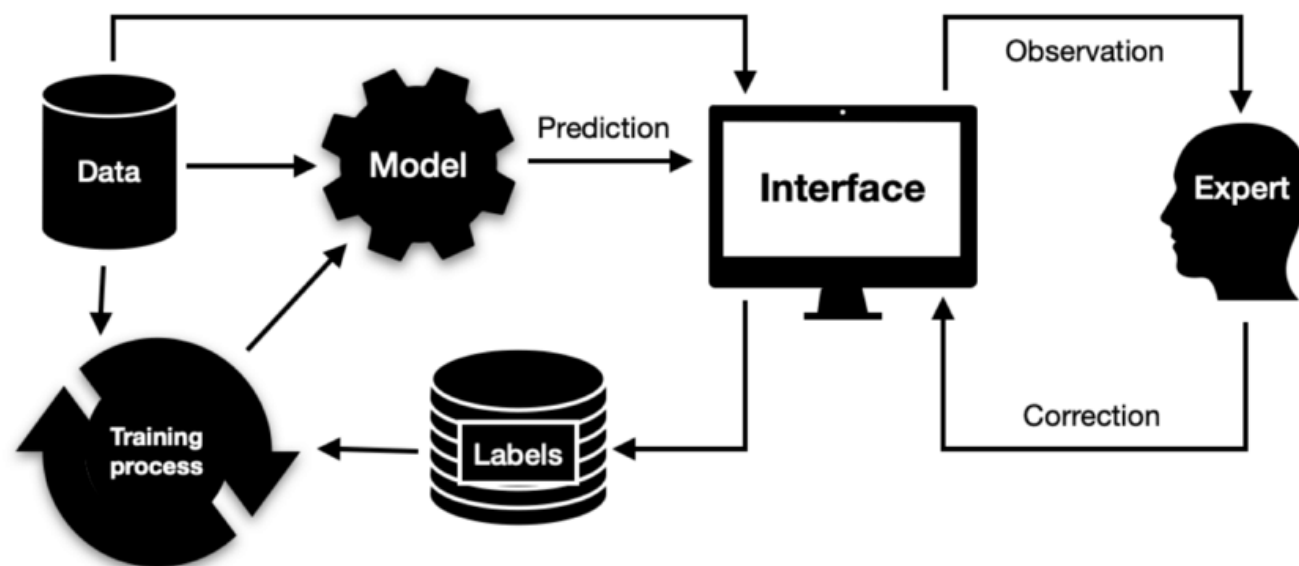
The Multidisciplinary Field of HCI





# وظایف مهندسين تعامل انسان و کامپیوتر در زمینه IML

- تعريف تکنیک‌های تعامل انسان و کامپیوتر
- پيدا کردن کاربردهای جديد





# ویژگیهای IRL

- سریعتر از Automatic RL
- Generalized نبودن پاسخها
- پیچیدگی تعامل با انسانهای متفاوت

