

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه صنعتی اصفهان

دانشکده برق و کامپیوتر

محمدعلی میرزایی بادیزی

استاد راهنما: دکتر پالهننگ

تسریع یادگیری مشارکتی در سیستم‌های چند عاملی

با بهره‌گیری از کوتاهترین مسیر تجربه شده

فهرست مطالب

مقدمه

- یادگیری تقویتی

مرور بر کارهای پیشین

- تقلید
- تخته سیاه
- حافظه جمعی
- پند دهی
- خبرگی
- خبرگی چند معیاره

ارائه روش پیشنهادی

- معرفی معیار شوک
- معرفی معیار کمترین فاصله
- تجربه شده
- نحوه استفاده از معیار ها

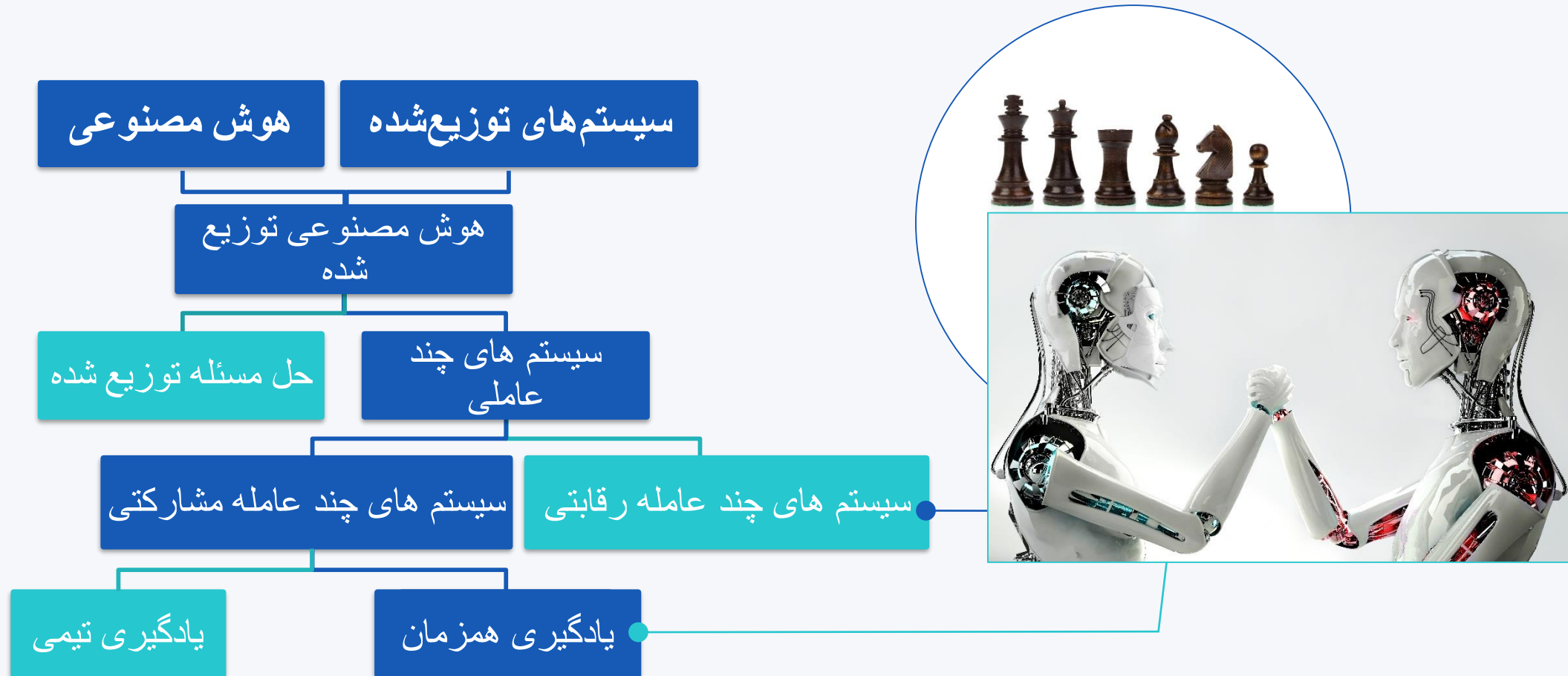
نتیجه گیری

- ارائه نتیجه آزمایشات
- ارائه پیشنهادهایی جهت کار های آتی

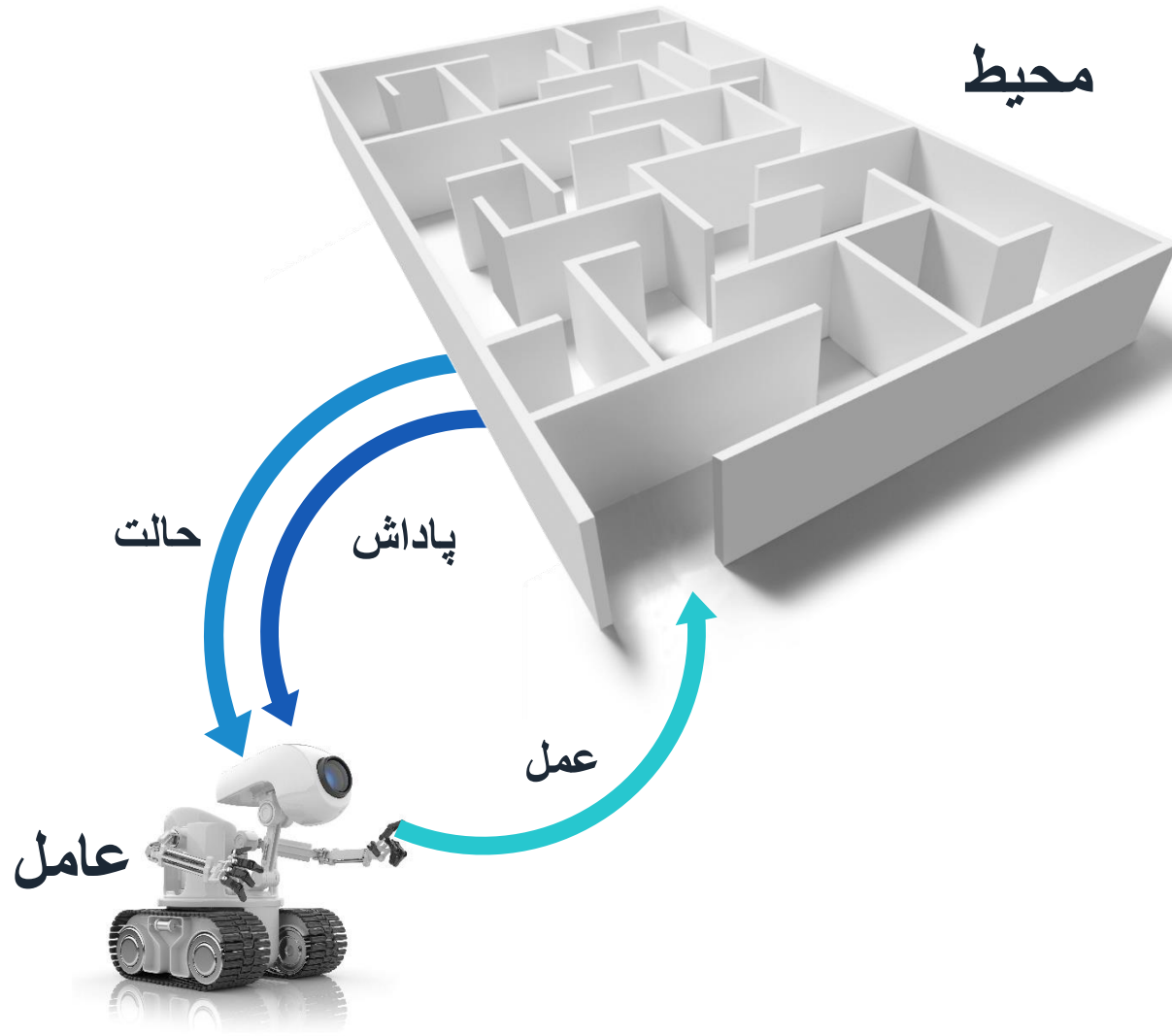


مقدمه

جایگاه پژوهش انجام شده



