Agent Learning from Interactions: Unifying Grounding and Large Language Models for Generalizable Agency





Agent Learning from Interactions: Unifying Grounding and Large Language Models for Generalizable Agency

Keynote Speaker: A/Prof. Wei Peng Principal Research Fellow in Al

School of Engineering RMIT University



Overview

- My Background
- Learning Pathway
- Learning from Interactions: Grounding Problem, Embodied Intelligence
- Towards more Generalizable LLM-based Agency: Agent Trajectory Learning
- Towards more Generalizable LLM-based Agency: Multi-agents Collective Decision-making





My Background

- 2025 now: A/Professor, Principal Research Fellow in Industrial AI, School of Engineering, RMIT University
- 2019 2025: Research Director and AI Expert, BPIT, Huawei Technologies
- > 2016 2019: Senior Lecturer, La Trobe University
- > 2010 2016: Senior Data Scientist, AUSTRAC, Lenovo, Telstra
- > 2008 2010: Research Fellow, RMIT University
- 2007 2008: Research Scientist, CSIRO ICT
- > 2007: PhD from the University of Sydney in Cognitive Agents

Research Interests: Large Language Models, Knowledge Representation, Cognitive Agents, and Embodied AI, along with their applications, 70+ papers, multiple patents

Career Highlights: *research, innovation, and commercialization across both industry and academia*



Bifurcated Pathway to AGI (Data-driven)

1. Human knowledge build (Cybernetics, KBS, to Data and Computing Power, LLM/scale law/transient learning on features for tasks)



We try to scale model in development to host all human knowledge and capability by feeding mega data ...

Bifurcated Pathway to AGI (Experience-driven)

2. Continual learning from experience

Rich Sutton's new path for AI : "… RL in AI, we don't have methods to learn continuously except for the linear case …" https://www.youtube.com/watch?v=NvfK1TkXmOQ



Learn from Interactions Challenges:



Generalization: beyond the env. trained

Inefficiency: learn from limited examples



Catastrophic Forgetting: learn without forgetting priori knowledge



Reward Misalignment: multi-objective, human values



Lack World Model and Abstraction: reason on concept

Learning from Interactions (Grounding Problem)

- Grounding: intrinsic process of assigning meanings to symbols/words/vectors/concepts by referencing to real world experience (objects, events).
- Symbolic grounding (Harnad, 1990): how can the semantic interpretation of a formal symbol system be made intrinsic to the system, rather than just parasitic on the meaning in our heads?
- Representation grounding (Chalmers, 1992): how can a representation in a computational system possess true meaning?
- Concept grounding (Dorffner & Prem, 1993): design a cognitive model (connectionism) only interfacing with its environment using sensor and motor signals; any concept of the system develops through self-organization based on adaptive interaction with the environment (besides given meta-level representation like innate architecture) – is grounded in Harnard's sense.



Searle's Chinese Room



Learning from Interactions (Embodied Intelligence)

Rodney Brooks' Intelligence without Representation (Brooks, 1991): no traditional representation, intelligence from sensor motor interaction with the environment, behavior-based model



Take Aways: Learn from situated sensor motor coordination to generate complex behaviors



Challenges: Limited memory, planning, reasoning capability, simplistic world model

Key Insight: ground disembodied intelligence (i.e., LLM) in interactions to develop concepts (levels of abstractions) in self-organized manner

Agent Trajectory Learning

> Agent Learning from Interaction:

AgentBank: Towards Generalized LLM Agents via Fine-Tuning on 50000+ Interaction Trajectories, *In Findings of the Association for Computational Linguistics*, pages 2124–2141, Association for Computational Linguistics (citation 13)



LLMs suffer even for SOTA commercial tools AgentBench (Liu et al., 2023)

We propose AgentBank to host 50000+ interaction trajectories

> Agent Learning from Interaction:

AgentBank: Towards Generalized LLM Agents via Fine-Tuning on 50000+ Interaction Trajectories, *In Findings of the Association for Computational Linguistics*, pages 2124–2141, Association for Computational Linguistics (citation 13)



Organize trajectories into multi-turn dialogues, mix general domain instructions and codes, utilize failure trajectories and propose the exploration-based trajectory optimization (ETO) method to learn the task-solving process, leading to significant performance gains.