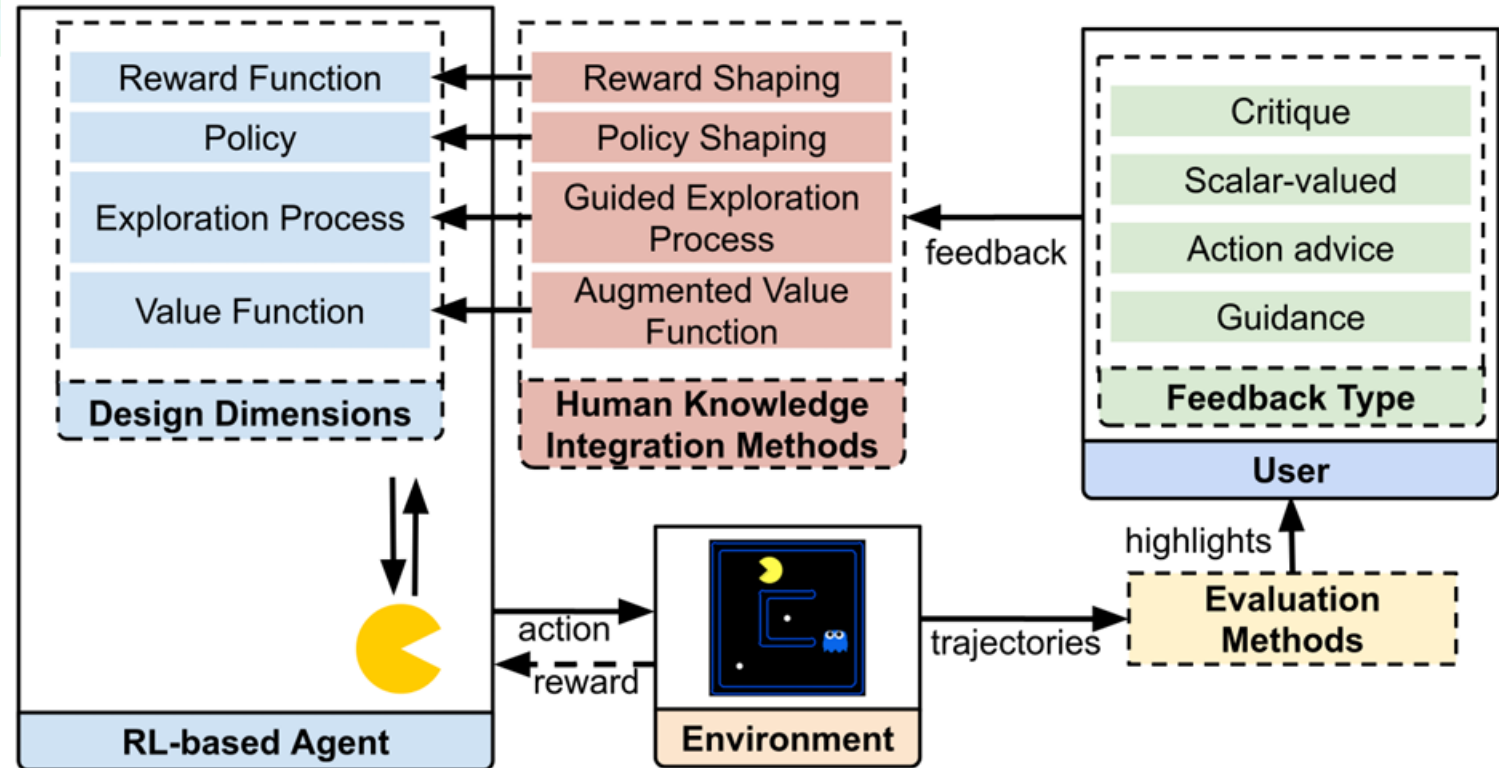
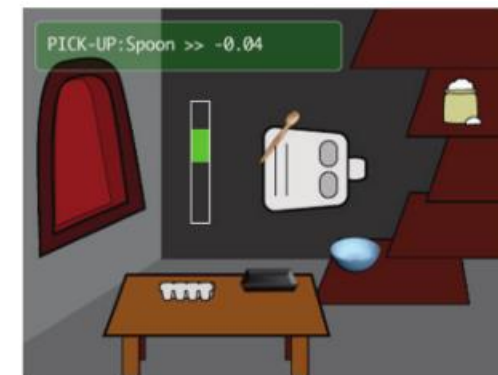
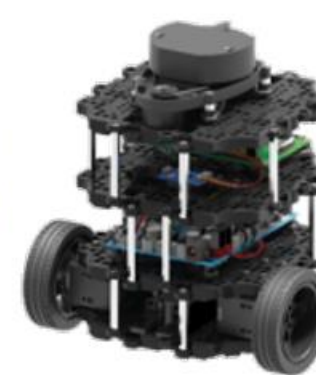
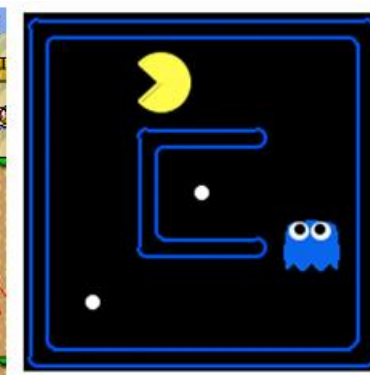
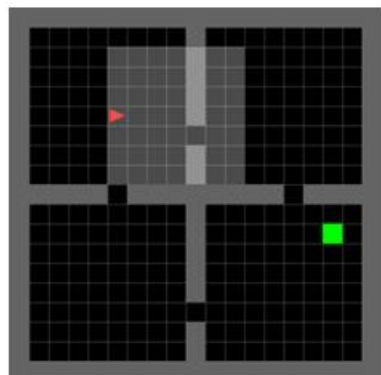


Figure 1. The interactive RL architecture.