یادگیری تقویتی



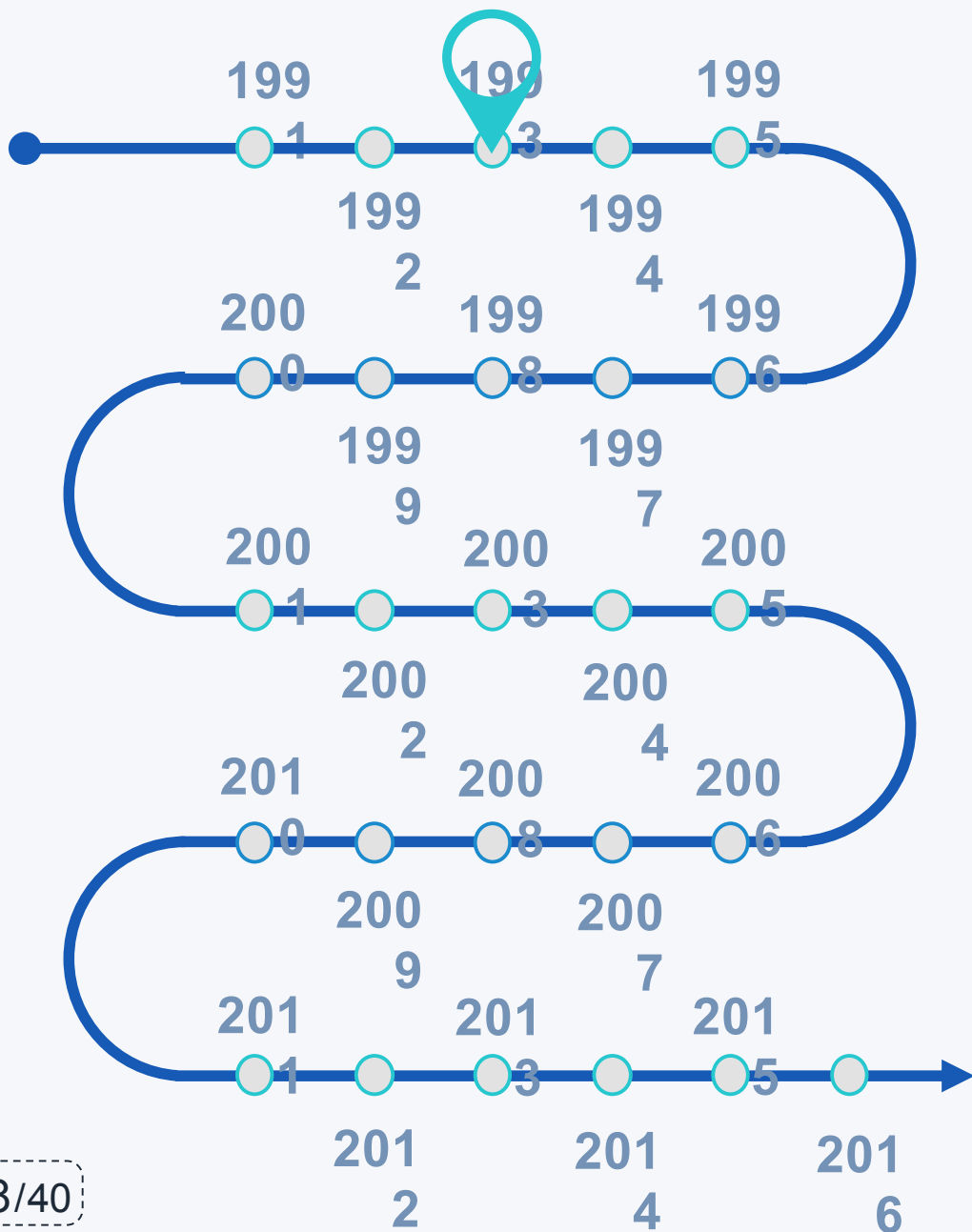
Q learning

1				
2				
3				
4				
5				
6				
	↓	↑	←	→

ارائه روشهای

پیشین



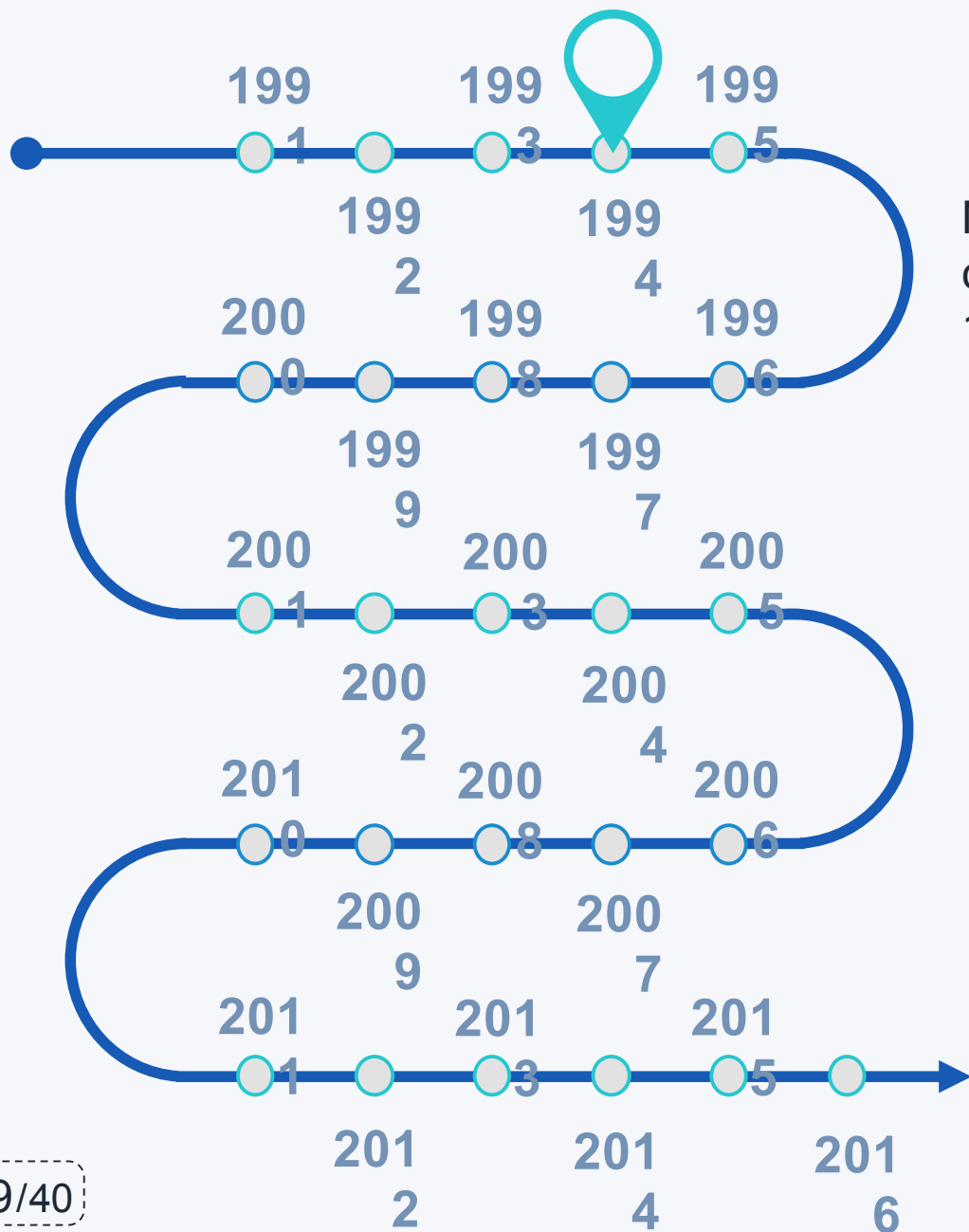


M. Tan, "Multi-Agent Reinforcement Learning: Independent vs. Cooperative Agents," Proc. Tenth Int. Conf. Mach. Learn., pp. 330–337, 1993.



اگر مشارکت به‌خوبی پیاده‌سازی شود هر عامل می‌تواند از تجربیات
عامل‌های دیگر استفاده بهینه نماید.



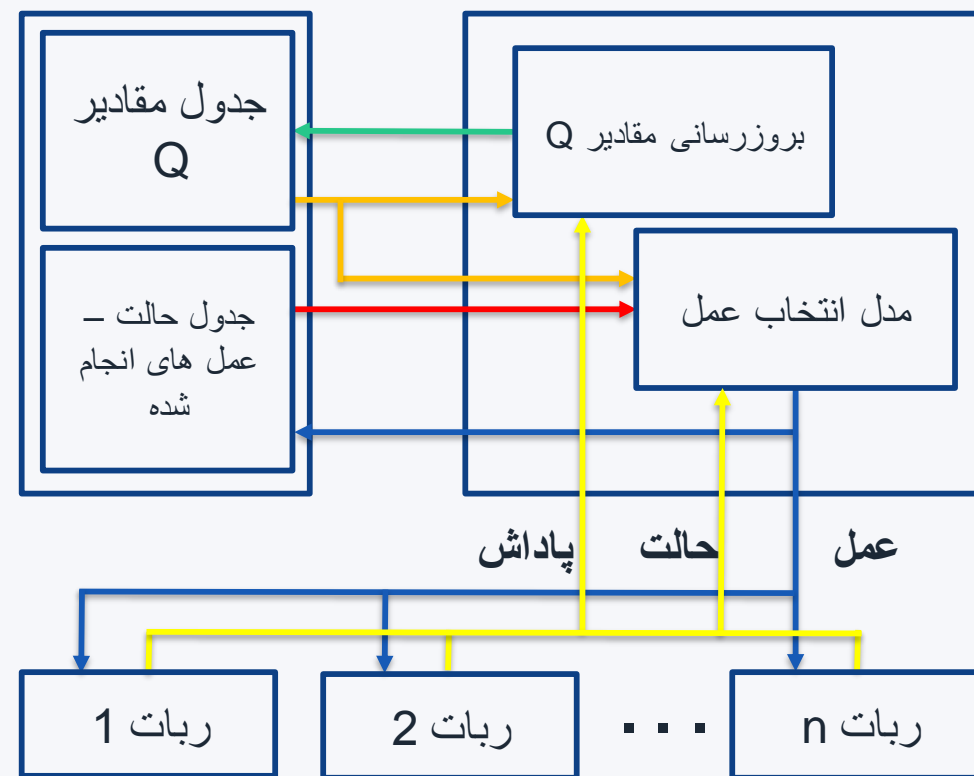


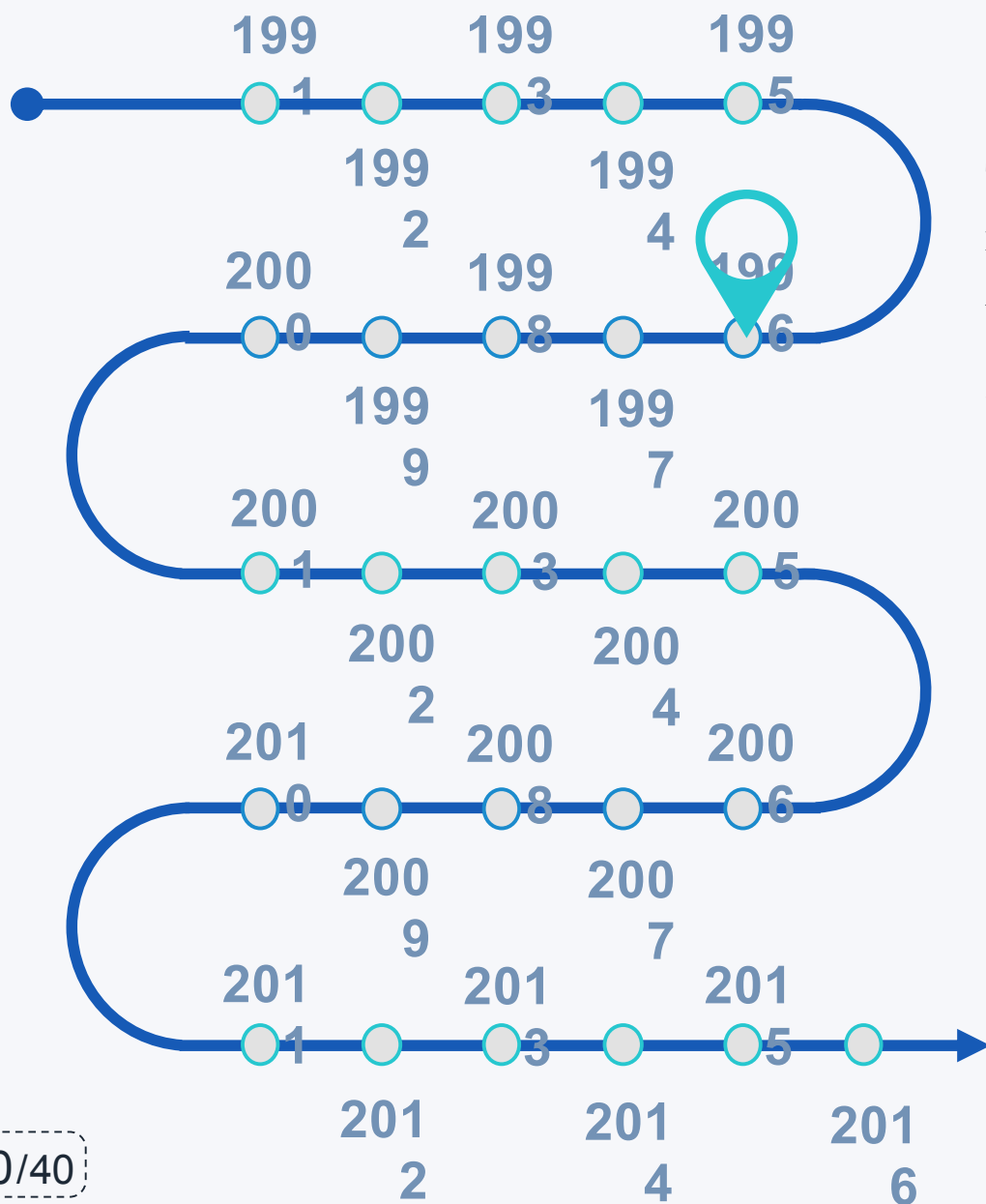
تخته سیاه



Blackboard

N. Carver and V. Lesser, "Evolution of blackboard control architectures," *Expert Syst. Appl.*, vol. 7, no. 1, pp. 1–30, Jan. 1994.