> Agent Learning from Interaction:

Watch Every Step! LLM Agent Learning via Iterative Step-Level Process Refinement, In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP 2024),* pages 1556-1572, Association for Computational Linguistics (citation 20)



Agents start to learn from Interactions and explorations: from SFT on trajectories to ETO (SAMOYED)

Treat an entire trajectory as single entity during training and prioritize the final reward of a trajectory over the process, thus overlooking exploitable information throughout interaction process.

We need to consider step level optimization



> Agent Learning from Interaction:

Watch Every Step! LLM Agent Learning via Iterative Step-Level Process Refinement, In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP 2024),* pages 1556-1572, Association for Computational Linguistics (citation 20)



> Agent Learning from Interaction:

Watch Every Step! LLM Agent Learning via Iterative Step-Level Process Refinement, In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP 2024),* pages 1556-1572, Association for Computational Linguistics (citation 20)

Paradigm	Models	WebShop	InterCodeSOL	ALFWorld		Average	
g				Seen	Unseen		
Prompt-based	GPT-4 (Achiam et al., 2023)	63.2	38.5	42.9	38.1	45.7	
	GPT-3.5-Turbo (Ouyang et al., 2022)	62.4	37.8	7.9	10.5	29.7	
	Llama-2-7B (Touvron et al., 2023)	17.9	4.0	0.0	0.0	5.5	
Outcome Refinement	Llama-2-7B + SFT (Chen et al., 2023)	60.2	54.9	60.0	67.2	60.6	
	Llama-2-7B + PPO (Schulman et al., 2017)	64.2	52.4	22.1	29.1	42.0	
	Llama-2-7B + RFT (Yuan et al., 2023)	63.6	56.3	62.9	66.4	62.3	
	Llama-2-7B + ETO (Song et al., 2024)	67.4	57.2	68.6	72.4	66.4	
Drogoog Dofinament	Llama-2-7B + Step-PPO	64.0	60.2	65.7	69.4	64.8	
riocess kennement	Llama-2-7B + IPR (ours)	71.3	61.3	70.3	74.7	69.4	



Conclusion:

- Agent learns from interaction via trajectory with step awards
- Learning from failure actions
- Automated process reward acquisition
- Step level process supervision via mixture trajectory optimization
- > Enhanced performance on three benchmarks
- Generalizable on unseen hold out

Limitation:

- Overfitting with limited data (need to leverage AgentBank data)
- > MC method constrained by sample size
- Consider GPT 4 to label process supervision data

Multi-Agents Collective Decision-making

> Multi-agents Collective Decision Making (CDM):

An Electoral Approach to Diversify LLM-based Multi-Agent Collective Decision-Making, *in Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP 2024),* pages 2712–2727, Association for Computational Linguistics, <u>https://aclanthology.org/2024.emnlp-main.158/</u>



Table 4: Full list of 52 surveyed LLM-based multi-agent collaboration works.



52 multi-agent collaboration frameworks: lack of diversity in Collective Decision-making (CDM)

CDM Method	Major -ity	Mono -tonic	Consis -tency	ПА	Cond -orcet	Ballot type	
Dictatorial (Blind)	×	1	1	1	×	Ranking	
Range Voting	×	1	1	1	×	Scores	
Plurality	1	1	1	X	×	Single*	
Borda Count	X	1	1	X	×	Ranking	
IRV	1	×	X	X	×	Ranking	
Ranked Pairs	1	1	×	X	1	Ranking	

Table 1: Criteria compliance of some typical CDM methods. *Range Voting* can be viewed as a special *utilitarian* method. **IIA** denotes *Independence from Irrelevant Alternatives.* *Single ballots can be derived from ranking ones. Find some examples in Appendix D.