# محیط‌های آزمایشی یادگیری تقویتی تعاملی

- رباتیک
- هوش مصنوعی در بازی‌ها
- تعامل انسان و کامپیوتر





# ابعاد طراحی

$R' = R + F$ , where  $F : S \times A \times S \rightarrow \mathbb{R}$  is the *shaping reward*

- تابع پاداش: فیدبک انسانی (توسط شخص خبره)، جهت تسریع یادگیری
- سیاست بهینه
- کاوش
- تابع ارزش

Design Dimension	Testbed	Interaction	Initiative	HKI	Feedback	Algorithms
Reward Function	Robot in maze-like environment [10]	FE	Passive	RS using HF + ER	Critique	DQL [67, 77]
	Navigation simulation [12]	GUI	Passive	Advantage Function	Critique	DAC
	Sophie's Kitchen game [105, 103]	GUI	Passive, Active	RS using HF + EF	Critique	QL [109], HRL
	Bowling game [108]	GUI	Passive	RS + HF	Scalar-valued	DQL
	Shopping assistant, GridWorld [76]	GUI	Active	Active IRD	Queries	Model-based RL
	Mario, GridWorld, Soccer simulation [85]	Coding	Passive	PBRs	Heuristic Function	QL, QSL, QL( $\lambda$ ), QSL( $\lambda$ )
	Navigation simulation [102]	VC	Passive	RS using HF + ER	AcAd	SARSA [86], SARSA( $\lambda$ )
	Atari, robotics simulation [20]	GUI	Active	RS using HF	Queries	DRL
Policy	GridWorld, TurtleBot robot [71]	GUI, GC	Passive	PS	AcAd	AC( $\lambda$ ) [15, 91]
	GridWorld [56]	VC	Passive	PS	Critique, AcAd	BQL [26]
	Pac-Man, Frogger [37]	GUI	Passive	PS	Critique	BQL
Exploration Process	Pac-Man, Cart-Pole simulation [114]	GUI	Passive	GEP	AcAd	QL
	Simulated cleaning Robot [24, 23]	VC	Passive	GEP	AcAd	SARSA
	Pac-Man [7]	GUI	Active	GEP	AcAd	SARSA( $\lambda$ )
	Pac-Man [32]	GUI	Active	Myopic Agent	AcAd	QL, QRL
	Sophie's Kitchen game [103]	GUI	Active	ACTG	Guidance	QL
	Street Fighter game [13]	Not apply	Passive	EB using Safe RL	Demonstration	HRL
	Nao Robot [93]	GUI	Passive	ACTG	Guidance	QL
	Nexi robot [54]	AT + CT	Passive	Myopic Agent	AcAd	SARSA( $\lambda$ )
Value Function	Mountain Car simulation [52]	GUI	Passive	Weighted VF	Demonstration	SARSA( $\lambda$ )
	Keepaway simulation [101]	GUI	Passive	Weighted VF	Demonstration	SARSA
	Mario, Cart Pole [18]	Not apply	Passive	Initialization of VF	Demonstration	QL( $\lambda$ )

# ابعاد طراحی

- طراحی پاداش:

The reward shaping (RS) method aims to mold the behavior of a learning agent by modifying its reward function to encourage the behavior the RL designer wants.

- طراحی سیاست:

The policy shaping (PS) approach consists of directly molding the policy of a learning agent to fit its behavior to what the RL designer envisions.

- فرآیند کاوش هدایت شده:

Guided exploration process methods aim to minimize the learning procedure by injecting human knowledge to guide the agent's exploration to states with a high reward.

- تابع ارزش تقویت شده:

The procedure to augment a value function consists of combining the value function of the agent with one created from human feedback.