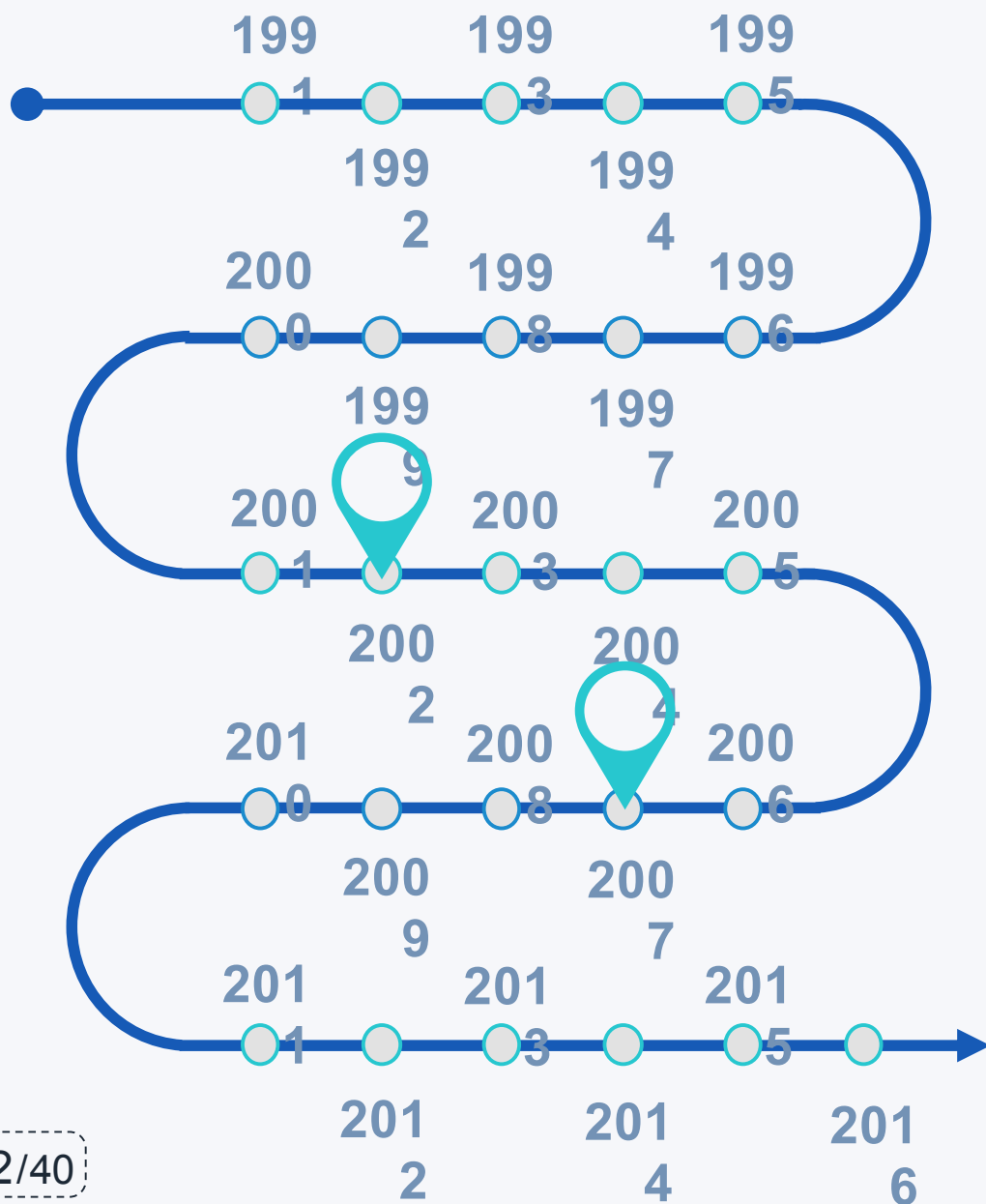


تقلید

Imitation

حافظہ جمعی

CM



پند



Advice

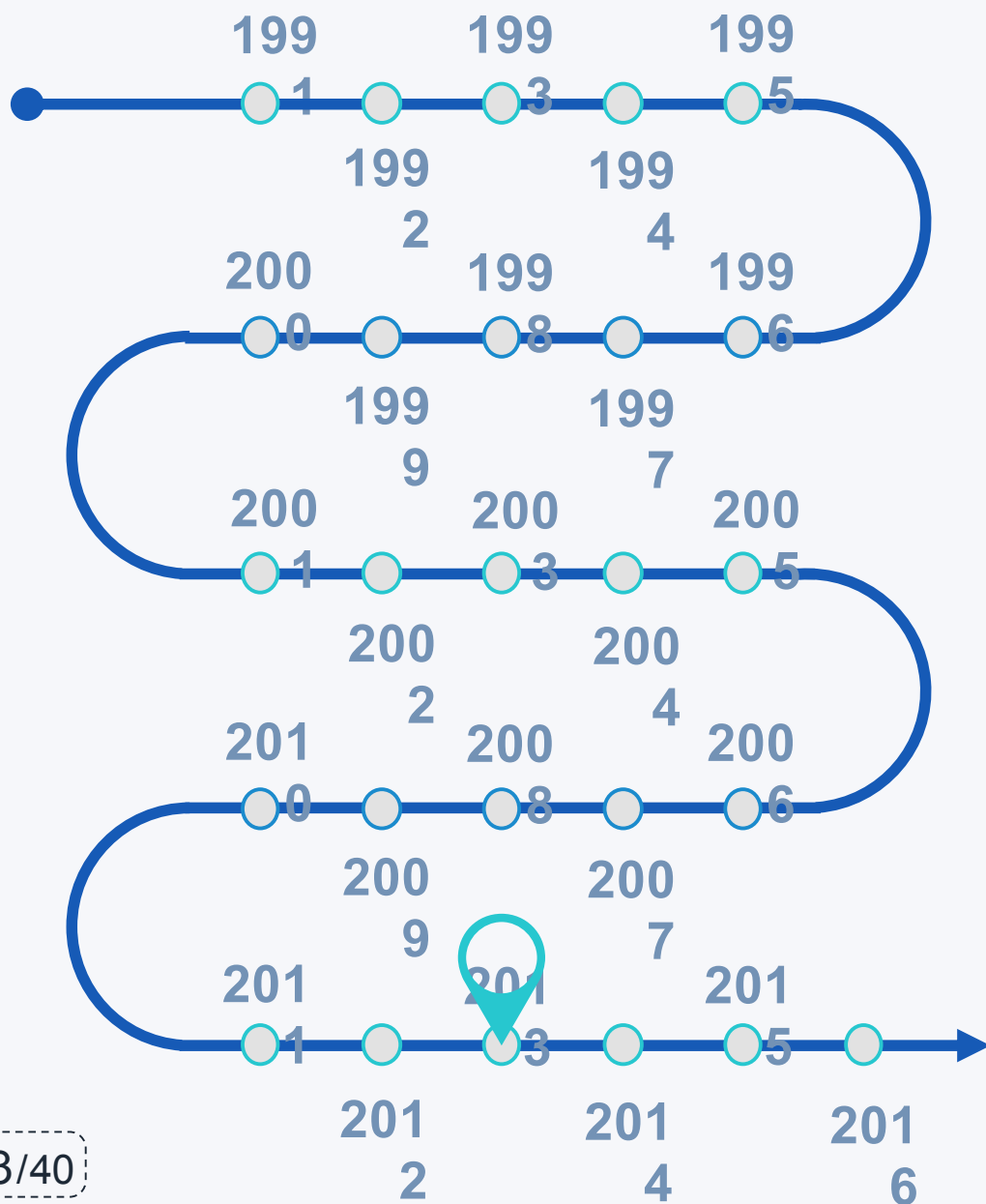
L. Nunes and E. Oliveira, "On Learning by Exchanging Advice," Proc. Artif. Intell. Simul. Behav. Conv. Symp. Adapt. agents multi-agent Syst. (AISB/AAMAS-II), Imp. Coll. london, vol. cs.LG/0203, pp. 583–599

مبتهی بر پختگی سیاست



Blackboard

Y. Yang, Y. Tian, and H. Mei, "Cooperative Q Learning Based on Blackboard Architecture," in 2007 International Conference on Computational Intelligence and Security Workshops (CISW 2007), 2007, pp. 224–227.

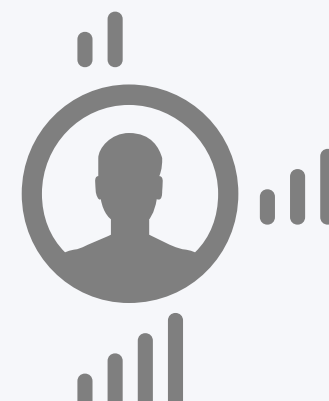
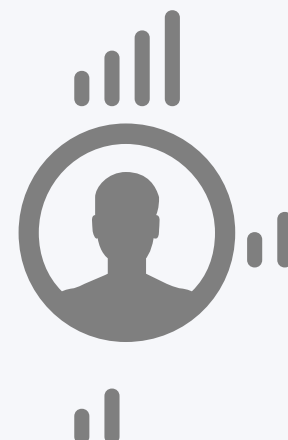


خبرگی چند معیاره

E. Pakizeh, M. Palhang, and M. M. Pedram, "Multi-criteria expertness based cooperative Q-learning," *Appl. Intell.*, vol. 39, no. 1, pp. 28–40, Jul. 2013.



MCE



1: معمولی

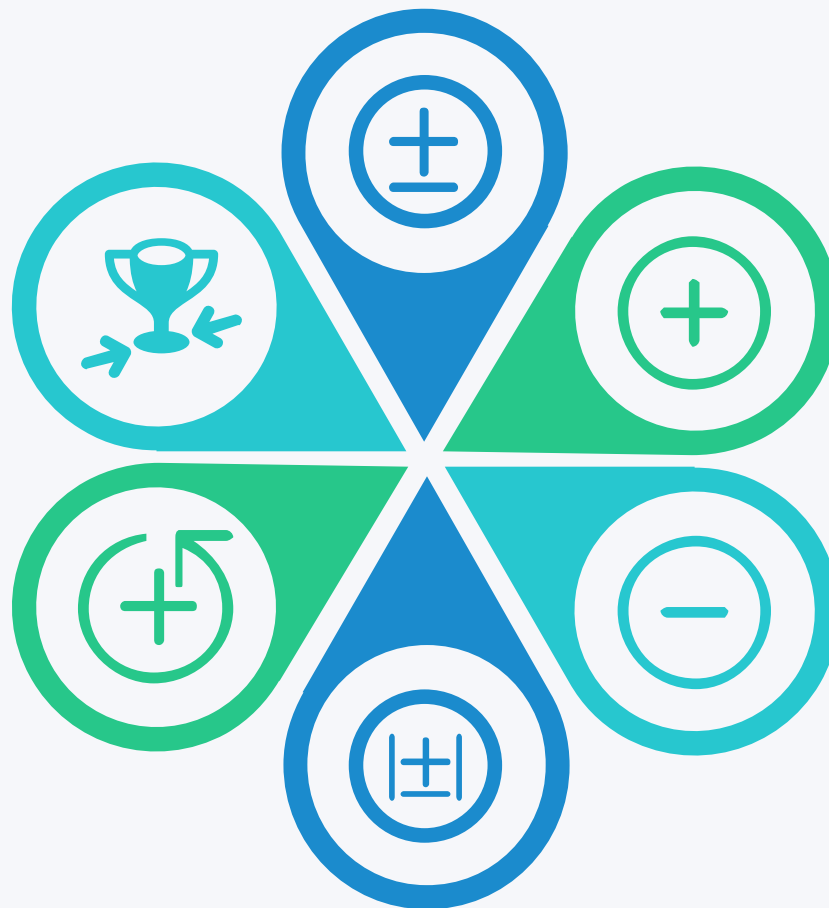
در نظر گرفتن پاداش و جریمه

2: مثبت

در نظر گرفتن پاداش

3: منفی

در نظر گرفتن جریمه



6: میانگین فاصله

ارزیابی بر اساس میانگین

تعداد قدم های برداشته

برای رسیدن به هدف

5: گرادیان

در نظر گرفتن پاداش تنها

در آخرین چرخه یادگیری

4: قدر مطلق

در نظر گرفتن جریمه و پاداش

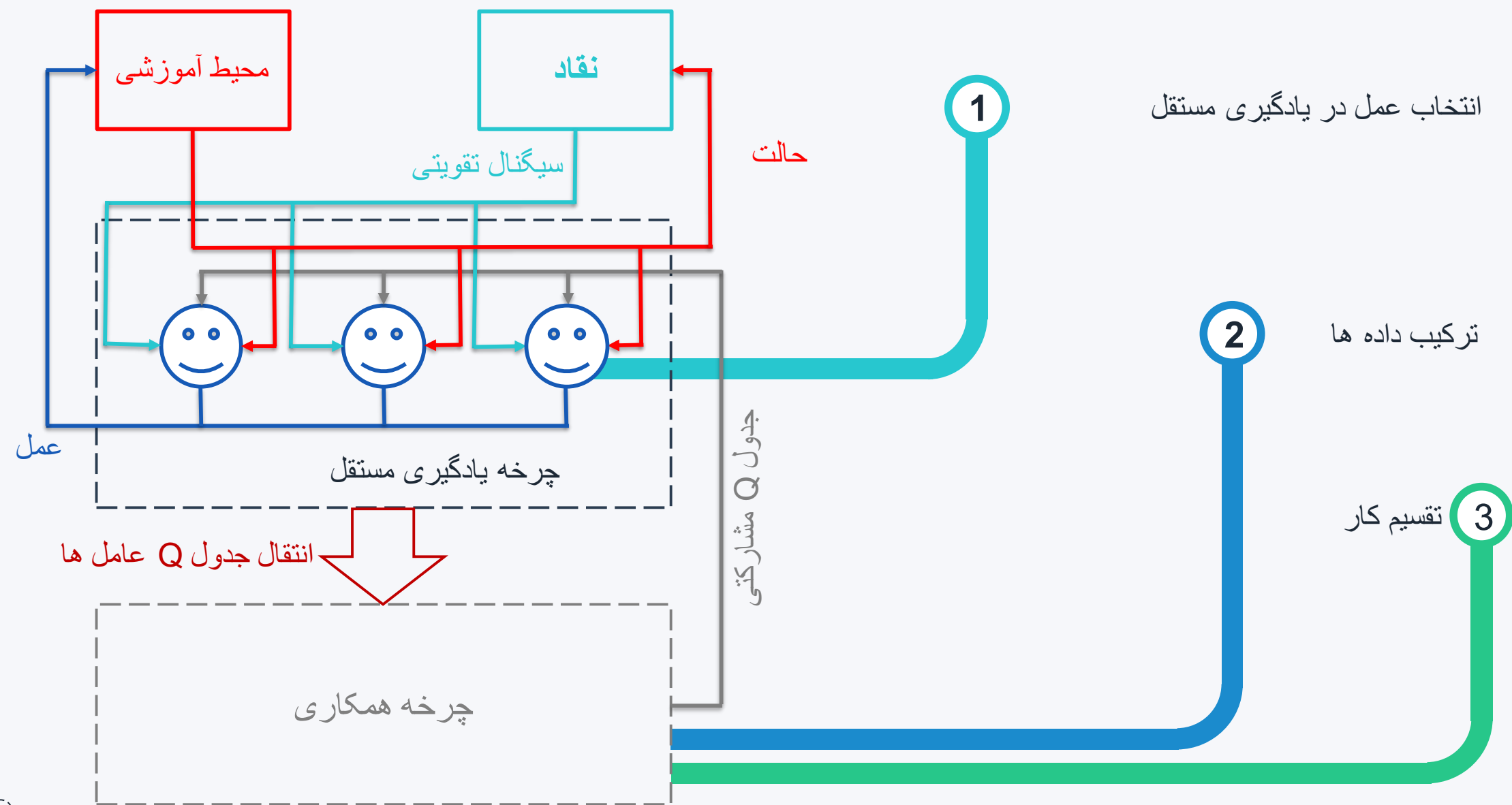
به عنوان دو پارامتر مثبت

تشریح

روش

پیشنهادی





معیار ارزیابی ارائه شده

عنوان: شوک

محاسبه بر اساس: تعداد دفعاتی که پاداش یکی از

خانه های جذب در یک خانه از ماتریس Q اثر داشته است.

