
Idea: Online Opponent Modeling with Foundation Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Opponent modeling (OM) is the ability to use prior knowledge and observations in
2 order to predict the behavior of an opponent. On the other hand, there has been
3 tremendous research at the intersection of foundation models (FM) and decision-
4 making which holds tremendous promise for creating powerful new systems that
5 can interact effectively across a diverse range of applications. This paper examines
6 the integration of foundation models with opponent modeling and tackles one
7 of the open problems in FMs for decision-making (i) leveraging and collecting
8 decision-making datasets D_{RL} ; specifically datasets for the opponent modeling
9 systems in the large-scale human demonstration, which is hard to scale., and (ii)
10 proposing a new framework for opponent modeling: Using FMs as a guiding tool
11 that enhances the agent capabilities in prediction. The goal is to train a policy
12 from a given environment without reward signals. I propose using foundation
13 models (FMs), i.e., large language models (LLMs) and vision-language models
14 (VLMs), to achieve this goal. The LLM generates instructions that help the agent
15 to learn features of the behavior of the opponent and ultimately enables the agent
16 to exploit the opponent’s strategy in the current environment $d(s_0)$. In contrast, the
17 VLM works as a policy-guided learning. The internet-scale knowledge capacity of
18 recent FMs enables automating impractical human effort in the RL framework [1].
19 Existing works query pre-trained LLMs for tasks to learn [2], language-level plans
20 [3], and language labels [4]; or use pre-trained VLMs to obtain visual feedback
21 [5]. ELLM [6] uses LLMs to propose new tasks for agents to learn. A line of work
22 [7] specifically focuses on using FMs for the Minecraft domain, while none of
23 the works integrate pre-trained LLM and VLM for opponent modeling. Inspired
24 by [8], this work is mainly motivated by two questions: How to leverage and
25 construct datasets for decision-making D_{RL} i.e. FMs and OM? And can we teach
26 RL agents to predict opponents’ actions and strategies accurately in opponent
27 modeling environments without human supervision?

28 1 Introduction

29 In a Partially-Observable Stochastic Game (POSG) [9] for a basic formalization of the competitive
30 environment. A POSG is defined by a tuple $\langle I, S, O^i, A, T, R^i, \Omega^i \rangle$, where $I = 1, 2, \dots, N$
31 is the set of agents. S is the state space. O^i is the observation space of agent i . $A = A^1 \times A^2 \times \dots \times A^N$
32 is the joint action space. $T : S \times A \times S \rightarrow [0, 1]$ denotes the transition dynamics, which defines the
33 probability distribution on the next state given the previous state and the joint action. $R_i : S \times A \times S \rightarrow$
34 \mathbb{R} denotes the reward function of agent i . $\Omega^i : S \times A \times O^i \rightarrow [0, 1]$ denotes the agent i ’s observation
35 function, which defines the probability distribution over its possible next observation given the
36 previous state and the joint action.

37 In this study, I utilize 1 to denote the controlled agent and 1 to denote the opponent and focus on
 38 modeling one opponent. I assume that the opponent’s policy originates from a set of fixed policies
 39 $\Pi = \{\pi^{1,k}(a^1|o^1)\}_{k=1,2,\dots,k}$, which are obtained by the scripts or RL algorithms pre-training.

40 1.1 Notations

41 The first step is to generate a set of imagined task instructions that are useful for learning behaviors.
 42 Given the proposed set of N-numbers of task instructions $\{\delta^{(i)}\}_{i=1\dots N}$ and their corresponding
 43 initial states, our goal is to train a multi-task policy $\pi(a|s, \delta)$ that follows the instructions. To
 44 accomplish this, VLM can be used as a policy-guided learning, which trains a multi-task policy in
 45 the training environment using the obtained instructions. The policy is trained to follow the given
 46 instruction by maximizing the VLM “alignment score” between the current observation and the
 47 instruction as its reward. Specifically, the reward is defined by:

$$r_t = r(o_{tH:t}, \delta) = \frac{\phi_v(o_{tH:t})^T \phi_T(\delta)}{|\phi_v(o_{tH:t})| \cdot |\phi_T(\delta)|} \quad (1)$$