Kenneth Arrow's Social Choice Theory

Diversifying CDM in LLM MAS



Key Findings

Base Model	Rand.	Score	Dictatorial-based			Ordinal Ranking					
MMLU	Rand.	Range Voting	Blind Dicta.	Informed Dicta.	Mis-Informed Dicta.	Plurality	Bucklin	Borda Count	IRV	Minimax	Ranked Pairs
mistral-7b	24.8	51.8 (-4.6)	56.4	55.9 (-0.5)	36.1 (-20.3)	56.8 (+0.4)	57.1 (+0.7)	56.9 (+0.5)	56.9 (+0.5)	57.0 (+0.6)	57.0 (+0.6)
11ama-3-8b	25.0	37.7 (-7.3)	45.0	36.5 (-8.5)	32.2 (-12.8)	45.9 (+0.9)	46.4 (+1.4)	46.3 (+1.3)	45.7 (+0.7)	45.9 (+0.9)	46.0 (+1.0)
glm-4-9b	25.2	61.3 (-0.4)	61.7	54.3 (-7.4)	53.0 (-8.7)	64.6 (+2.9)	64.5 (+2.8)	64.1 (+2.4)	64.9 (+3.2)	64.4 (+2.7)	64.6 (+2.9)
11ama-3-70b	25.3	74.9 (+1.6)	73.3	70.1 (-3.2)	62.6 (-10.7)	73.9 (+0.6)	73.8 (+0.5)	73.7 (+0.4)	73.9 (+0.6)	73.9 (+0.6)	73.9 (+0.6)
qwen-2-72b	25.1	69.2 (-0.5)	69.7	69.7 (±0.0)	39.5 (-30.2)	70.0 (+0.3)	69.9 (+0.2)	70.0 (+0.3)	69.9 (+0.2)	69.9 (+0.2)	69.9 (+0.2)
qwen-1.5-110b	25.0	71.3 (-1.5)	72.8	73.0 (+0.2)	46.3 (-26.5)	72.9 (+0.1)	72.9 (+0.1)	72.7 (-0.1)	72.9 (+0.1)	72.9 (+0.1)	72.9 (+0.1)
gpt-3.5	24.9	63.0 (+2.2)	60.8	64.7 (+3.9)	36.9 (-23.9)	65.9 (+5.1)	65.5 (+4.7)	65.6 (+4.8)	65.6 (+4.8)	65.6 (+4.8)	65.6 (+4.8)
gpt-4	25.0	80.7 (+5.1)	75.6	82.1 (+6.5)	70.9 (-4.7)	82.5 (+6.9)	81.9 (+6.3)	81.9 (+6.3)	81.9 (+6.3)	81.9 (+6.3)	81.9 (+6.3)
MMLU-Pro											
mistral-7b	9.6	20.9 (-9.0)	29.9	27.7 (-2.2)	15.6 (-14.3)	31.7 (+1.8)	30.7 (+0.8)	31.4 (+1.5)	31.2 (+1.3)	31.7 (+1.8)	31.7 (+1.8)
11ama-3-8b	9.7	18.9 (-2.4)*	21.3	23.8 (+2.5)	19.3 (-2.0)	22.2 (+0.9)	23.8 (+2.5)	24.5 (+3.2)	22.6 (+1.3)	23.0 (+1.7)	23.4 (+2.1)
glm-4-9b	9.6	26.2 (-5.7)*	31.9	28.2 (-3.7)	23.9 (-8.0)	36.4 (+4.5)	35.9 (+4.0)	34.8 (+2.9)	36.7 (+4.8)	35.6 (+3.7)	36.2 (+4.3)
11ama-3-70b	10.3	46.7 (+3.5)	43.2	44.6 (+1.4)	24.6 (-18.6)	42.8 (-0.4)	43.5 (+0.3)	43.6 (+0.4)	43.0 (-0.2)	43.2 (±0.0)	43.5 (+0.3)
qwen-2-72b	10.4	35.1 (-1.7)	36.8	37.4 (+0.6)	19.5 (-17.3)	37.2 (+0.4)	36.7 (-0.1)	36.7 (-0.1)	37.2 (+0.4)	37.3 (+0.5)	37.2 (+0.4)