معیار ارزیابی ارائه شده



0	0	0	0	1
0	0	0	0	2
0	0	0	0	3
0	0	0	0	4
0	0	0	0	5
0	0	0	0	6
0	0	0	0	7
0	0	0	0	8
0	0	0	0	9

جدول Q



0	0	0	1	1
0	0	0	2	2
0	1	1	1	3
0	0	0	0	4
1	0	1	0	5
0	0	0	0	6
0	0	0	0	7
1	0	0	0	8
0	0	0	0	9

جدول مشاهده



0	0	0	0	1
0	0	0	0	2
0	0	0	1	3
0	0	0	0	4
0	0	0	0	5
0	0	0	0	6
0	0	0	0	7
0	0	0	0	8
0	0	0	0	9

جدول شوک

1 	2	3	4
		5	
6	7	8	9

معیار ارزیابی ارائه شده

عنوان: کوتاهترین فاصله تجربه شده

محاسبه بر اساس: کوتاهترین طول مسیر تا هر یک

از اهداف وبا انتخاب هر عمل که مشاهده شده است

معیار ارزیابی ارائه شده

عنوان: کوتاهترین فاصله تجربه شده

محاسبه بر اساس: کوتاهترین طول مسیر تا هر یک

از اهداف وبا انتخاب هر عمل که مشاهده شده است

4	3	2	1	
∞	∞	∞	3	1
∞	∞	∞	2	2
∞	∞	∞	1	3
∞	∞	∞	∞	4
2	∞	∞	∞	5
∞	∞	∞	∞	6
∞	∞	∞	∞	7
3	∞	∞	∞	8
∞	∞	∞	∞	9



2	3	4	0	3	0	0	5	0
1	1	1	0	4	0	0	4	0
0	0	0	0	0	0	0	0	0
1	2	3	4	5	6	7	8	9

انتخاب عمل در یادگیری مستقل

$$\pi(s)_t = \text{Boltzmann}(Q_t, \tau_1)$$

$$p_i = \frac{e^{-\varepsilon_i/kT}}{\sum_{j=1}^M e^{-\varepsilon_j/kT}}$$

$$\pi(s)_t = (1 - \mu) \text{Boltzmann}(Q_t, \tau_1) + \mu \text{Boltzmann}\left(\frac{1}{\text{SEP}_t}, \tau_2\right)$$

تقسیم کار



$$\pi(SEP) == \pi(Q)$$



$$\pi(SEP) \neq \pi(Q)$$



$$\pi(SEP) \neq \pi(Q)$$



$$\pi(SEP) == \pi(Q)$$



$$\pi(SEP) == \pi(Q)$$

$Group_1 =$

$Group_2 =$

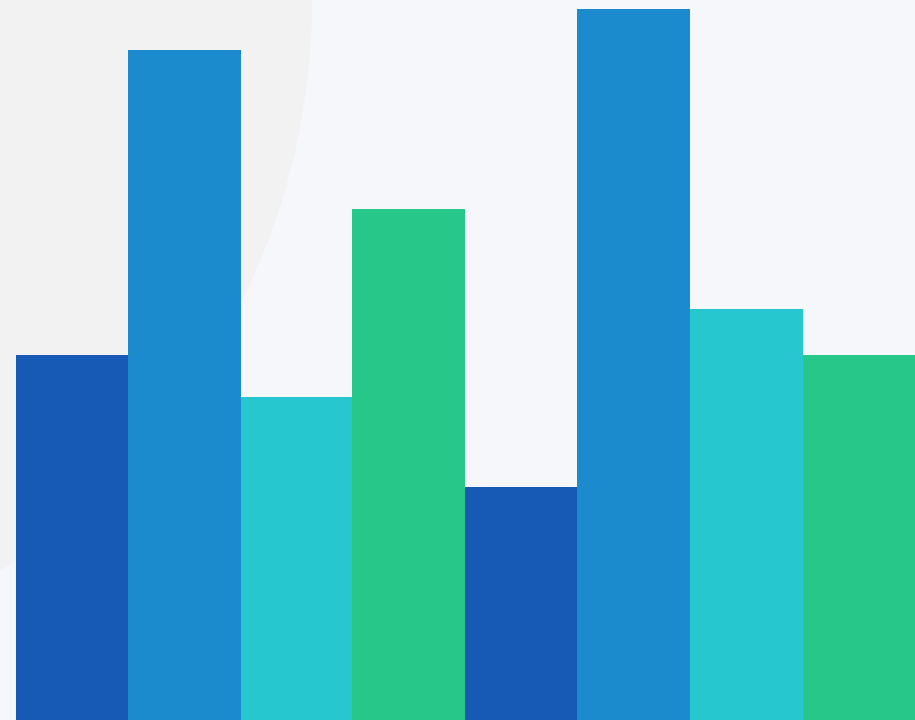
ترکیب داده ها

$$Q_{co}(s) = \sum_{i \in G} \frac{\sum_{k \in actionSet} shock_i(s, k)}{\sum_{j \in G} \sum_{k \in actionSet} shock_j(s, k)} Q_i(s)$$