# مسائل تحقیقاتی

High-dimensional  
Environments

محیط‌های با ابعاد بالا

Lack of Evaluation  
Techniques

کمبود روش‌های ارزیابی

Lack of Human-like  
Oracles

فقدان مشاوران شبیه به انسان

Modeling Users

مدل‌سازی کاربران

Combining Different  
Design Dimensions

ترکیب ابعاد مختلف طراحی

Safe Interactive RL

یادگیری تقویتی تعاملی ایمن

Fast Evaluation of  
Behaviors

ارزیابی سریع رفتارها

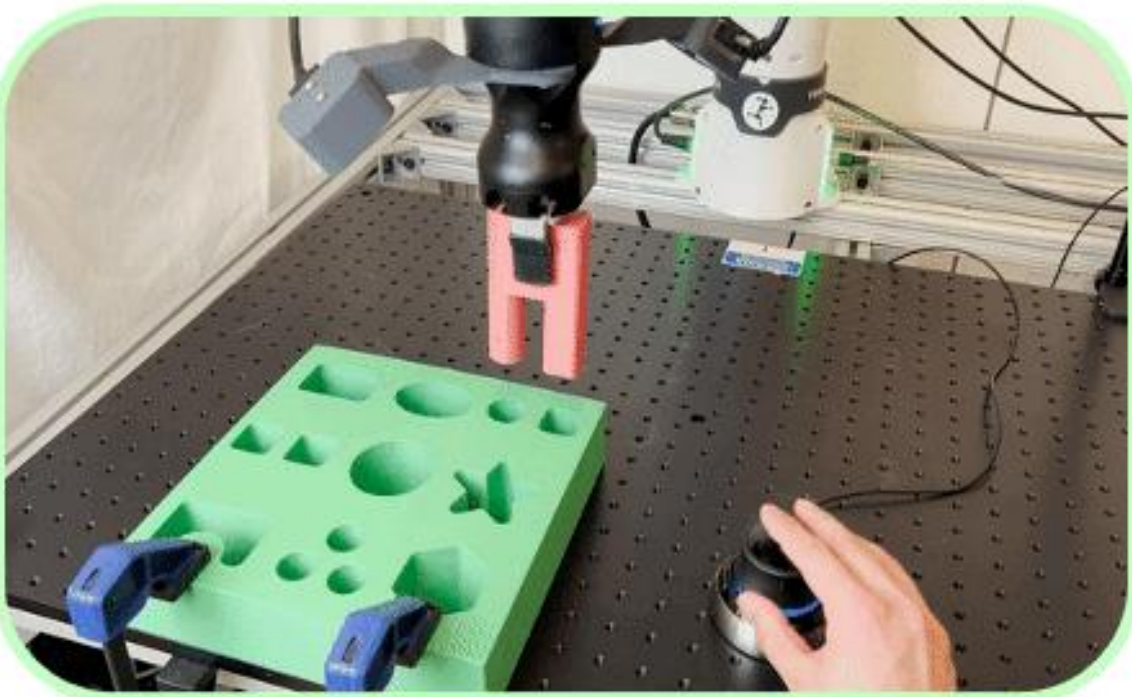
Explainable Interactive  
RL

یادگیری تقویتی تعاملی قابل  
توضیح



# یادگیری تقلیدی تعاملی در یادگیری تقویتی

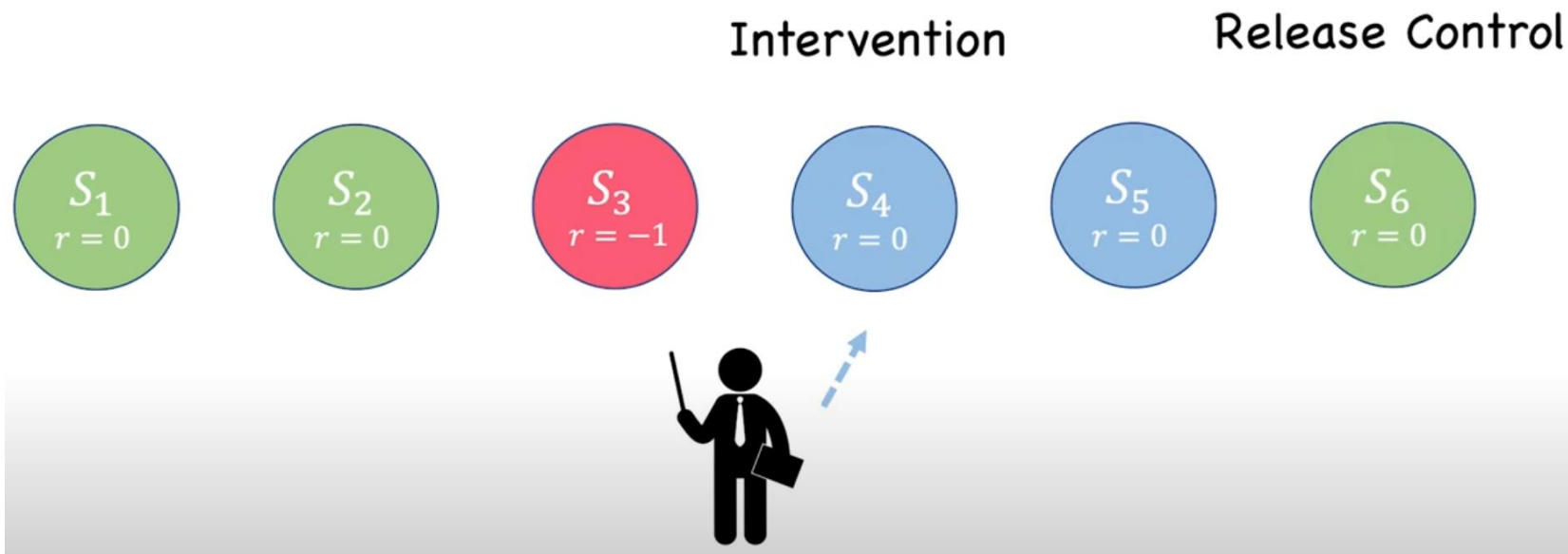
- اعمال فیدبک توسط انسان در صورت بهینه نبودن پالیسی
- اعمال یادگیری تقویتی با فیدبک تعاملی بعنوان سیگنال **ریوارد**