qwen-1.5-110b	10.1	45.7 (+0.9)	44.8	42.8 (-2.0)	16.6 (-28.2)	44.7 (-0.4)	44.9 (+0.1)	44.6 (-0.2)	45.1 (+0.3)	45.0 (+0.2)	44.8 (±0.0
gpt-3.5	9.9	28.5 (+2.6)	25.9	27.1 (+1.2)	13.0 (-12.9)	26.5 (+0.6)	27.0 (+1.1)	28.5 (+2.6)	26.5 (+0.6)	26.7 (+0.8)	27.2 (+1.3)
gpt-4	9.9	46.4 (-0.5)	46.9	46.9 (±0.0)	34.6 (-12.3)	47.3 (+0.4)	47.5 (+0.6)	47.7 (+0.8)	47.5 (+0.6)	47.8 (+0.9)	47.7 (+0.8)
ARC-Challenge											
mistral-7b	24.9	53.1 (-17.9)	71.0	70.3 (-0.7)	47.7 (-23.3)	71.7 (+0.7)	71.7 (+0.7)	71.6 (+0.6)	71.7 (+0.7)	71.7 (+0.7)	71.6 (+0.6)
11ama-3-8b	25.2	44.4 (-21.8)	66.2	52.8 (-13.4)	41.1 (-25.1)	71.3 (+5.1)	70.0 (+3.8)	70.0 (+3.8)	71.6 (+5.4)	71.3 (+5.1)	71.3 (+5.1)
glm-4-9b	24.8	69.9 (-9.7)*	79.3	80.1 (+0.8)	65.1 (-14.2)	82.7 (+3.4)	82.3 (+3.0)	82.0 (+2.7)	82.8 (+3.5)	83.0 (+3.7)	82.7 (+3.4)
11ama-3-70b	25.3	88.9 (+1.1)	87.8	87.9 (+0.1)	80.8 (-7.0)	88.5 (+0.7)	88.4 (+0.6)	88.1 (+0.3)	88.5 (+0.7)	88.4 (+0.6)	88.4 (+0.6)
gwen-2-72b	24.8	84.7 (-1.1)	85.8	86.0 (+0.2)	36.7 (-49.1)	86.3 (+0.5)	86.2 (+1.3)	85.8 (±0.0)	86.3 (+0.5)	86.3 (+0.5)	86.2 (+0.4)
gwen-1.5-110b	24.7	87.0 (-0.7)	87.7	88.3 (+0.6)	53.4 (-34.3)	88.1 (+0.4)	88.1 (+0.4)	88.0 (+0.3)	88.1 (+0.4)	88.1 (+0.4)	88.1 (+0.4)
gpt-3.5	25.2	78.1 (+1.2)	76.9	77.0 (+0.1)	29.9 (-47.0)	78.2 (+1.3)	77.9 (+1.0)	78.2 (+1.3)	78.1 (+1.2)	77.9 (+1.0)	77.9 (+1.0)
gpt-4	25.0	92.9 (+0.4)	92.5	92.8 (+0.3)	87.3 (-5.2)	92.9 (+0.4)	92.7 (+0.2)	92.8 (+0.3)	92.8 (+0.3)	92.8 (+0.3)	92.9 (+0.4)

Table 2: Overall accuracy results on MMLU, MMLU-Pro and ARC-Challenge benchmarks. 'Rand.' and 'Dicta.' denote 'random' and 'dictatorial', respectively. The numbers in parentheses are relative to the *blind dictatorial* baselines. Performance gains are marked in red, and loss in blue. Notable cases are marked in **bold**. *Results marked with asterisk are calculated utilizing partial profiles (see Appendix C).