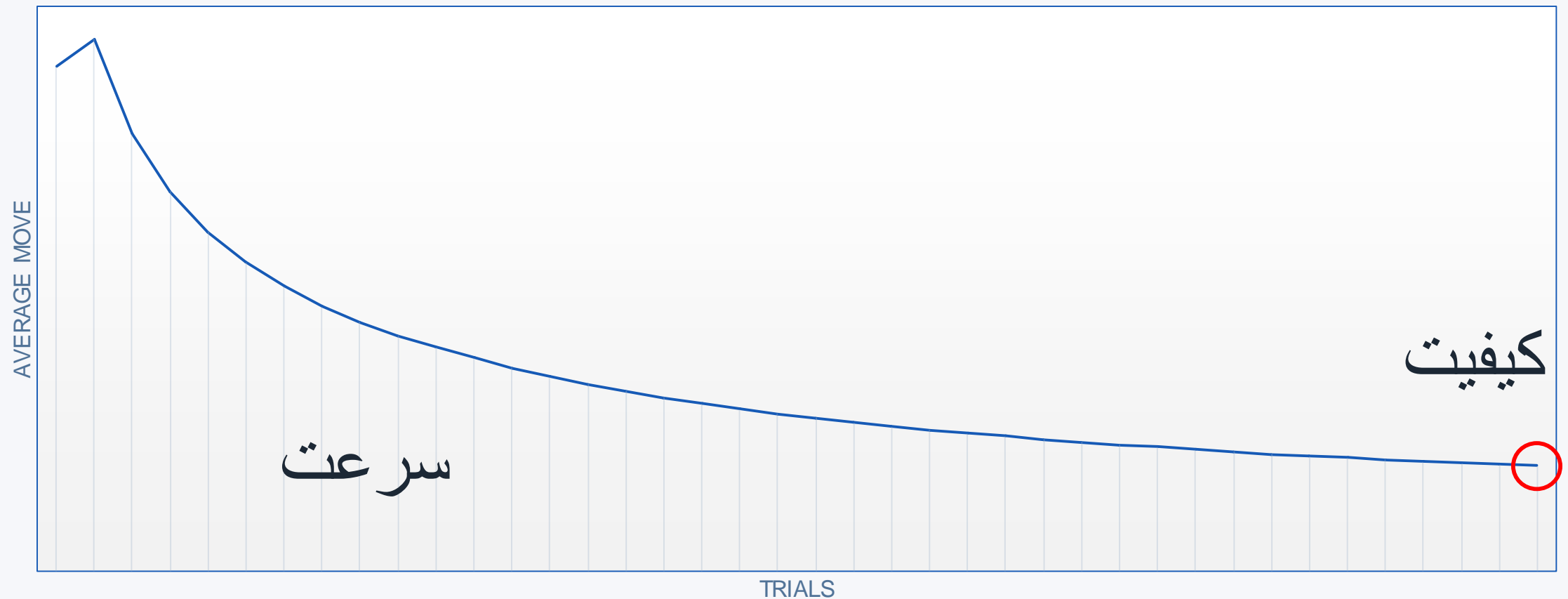
SEP

نتیجه

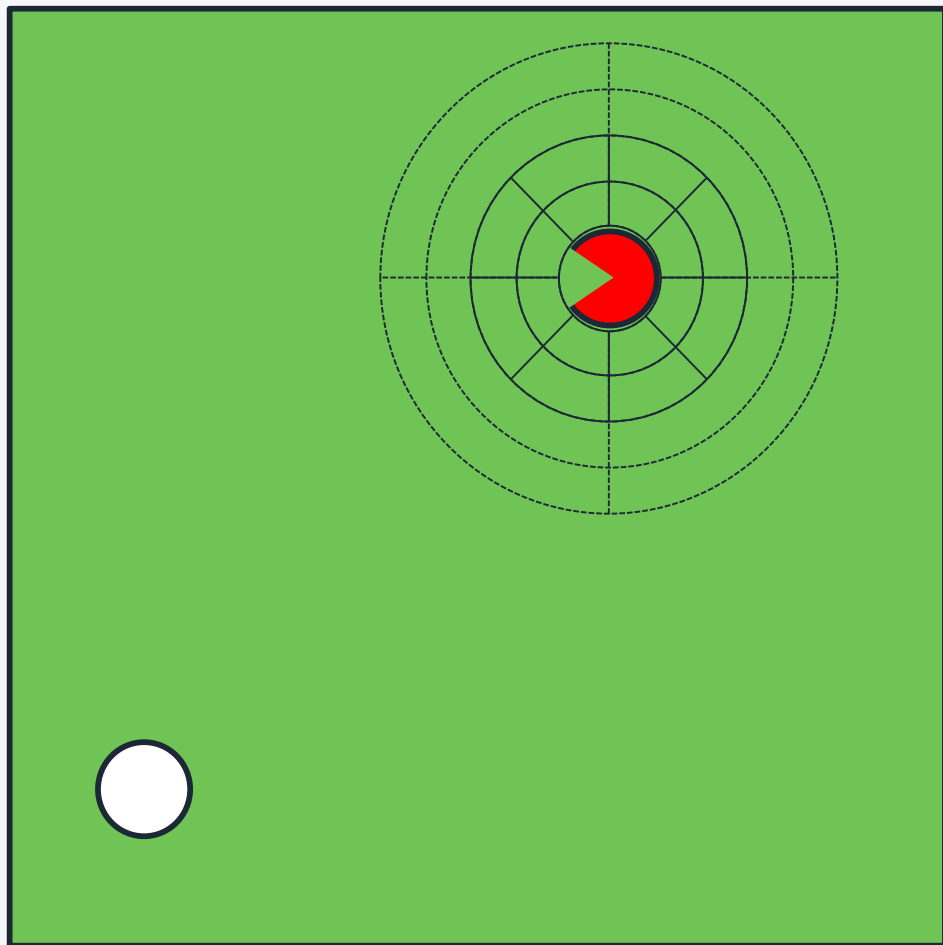
آزمایشات



بهره گیری از کوتاهترین فاصله تجربه شده در یادگیری مستقل



صید و صیاد



16 عمل

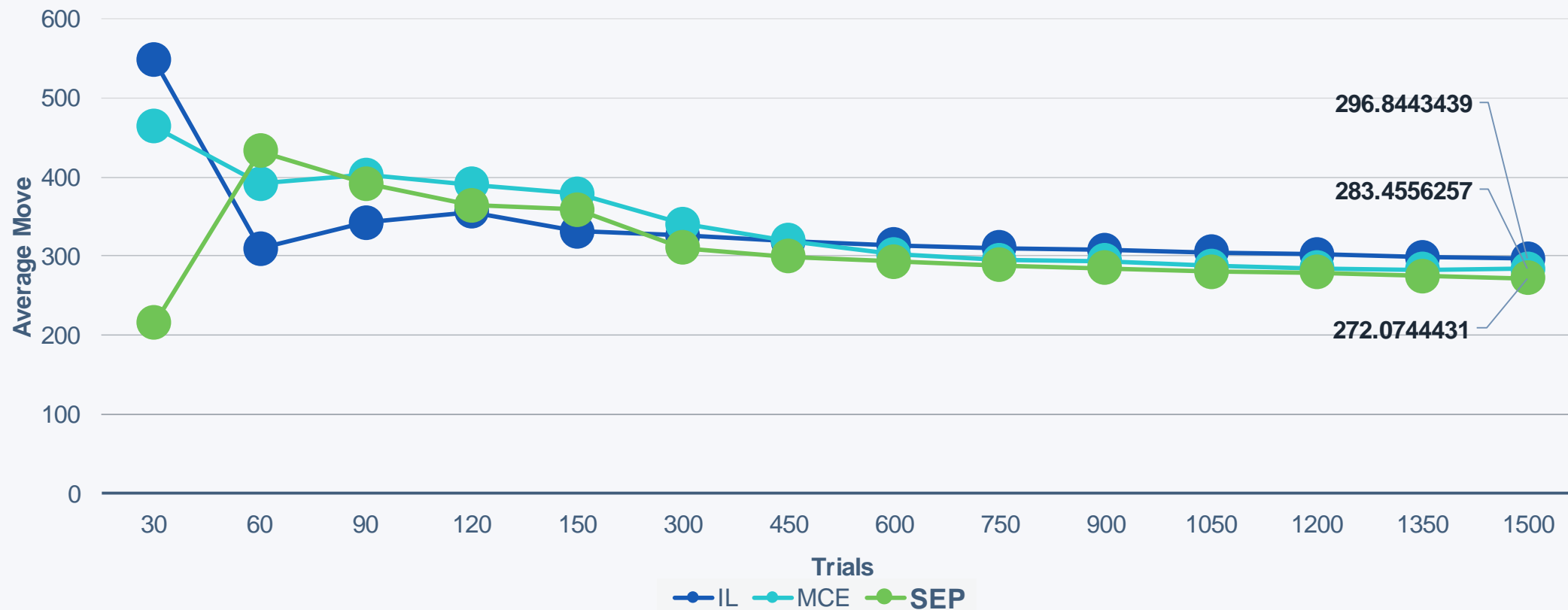
16 حالت

پاداش هر شکار 10 و جریمه 0.1 برای هر حرکت بی حاصل

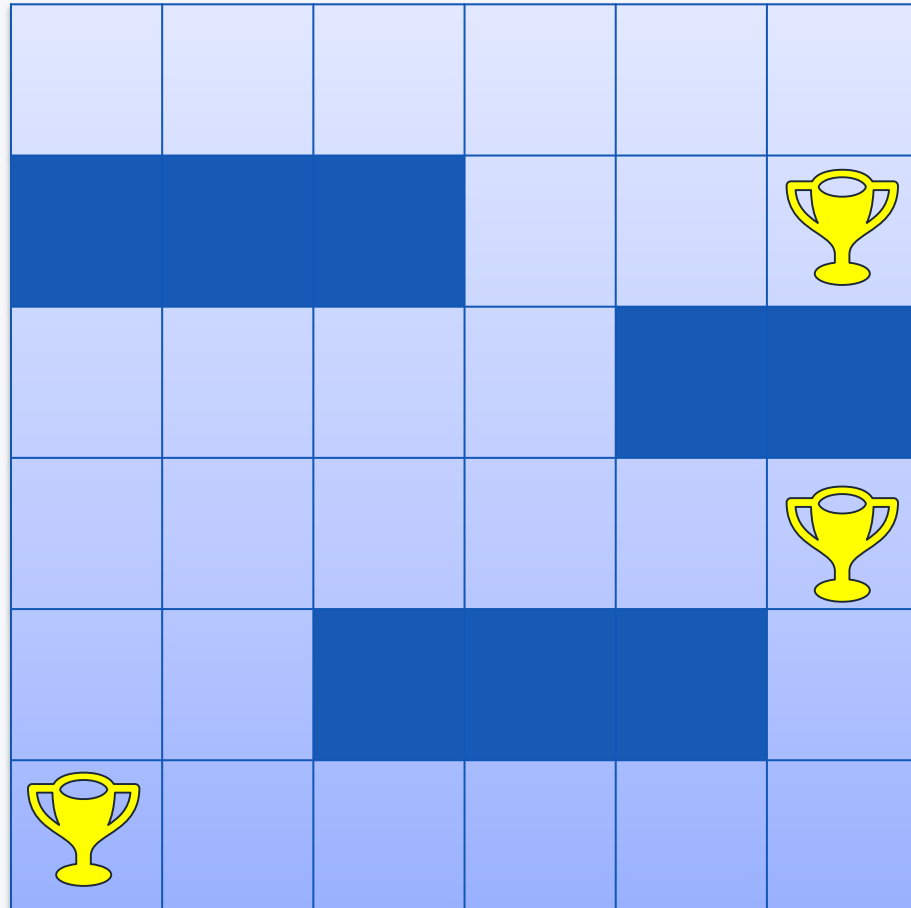
صید  صیاد 

بهره گیری از روش پیشنهادی در مقایسه با روش خبرگی چند معیاره

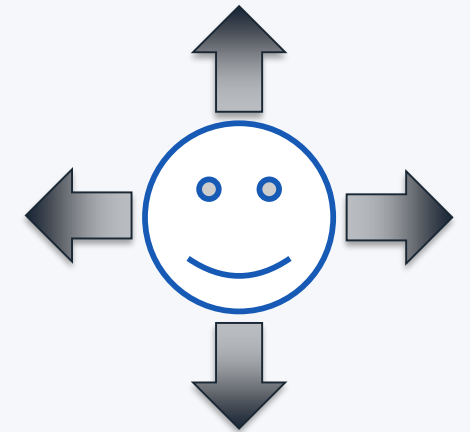
محیط صید و صیاد



پلکان مارپیچ



حالت 28



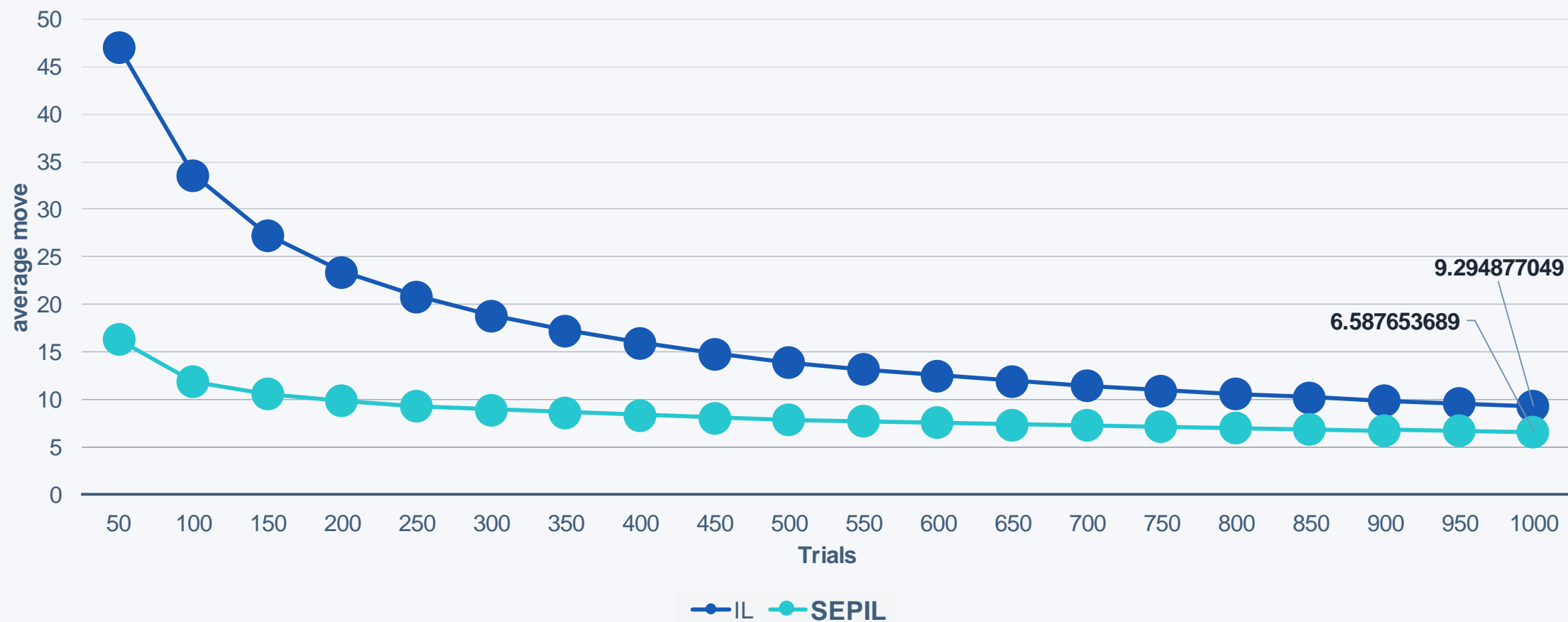
حالت 4

پاداش کشف طلا 10 و جریمه بر خورد یا موانع 1- پاداش دیگر حرکت ها بر

اساس رابطه زیر

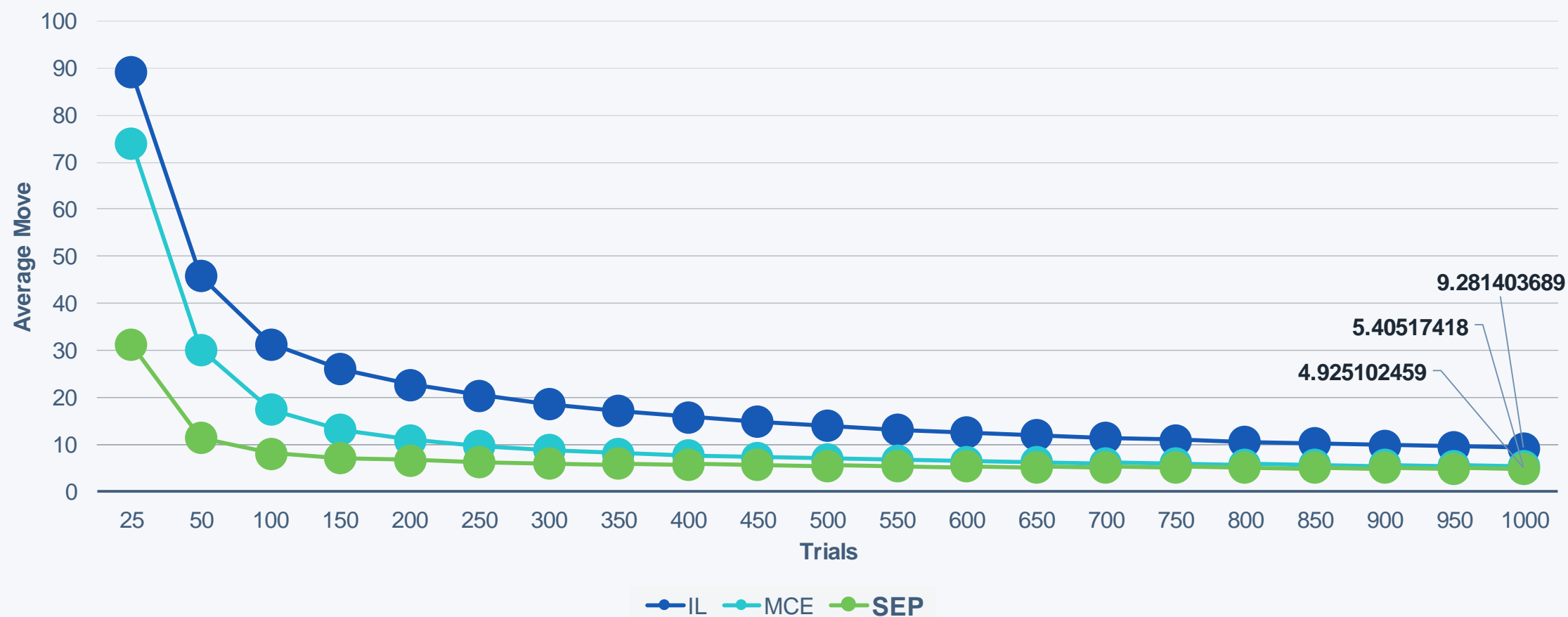
$$Reward = \frac{1}{\text{distance between the agent and the goal}}$$