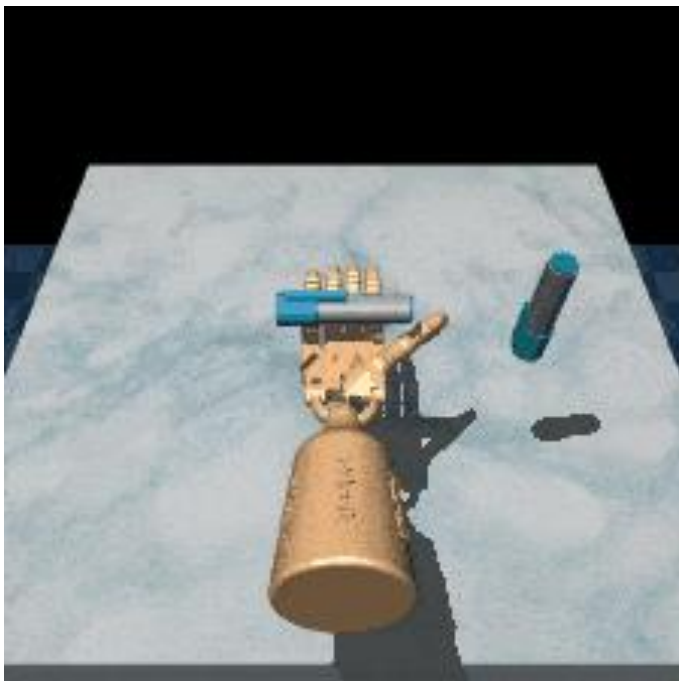


## Algorithm 2 RLIF

**Require:**  $\pi, \pi^{\text{exp}}, D$

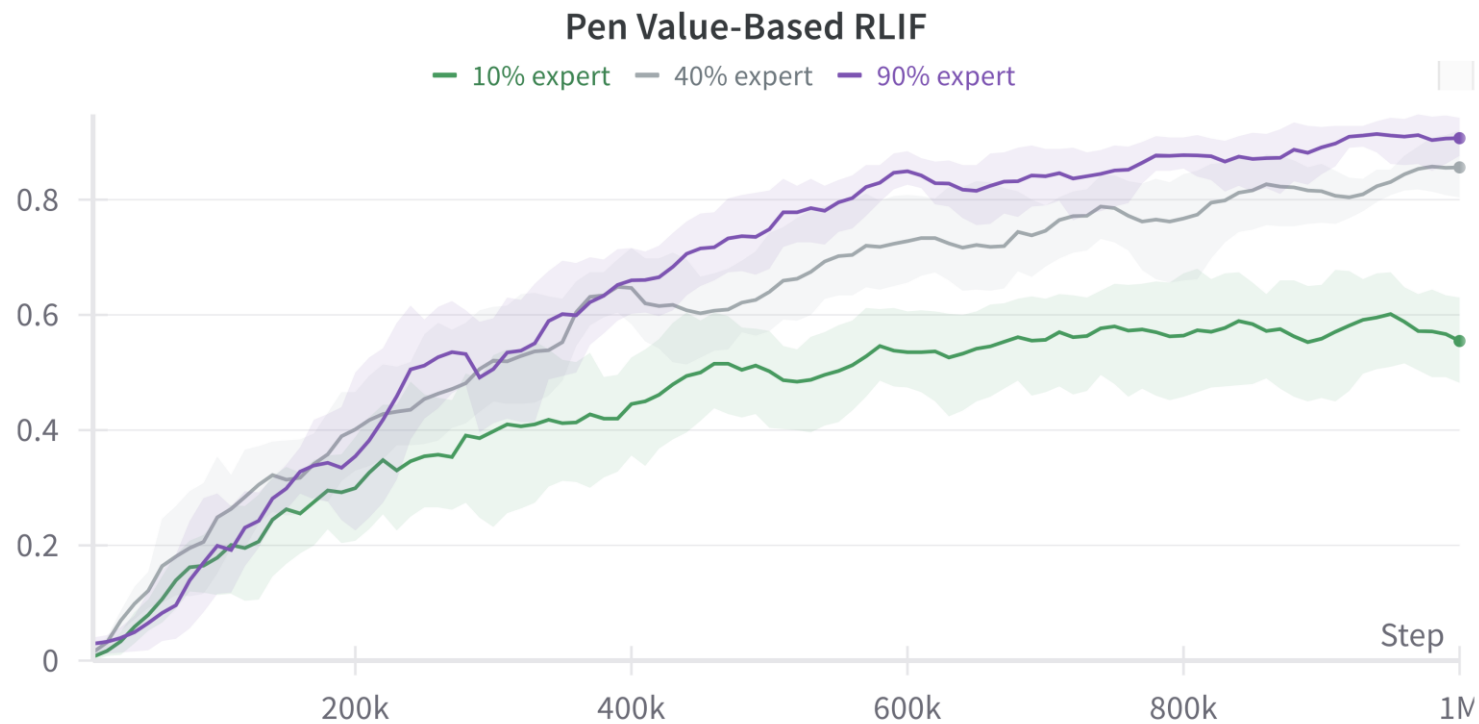
```
1: for trial  $i = 1$  to  $N$  do
2:   Train  $\pi$  on  $D$  via reinforcement learning.
3:   for timestep  $t = 1$  to  $T$  do
4:     if  $\pi^{\text{exp}}$  intervenes at  $t$  then
5:       label  $(s_{t-1}, a_{t-1}, s_t)$  with -1 reward,
       append to  $D_i$ 
6:     else
7:       label  $(s_{t-1}, a_{t-1}, s_t)$  with 0 reward,
       append to  $D_i$ 
8:     end if
9:   end for
10:   $D \leftarrow D \cup D_i$ 
11: end for
```