Robustness against Unreliable Agents



Figure 4: Accuracy impact of increasing number of unreliable agents built on gpt-3.5 and gpt-4.

Limitation:

- MCQA is a limited scenario of CDM (preference over correctness)
- Limited CDM methods in GEDI, no compound of multiple voting strategies
- Voting Tax: computation cost of inter-agent communication is high

Major Literature

- > John, Searle. 1980. Minds, Brains and Programs. Behavioral and Brain Sciences, 3: 417-457.
- Stevan, Harnad. 1990. The Symbol Grounding Problem. Physica D: Nonlinear Phenomena, 42(1-3): 335–346.
- Rodney A. Brooks. 1991. Intelligence without Representation. Artificial Intelligence, 47:139-159.
- David J. Chalmers. 1992. Subsymbolic Computation and the Chinese Room. In Dinsmore, J. (ed.) The Symbolic and Connectionist Paradigms: Closing the Gap. Hillsdale, NJ: Lawrence Erlbaum.
- Georg Dorffner and Erich Prem. 1993. Connectionism, Symbol Grounding, and Autonomous Agents. In Proceedings of the Fifteenth Annual Meeting of the Cognitive Science Society, pages 144-148. Hillsdale, NJ: Lawrence Erlbaum.
- Richard Sutton, 2024, New Path for AI: Approximately Correct Podcast, URL: <u>https://www.youtube.com/watch?v=NvfK1TkXmOQ</u>
- Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating Ilms as agents. ArXiv preprint, abs/2308.03688, 2023. URL https://arxiv.org/abs/2308.03688.
- Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. 2024. AgentTuning: Enabling Generalized Agent Abilities for LLMs. In Findings of the Association for Computational Linguistics: ACL 2024, pages 3053–3077, Bangkok, Thailand. Association for Computational Linguistics.
- Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. 2024. Trial and Error: Exploration-Based Trajectory Optimization of LLM Agents. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 7584–7600, Bangkok, Thailand. Association for Computational Linguistics.
- Yifan Song, Weimin Xiong, Xiutian Zhao, Dawei Zhu, Wenhao Wu, Ke Wang, Cheng Li, Wei Peng, and Sujian Li. 2024. AgentBank: Towards Generalized LLM Agents via Fine-Tuning on 50000+ Interaction Trajectories. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 2124–2141, Miami, Florida, USA. Association for Computational Linguistics.
- Weimin Xiong, Yifan Song, Xiutian Zhao, Wenhao Wu, Xun Wang, Ke Wang, Cheng Li, Wei Peng, and Sujian Li. 2024. Watch Every Step! LLM Agent Learning via Iterative Step-level Process Refinement. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1556–1572, Miami, Florida, USA. Association for Computational Linguistics.
- Xiutian Zhao, Ke Wang, and Wei Peng. 2024. An Electoral Approach to Diversify LLM-based Multi-Agent Collective Decision-Making. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 2712–2727, Miami, Florida, USA. Association for Computational Linguistics.



Thanks Questions & discussion welcome.



Acknowledgement of Country

RMIT University acknowledges the people of the Woi wurrung and Boon wurrung language groups of the eastern Kulin Nation on whose unceded lands we conduct the business of the University.

RMIT University respectfully acknowledges their Ancestors and Elders, past and present.

RMIT also acknowledges the Traditional Custodians and their Ancestors of the lands and waters across Australia where we conduct our business.

Artwork 'Sentient' by Hollie Johnson

Hollie is a Gunaikurnai and Monero Ngarigo woman from Gippsland who graduated from RMIT with a BA in Photography in 2016.