بهره گیری از کوتاهترین فاصله تجربه شده در یادگیری مستقل

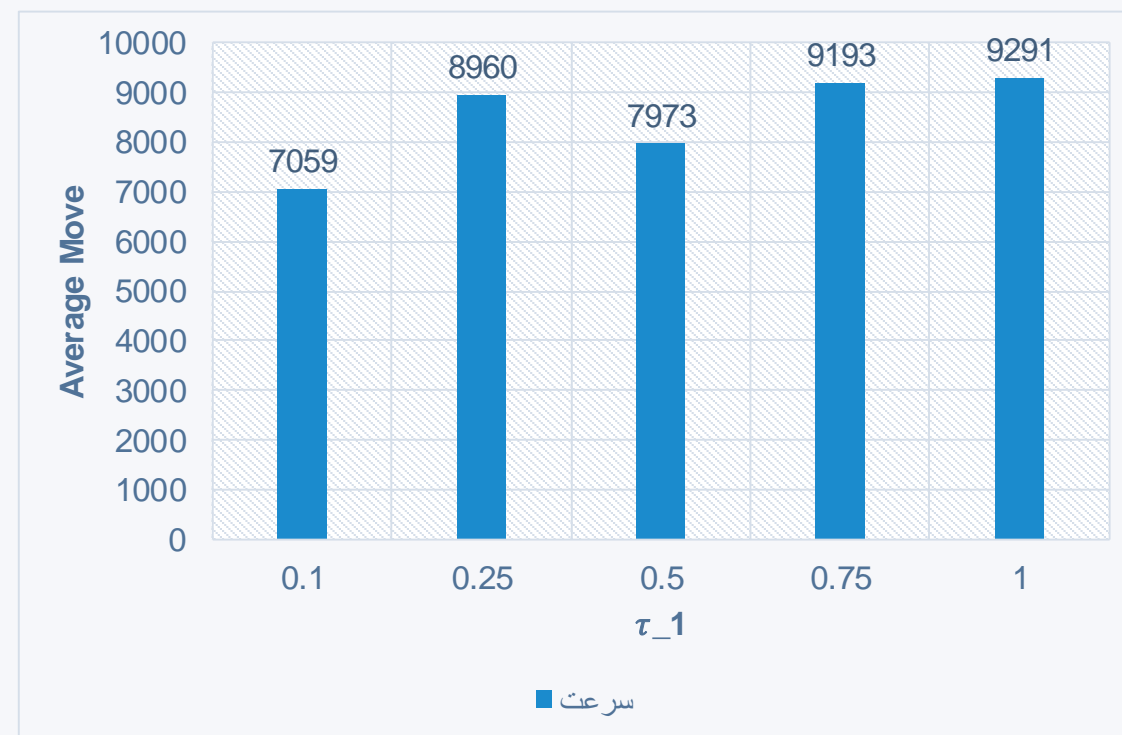
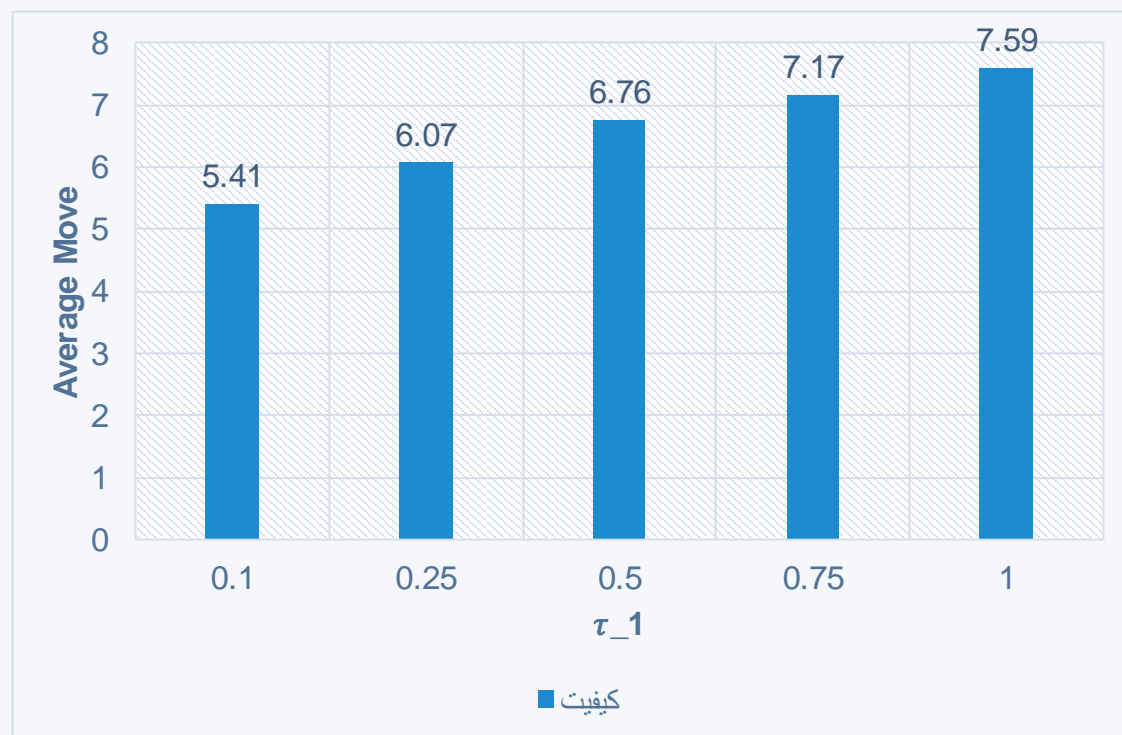


بهره گیری از روش پیشنهادی در مقایسه با روش خبرگی چند معیاره

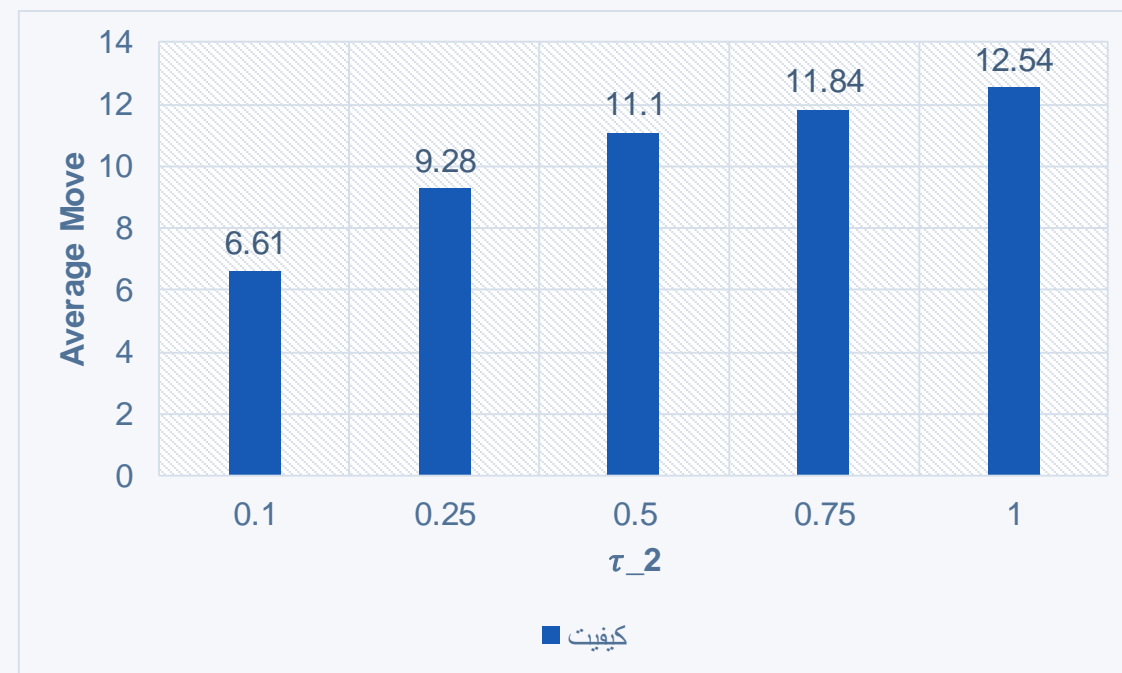
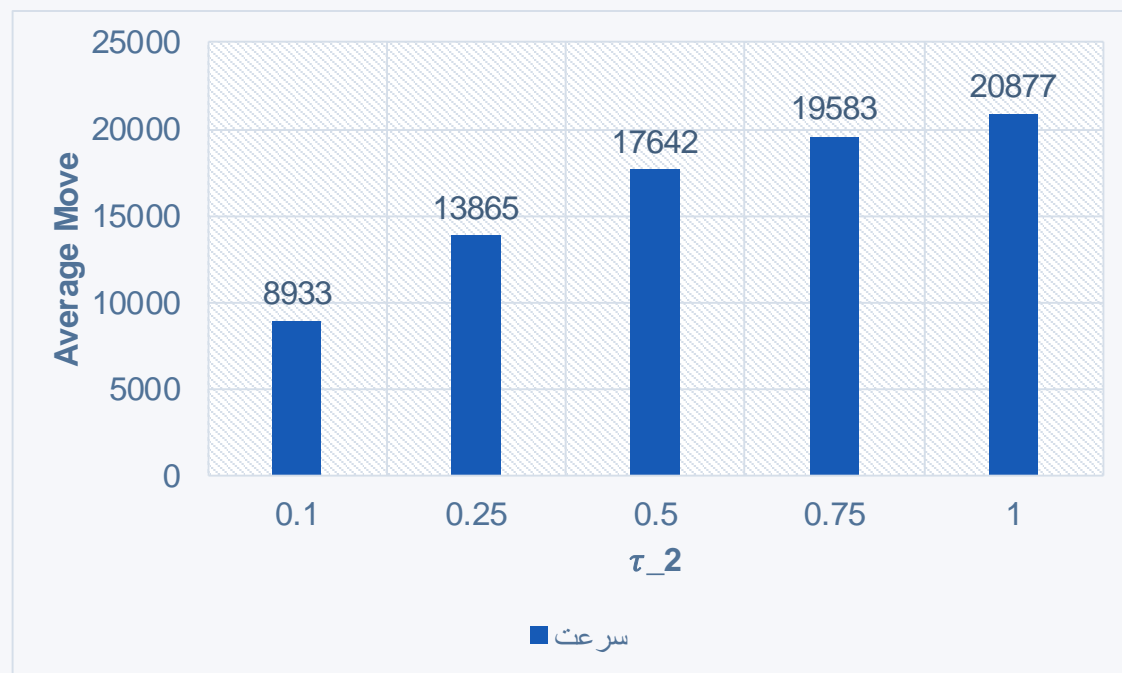
محیط پلکان مارپیچ