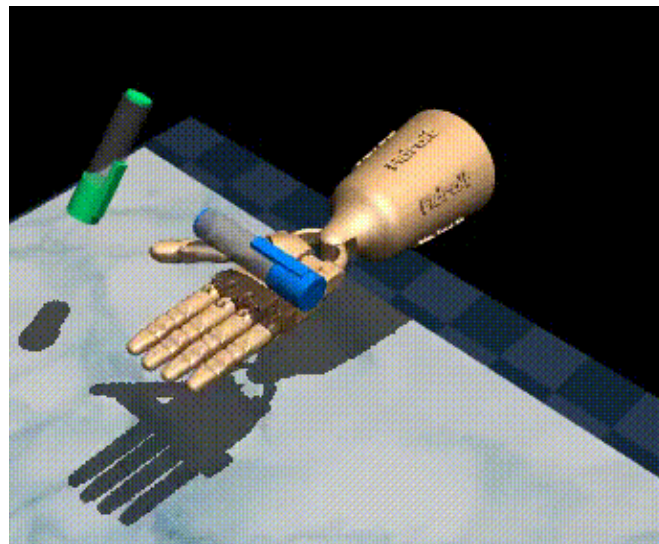


# یادگیری تقلیدی تعاملی در یادگیری تقویتی

- **D4RL** : ابزار متن-باز برای بنچ-مارک الگوریتم‌های آفلاین یادگیری تقویتی
- محیط **Adroit Pen**







# یادگیری تقلیدی تعاملی در یادگیری تقویتی

- استراتژی مداخله:

- random : مداخله در ۳۰ و ۵۰ و ۸۵ درصد زمان‌ها

- بر پایه value : مقایسه با پالیسی خبره (expert)

$$\mathbb{P}(\text{Intervention}|s) = \begin{cases} \beta, & \text{if } Q^{\pi^{\text{ref}}}(s, \pi^{\text{exp}}(s)) * \alpha > Q^{\pi^{\text{ref}}}(s, \pi(s)) \\ 1 - \beta, & \text{otherwise.} \end{cases}$$

We choose a value for  $\beta$  close to 1 such as 0.95 and a value of  $\alpha$  close to 1 such as 0.97.

- پالیسی خبره:

- ترین روی دیتاست D4RL در ابعاد کوچکتر (۱۰ و ۴۰ و ۹۰ درصد)

Tasks	Dataset	Subsampled Size
Adroit Pen	pen-expert-v1	50 trajectories

# RLIF: Interactive Imitation Learning as Reinforcement Learning

## یادگیری تقلیدی تعاملی در یادگیری تقویتی

Jianlan Luo\* Perry Dong\* Yuexiang Zhai Yi Ma Sergey Levine

International Conference on Learning Representations (ICLR) 2024

Vienna, Austria

• نتیجه شبیه سازی



Paper



Code

Domain	Expert Level	RLIF with Value Based Intervention	RLIF with Random Intervention	HG-DAgger	HG-DAgger with 85% Random Intervention	DAgger	DAgger with 85% Random Intervention	BC
adroit-pen	~90%	<b>88.47</b>	42.87	73.47	74.27	78.13	79.07	54.13
	~40%	<b>80.87</b>	34.13	60	29.33	35.73	38.67	
	~10%	<b>64.04</b>	28.33	28.53	9.47	8.93	12.8	
	average	<b>77.79</b>	35.11	54	37.69	40.93	43.51	
locomotion-hopper	~110%	108.99	106.51	53.55	<b>112.7</b>	57.94	76.13	44.46
	~70%	<b>99.66</b>	75.62	44.75	69.73	20.49	43.59	
	~20%	<b>102.85</b>	19.11	11.94	19.66	12.37	20.1	
	average	<b>103.83</b>	67.08	36.75	67.36	30.27	46.61	
locomotion-walker	~110%	<b>109.17</b>	93.76	80.3	86.93	70.58	61.64	64.77
	~40%	<b>108.42</b>	103.9	40.66	42.65	38.7	19.63	
	~15%	<b>108.01</b>	75.12	25.2	24.37	19.54	10.29	
	average	<b>108.53</b>	90.93	48.72	51.32	42.94	30.46	