48 where o_t is the visual observation of time step t with $o_{tH:t}$ implying the sequence of observations
 49 with size of H, $|\cdot|$ refers to L2-norm of a vector, ϕ_T and ϕ_v are the text and video encoder of the
 50 VLM, δ is the language instruction, and H is the length of video that the VLM takes.

$$\Sigma_{i=1\dots N} E_{o_t \sim \pi, \rho, P[\Sigma_t \hat{r}(o_{tH:t}) \delta_i]} \quad (2)$$

51 .

52 2 Opponent Modeling

53 Assuming that the interaction between the controlled agent and the opponent policy $\pi^{1,k}$ generates
 54 their respective trajectories, denoted as $\tau^{1,k} = (o_0^{1,k}, a_0^{1,k}, r_0^{1,k}, o_1^{1,k}, a_1^{1,k}, r_1^{1,k}, \dots) \in \tau^{1,k}$ and
 55 $\tau^{-1,k} = (o_0^{-1,k}, a_0^{-1,k}, r_0^{-1,k}, o_1^{-1,k}, a_1^{-1,k}, r_1^{-1,k}, \dots) \in \tau^{-1,k}$. The resultant dataset is thus denoted as $D^k =$
 56 $\tau^{1,k}, \tau^{-1,k}$. Within the context of offline learning, we presume the availability of the dataset $D^{off} =$
 57 $D^k_{k=1,2,\dots,k}$. Specifically, $\tau^{1,k}$ is acquired through interactions with $\pi^{-1,k}$ while employing its
 58 approximate best response policy $\pi^{1,k,*}$, usually with certain noise.

59 The objective of OM is to use D^{off} to pre-train an opponent-aware adaptive controlled agent policy
 60 $M_\theta(a^1|o^1; D)$ and deploy M into a new environment with an unknown test opponent policy set Π^{test} ,
 61 such that the controlled agent achieves the maximum expected return (i.e., cumulative reward):

$$\max E_{\pi^{-1}} \sim \Pi^{test}, D^{off}, T, \Omega[\Sigma_{t=0} \dots \infty R_t^1 | a_t^1 \sim M_\theta \cdot \pi^{-1}] \quad (3)$$

62 D is the opponent’s information data, sampled from Doff during offline pre-training and must be
 63 collected during deployment.

64 3 How to Leverage or Collect Datasets

65 One key challenge in applying foundation models to decision-making lies in the dataset gap: the
 66 broad datasets from vision and language D and the task-specific interactive datasets D_{RL} can be
 67 of distinct modalities and structures. For instance, when D consists of videos, it generally does not
 68 contain explicit action labels indicating the cause-effect relationship between different frames, nor
 69 does it contain explicit reward labels indicating which videos are better than others, whereas actions
 70 and rewards are key components of D_{RL} . Despite this gap, broad video and text data can be made
 71 more task-specific through post-processing ($D \rightarrow D_{RL}$), leveraging hindsight relabeling of actions
 72 and rewards (e.g., using human feedback). Meanwhile, decision-making datasets can be made more
 73 broad and general ($D \rightarrow D_{RL}$) by combining a wide range of tasks-specific datasets (e.g., Gato).
 74 Below we provide a list of examples of D and D_{RL} that can be used for research in foundation
 75 models for decision-making, and propose additional approaches for bridging the gap between D and
 76 D_{RL} . In the manuscript of [10] proposed bridging D and D_{RL} . To enable better datasets tailored
 77 for decision-making, one can either increase the scale of D_{RL} by large-scale logging and merging

78 task-specific sets of interactive data or by relabeling D with action and reward information. One
79 could also consider augmenting D_{RL} with metadata, such as informational and instructional texts
80 and videos.