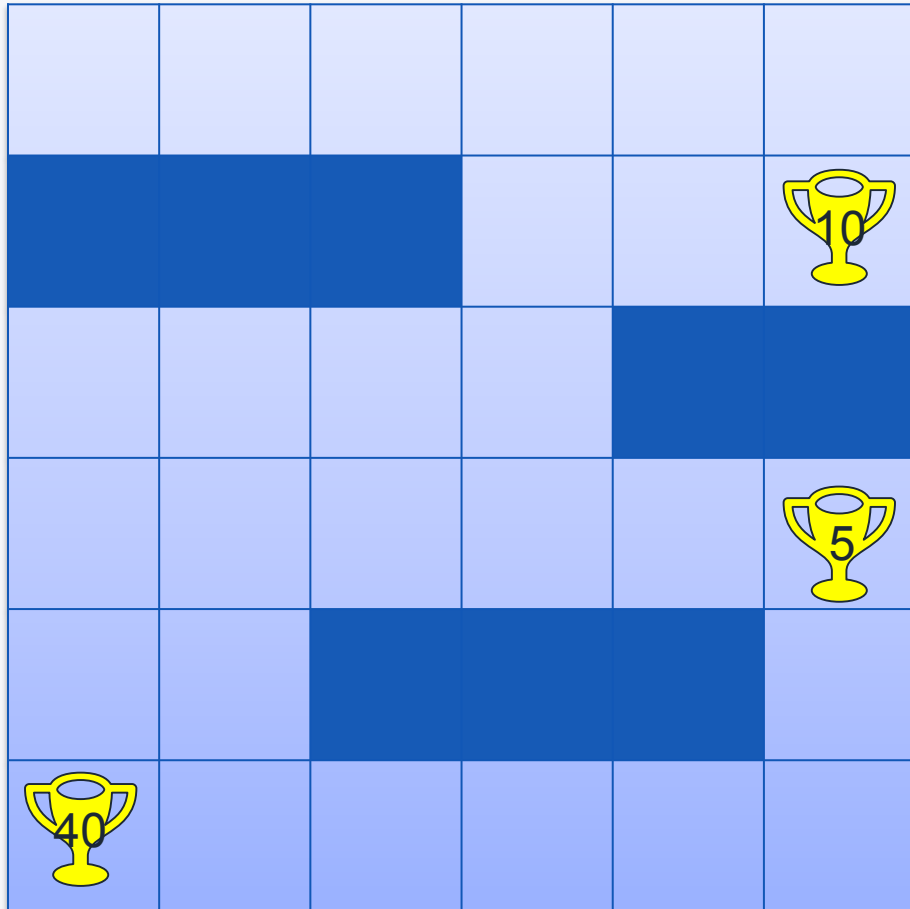
بررسی حساسیت روش پیشنهادی در برابر پارامتر τ_1



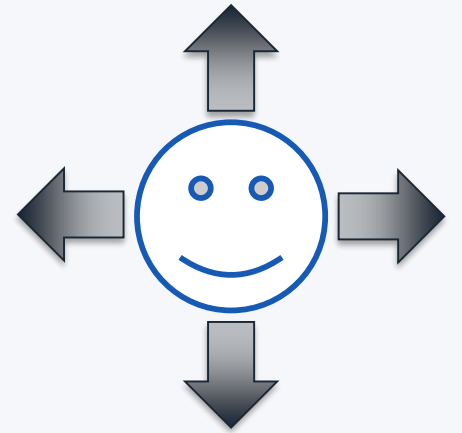
بررسی حساسیت روش پیشنهادی در برابر پارامتر τ_2



پلکان مارپیچ



حالت 28

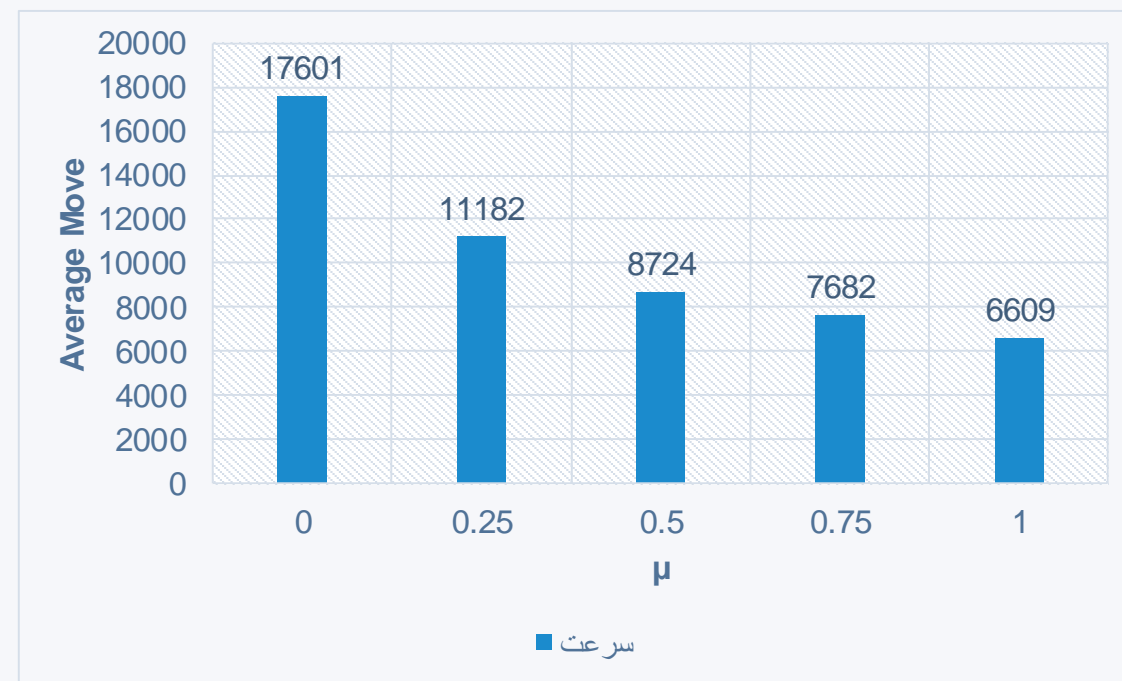
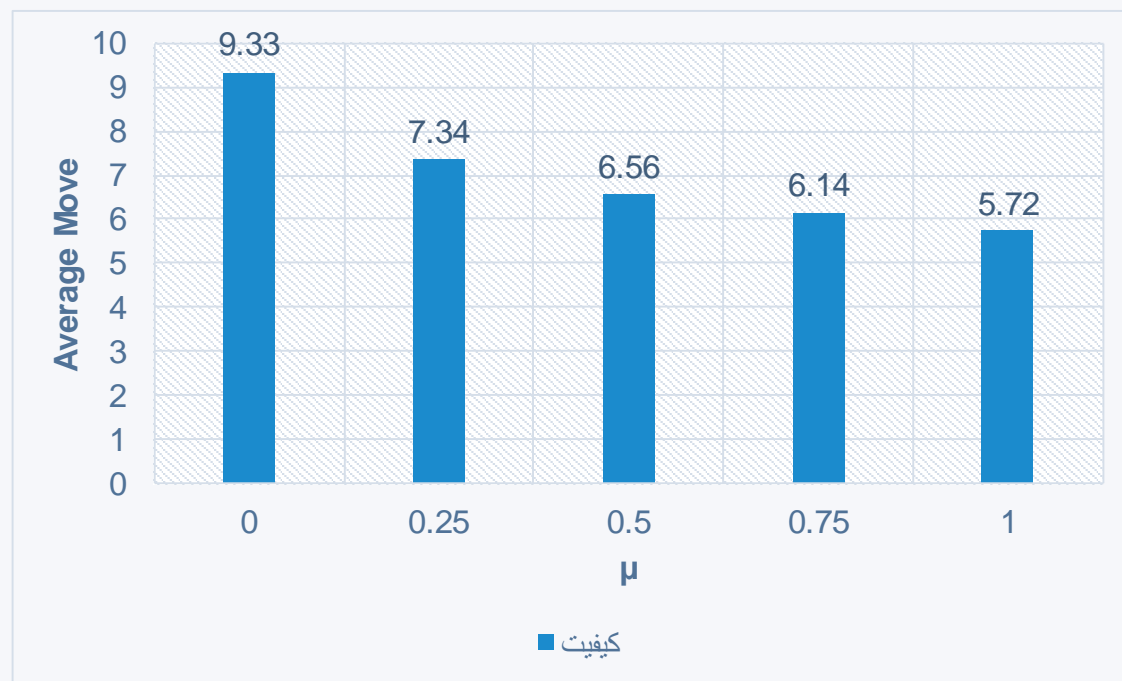


عمل 4

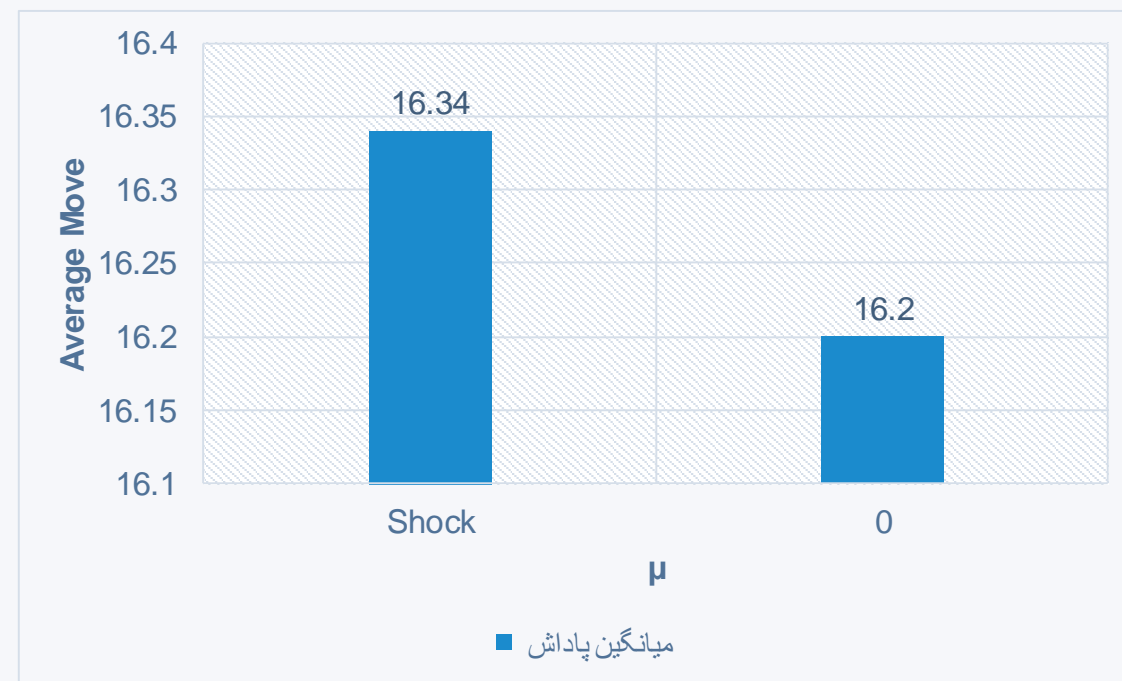
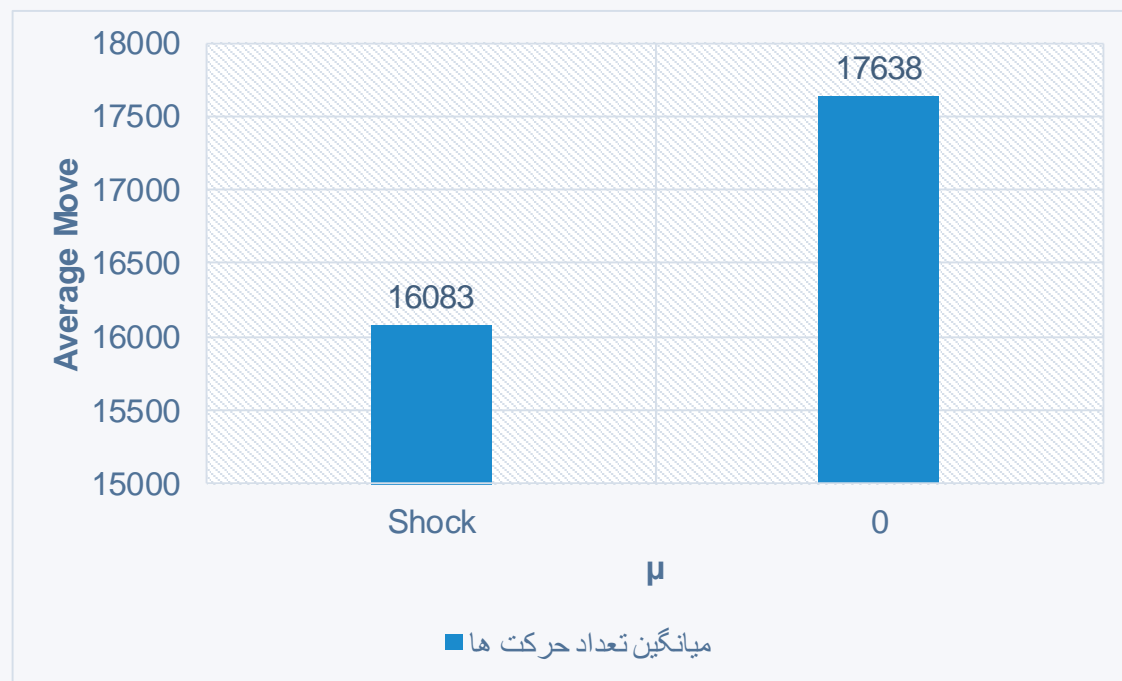
جریمه بر خورد با موانع 1- پاداش دیگر حرکت ها بر اساس رابطه زیر

$$Reward = \frac{1}{\text{distance between the agent and the goal}}$$

بررسی حساسیت روش پیشنهادی در برابر پارامتر μ



بررسی حساسیت روش پیشنهادی در برابر پارامتر μ



جمع بندی



جمع بندی

- ارائه معیار هایی جهت ارزیابی عامل های یادگیری مشارکتی
- بهره گیری از معیار های ارائه شده در راستای تسریع یادگیری مشارکتی
- بهبود یادگیری تقویتی با بهره گیری از اطلاعات تجربی

پیشنهادهات

- متغیر کردن میزان بهره گیری از حداقل فاصله تجربه شده.
- تقسیم کار مناسب.
- تهیه معیارهایی مشابه معیار کوتاهترین فاصله تجربه شده.