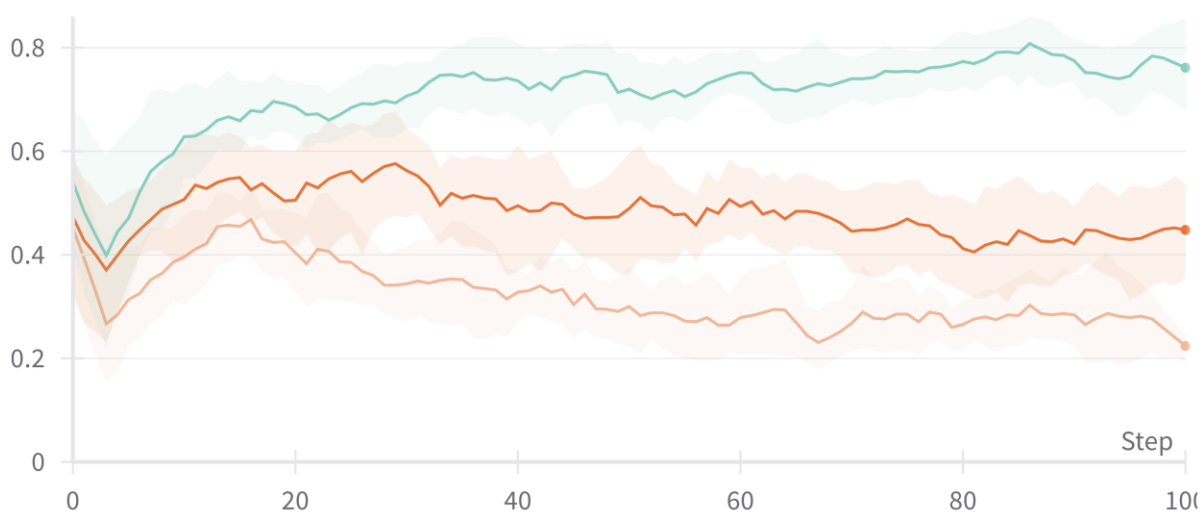
# یادگیری تقلیدی تعاملی در یادگیری تقویتی

```
bash scripts/run_pen_hgdagger.sh
bash scripts/run_pen_value_based_rlif.sh
```