81 **References**

- 82 [1] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli
83 Ding, Danny Driess, Avinava Dubey, Chelsea Finn, Pete Florence, Chuyuan Fu, Montse Gonzalez Arenas,
84 Keerthana Gopalakrishnan, Kehang Han, Karol Hausman, Alex Herzog, Jasmine Hsu, Brian Ichter, Alex Irpan,
85 Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Lisa Lee, Tsang-Wei Edward Lee,
86 Sergey Levine, Yao Lu, Henryk Michalewski, Igor Mordatch, Karl Pertsch, Kanishka Rao, Krista Reymann,
87 Michael Ryoo, Grecia Salazar, Pannag Sanketi, Pierre Sermanet, Jaspier Singh, Anikait Singh, Radu Soricut,
88 Huong Tran, Vincent Vanhoucke, Quan Vuong, Ayzaan Wahid, Stefan Welker, Paul Wohlhart, Jialin Wu, Fei
89 Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-2: Vision-language-action models
90 transfer web knowledge to robotic control. In arXiv preprint arXiv:2307.15818, 2023.
- 91 [2] Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek Gupta,
92 and Jacob Andreas. Guiding pretraining in reinforcement learning with large language models. International
93 Conference on Machine Learning, 2023.
- 94 [3] Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox,
95 Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large language
96 models. In 2023 IEEE International Conference on Robotics and Automation (ICRA), pages 11523–11530,
97 2023. doi: 10.1109/ICRA48891.2023.10161317.
- 98 [4] Jesse Zhang, Jiahui Zhang, Karl Pertsch, Ziyi Liu, Xiang Ren, Minsuk Chang, Shao-Hua Sun, and Joseph J
99 Lim. Bootstrap your own skills: Learning to solve new tasks with large language model guidance. In 7th Annual
100 Conference on Robot Learning, 2023. URL <https://openreview.net/forum?id=a0mFRgadGO>.
- 101 [5] Yuchen Cui, Scott Niekum, Abhinav Gupta, Vikash Kumar, and Aravind Rajeswaran. Can foundation models
102 perform zero-shot task specification for robot manipulation? Learning for Dynamics and Control Conference,
103 2022.
- 104 [6] Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek Gupta,
105 and Jacob Andreas. Guiding pretraining in reinforcement learning with large language models. International
106 Conference on Machine Learning, 2023.
- 107 [7] Yevgen Chebotar, Quan Vuong, Karol Hausman, Fei Xia, Yao Lu, Alex Irpan, Aviral Kumar, Tianhe
108 Yu, Alexander Herzog, Karl Pertsch, Keerthana Gopalakrishnan, Julian Ibarz, Ofir Nachum, Sumedh Anand
109 Sontakke, Grecia Salazar, Huong T Tran, Jodilyn Peralta, Clayton Tan, Deeksha Manjunath, Jaspier Singh,
110 Brianna Zitkovich, Tomas Jackson, Kanishka Rao, Chelsea Finn, and Sergey Levine. Q-transformer: Scalable
111 offline reinforcement learning via autoregressive q-functions. In 7th Annual Conference on Robot Learning,
112 2023. URL <https://openreview.net/forum?id=0I3su3mkuL>.
- 113 [8] Taewook Nam, Juyong Lee, Jesse Zhang, Sung Ju Hwang, Joseph J. Lim, Karl Pertsch. LiFT: Unsupervised
114 Reinforcement Learning with Foundation Models as Teachers. 2023. arXiv:2312.08958v1
- 115 [9] Yaodong Yang and Jun Wang. An overview of multi-agent reinforcement learning from game theoretical
116 perspective. 2023. arXiv preprint arXiv:2011.00583, 2020
- 117 [10] Sherry Yang, Ofir Nachum, Yilun Du, Jason Wei, Pieter Abbeel, Dale Schuurmans. Foundation Models for
118 Decision Making: Problems, Methods, and Opportunities. arXiv:2303.04129v1