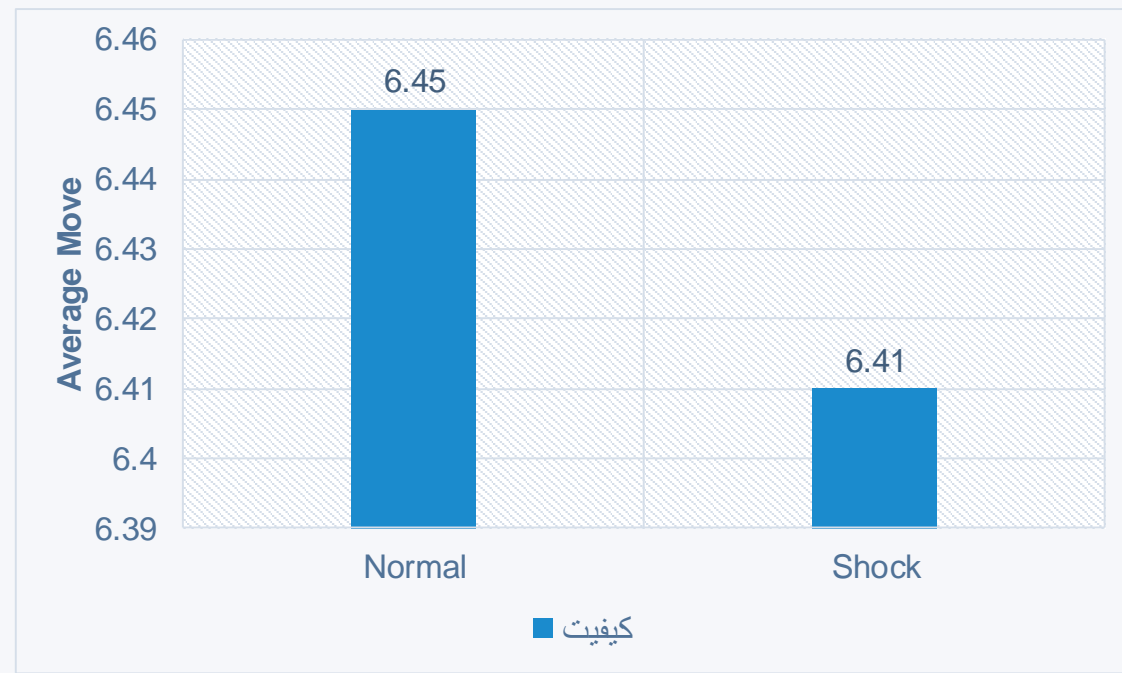
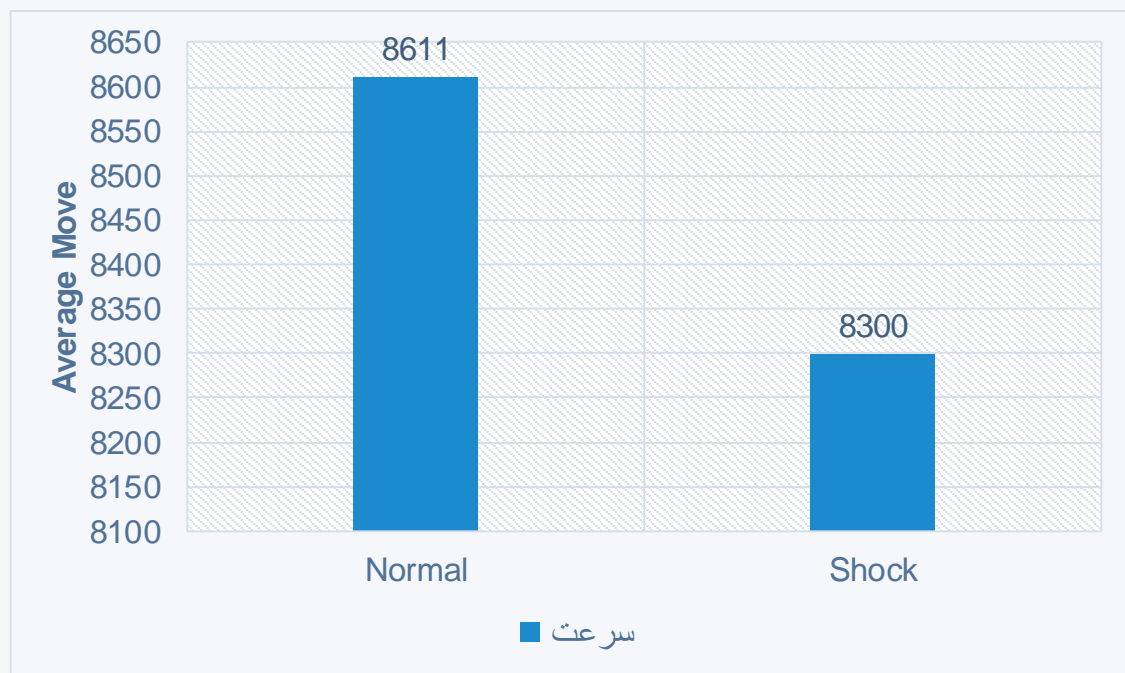
باتشکر از

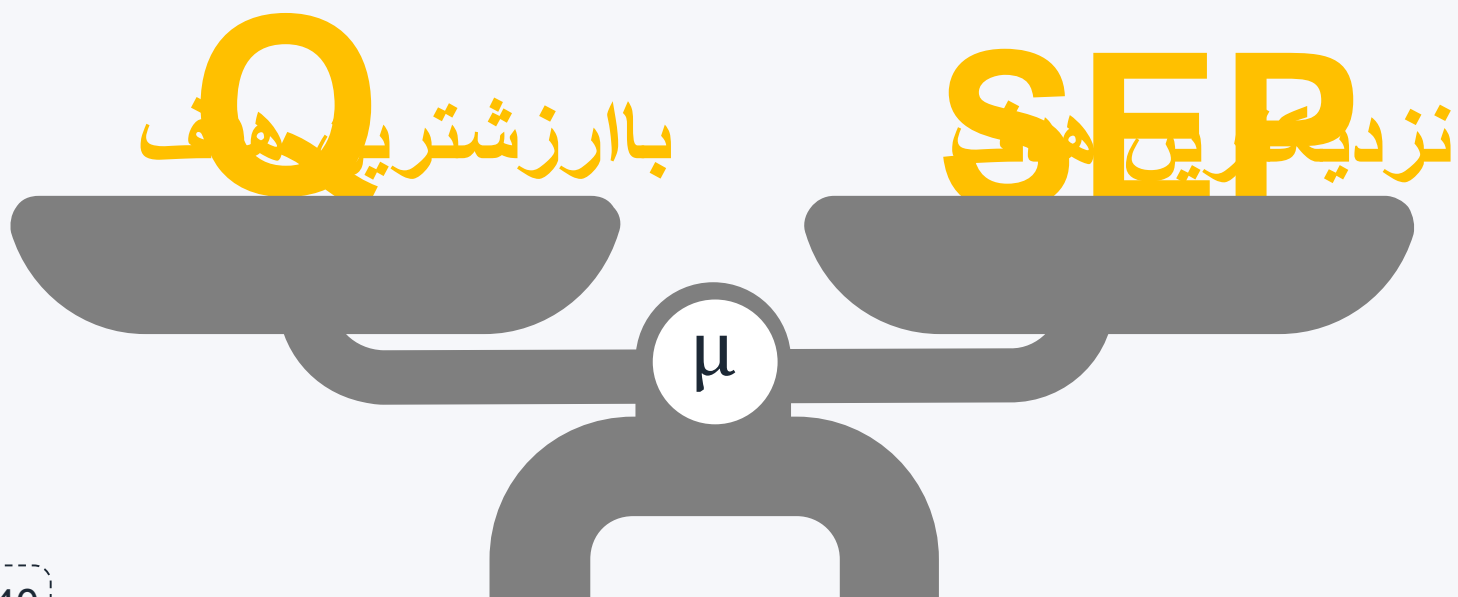
توجه شما

پرسش و پاسخ

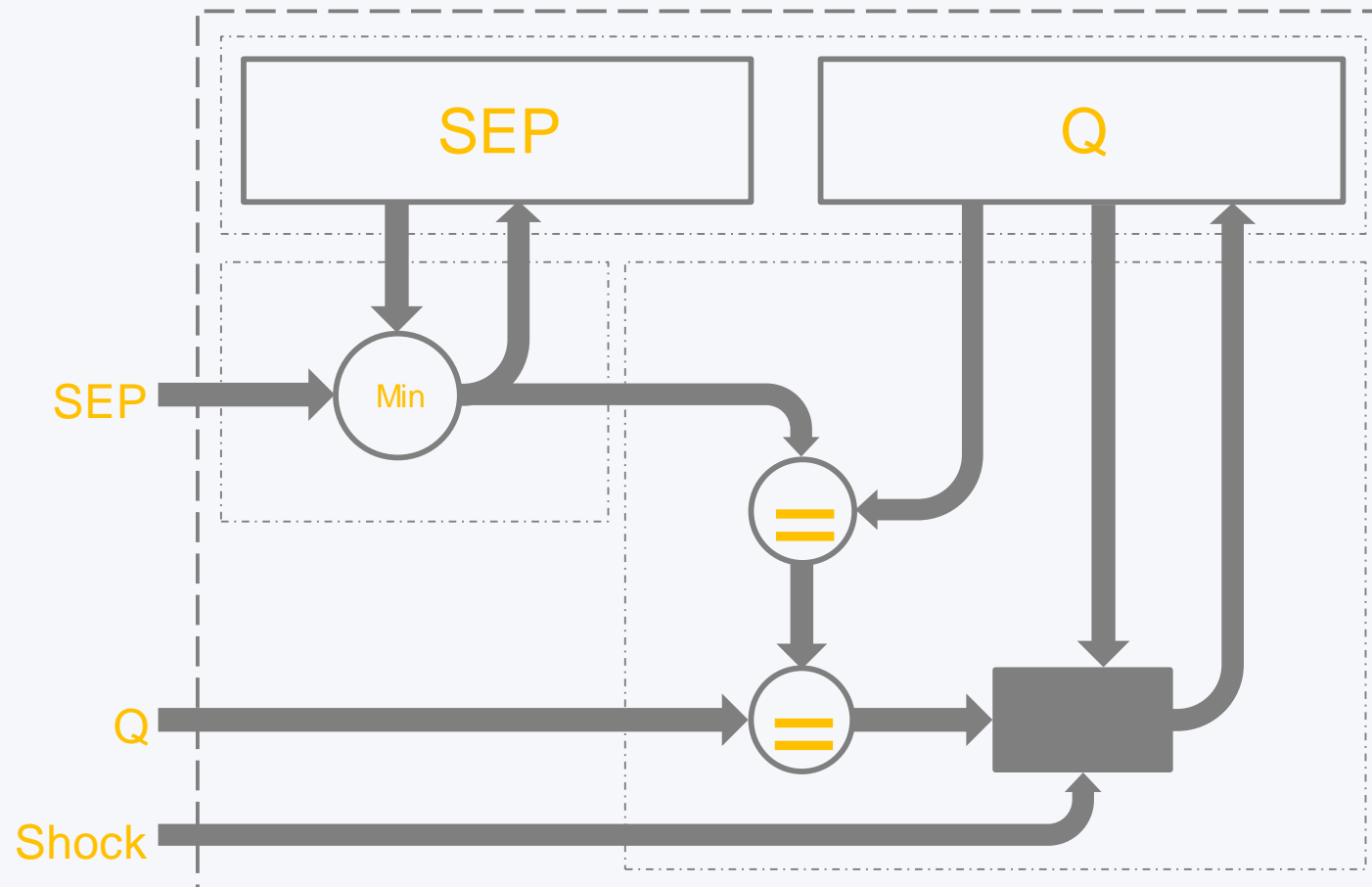


تاثیر استفاده از شوک





ترکیب داده ها



محیط های

آزمایشی