• نتیجه شبیه‌سازی

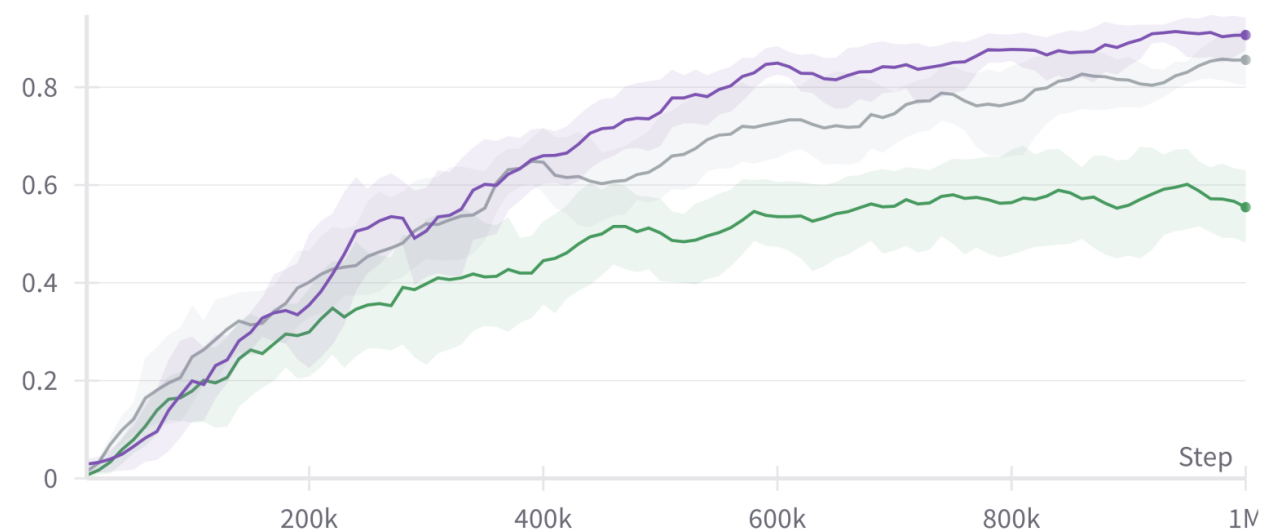
Pen Random Intervention DAGger

10% expert 40% expert 90% expert

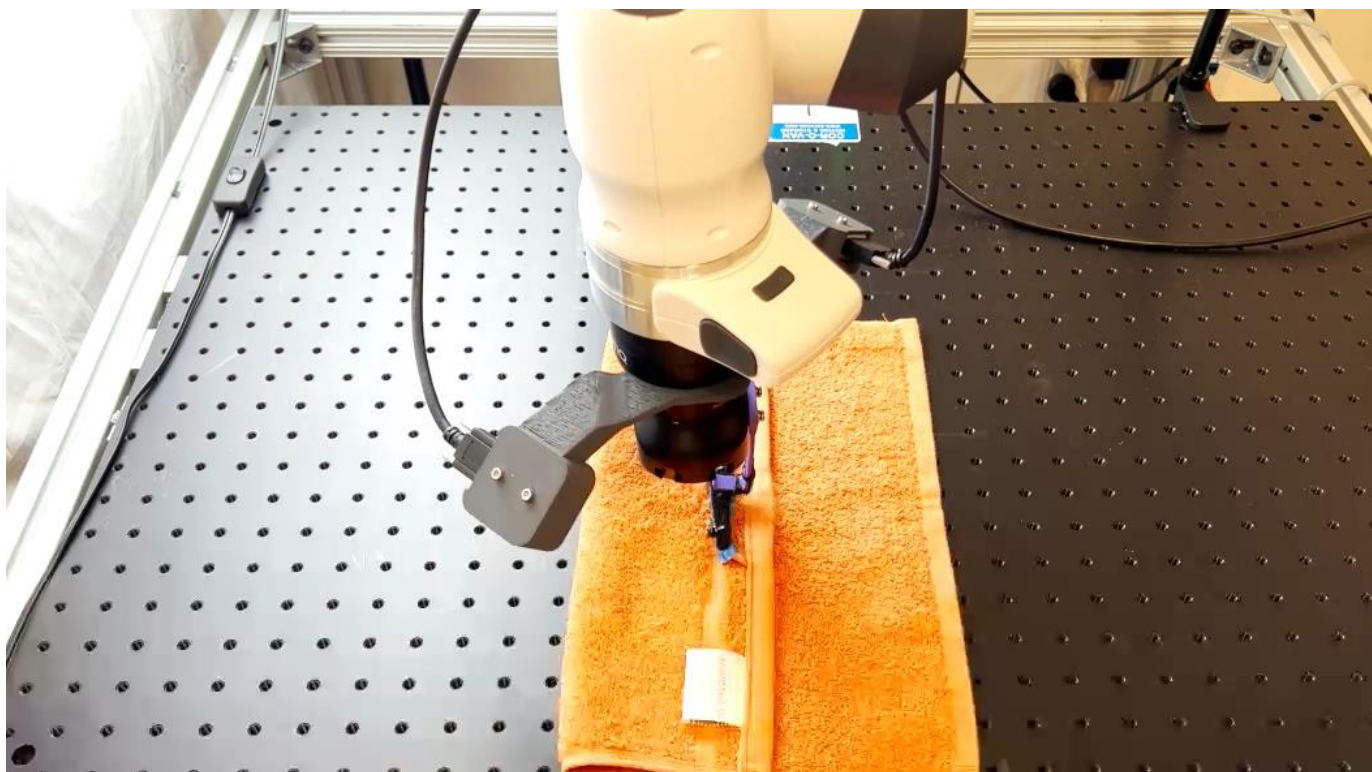


Pen Value-Based RLIF

10% expert 40% expert 90% expert



# یادگیری تقلیدی تعاملی در یادگیری تقویتی



- اعمال روی ربات حقیقی
  - داینامیک غیرپیوسته و غیرقابل پیش‌بینی برخورد
  - ورودی سنگین (تصویر دوربین)
  - اعمال روی پلتفرم با شکل متغیر (حوله)
  - اپراتور مداخله‌گر (انسان) غیر ایده‌آل
- انتخاب ریوارد بصورت دستی بسیار دشوار است.
- نتیجه: ۹۵ درصد نرخ موفقیت در ۲۰ مسیر سنجش