
Tackling AlfWorld with Action Attention and Common Sense from Pretrained LMs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Pre-trained language models (LMs) capture strong prior knowledge about the
2 world. This common sense knowledge can be used in control tasks. However,
3 directly generating actions from LMs may result in a reasonable narrative, but
4 not executable by a low level agent. We propose to instead use the knowledge in
5 LMs to simplify the control problem, and assist the low-level actor training. We
6 implement a novel question answering framework to simplify observations and an
7 agent that handles arbitrary roll-out length and action space size based on action
8 attention. On the Alfworld benchmark for indoor instruction following, we achieve
9 a significantly higher success rate (50% over the baseline) with our novel object
10 masking - action attention method.

11 1 Introduction

12 Humans can abstractly plan for their everyday tasks without execution ; for example, given the task
13 “Make breakfast”, we can roughly plan to first pick up mug and make coffee, then pick up egg and
14 scramble it, etc. This capability, if endowed to embodied agents, can help induce generalizable
15 common-sense and reasoning. Recently, a few works Huang et al. (2022a,b); Ahn et al. (2022);
16 Yao et al. (2020) have used large language models (LLM) Bommasani et al. (2021) for abstract
17 planning for embodied or gaming agents. These works have shown incipient success in extracting
18 procedural world knowledge from LLMs in linguistic forms and matching them with executable
19 actions conditioned on the environment.

20 However, most of recent progress neglects two aspects of abstract planning that are essential at
21 execution time. First, Huang et al. (2022a,b); Ahn et al. (2022); Yao et al. (2020) only deal with
22 open-loop planning. Since closed-loop planning enables the agent to adapt/ correct its policies upon
23 observations, it provides more flexibility at execution. Second, recent works do not address the
24 intermediate planning for finding certain objects; for example, to actually execute the planned subtask
25 “Find toothbrush”, there is the concern of “where.” The problem of “where to look” is non-trivial at
26 the execution time of a mobile agent Chaplot et al. (2020); Min et al. (2021); Blukis et al. (2021),
27 since there can be many receptacles and some, such as cabinets and drawers, even occlude small
28 objects. While it is ideal that these two aspects are considered for abstract planning, addressing
29 these concerns lead to long rollouts and large action spaces (Fig 1). More specifically, the rollouts
30 accumulated in a closed-loop setting may be too long to fit into any LM, and the large number of
31 actions makes learning with behavior cloning or reinforcement learning extremely difficult; these two
32 challenges make closed-loop, intermediate planning very difficult in the textual domain.

33 In this work, we address these two major problems with (1) a novel question answering framework to
34 filter irrelevant objects (**Object Masking**) and (2) querying long/variable length of actions (**Action**
35 **Attention**). We focus on instruction following in indoor households; on the Alfworld benchmark, we
36 achieve a significantly higher success rate (absolute 50% over the baseline) with our novel object

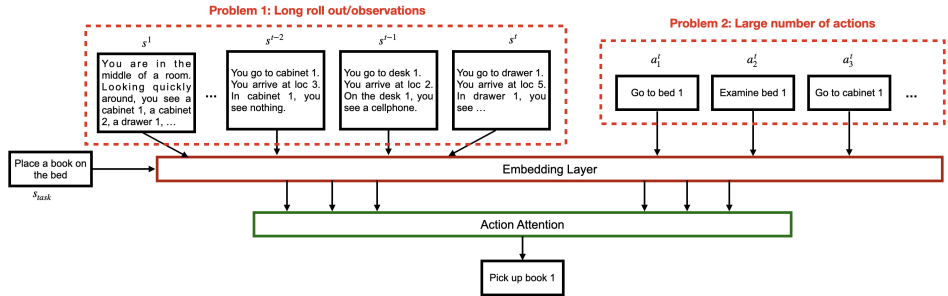


Figure 1: Overview of Action Attention method. Action Attention block is a transformer-based framework that computes a key k for each permissible action and output action scores as dot-product between key and query q from the observations. This method addresses the two problems of: (1) long roll outs and (2) large number of actions.

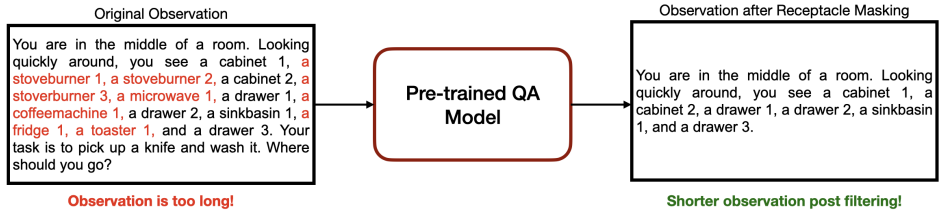


Figure 2: Overview of Receptacle Masking method. We use a pre-trained question answering model to filter irrelevant receptacles/objects in the observation of each scene. As we can see, the original observation is too long and the receptacles shown in red are not relevant for task completion. These receptacles are filtered by the QA model making the observation shorter.

37 masking - action attention method. The strong performance of our method demonstrates that large
 38 language models can be used as knowledge bases to query common sense for closed-loop intermediate
 39 planning.

40 2 Related Work

41 **Text Games** Text-based games are complex, interactive simulations where the game state and action
 42 space are in text. They are fertile ground for language-focused machine learning research. In addition
 43 to language understanding, successful play requires skills like memory and planning, exploration
 44 (trial and error), and common sense. The AlWorld simulator extends a common text-based game
 45 simulator, TextWorld Côté et al. (2018), to create text-based analogs of each ALFRED scene.

46 **LMs for Control** LMs have been used for planning high-level policies Huang et al. (2022a); Ahn
 47 et al. (2022). Huang et al. (2022a) focus on high-level sub-goals that are not executable directly in
 48 most control environment. Ahn et al. (2022) on the other hand, require few-shot demonstrations of
 49 up to 17 examples, making the length of prompt infeasible for AlWorld.

50 3 Methodology

51 Our method consists of action attention (Fig 1) and receptacle/object masking (Fig 2). The action
 52 attention module scores each permissible action with a transformer-based architecture and is trained
 53 on imitation learning on the expert. Receptacle/object masking uses a zero-shot QA model to filter
 54 out irrelevant objects in the observation.

55 **Problem Setting** We define the task description as s_{task} , the observation string at time step t as
 56 s^t , the list of permissible actions $\{a_i^t | a_i^t \text{ can be executed}\}$ as A^t . For each observation string s^t , we
 57 define the receptacles and objects within the observation as r_i^t and o_i^t respectively. We are interested
 58 in learning a policy π that outputs the optimal action among permissible actions.

59 **Action Attention** Since the number of permissible actions can vary a lot by environment, the
 60 agent needs to handle arbitrary dimensions of action space. While Shridhar et al. (2020) addresses

61 this challenge by generating actions token-by-token, such generation process leads to degenerate
62 performance even on the training set.

63 We eschew the long roll out/ large action space problems by (1) representing observations by
64 averaging over history, and (2) individually encoding actions (Fig 1). In our proposed action attention
65 framework, we first represent historical observations H^t as the average of the embeddings of all
66 individual observations through history, and H^A is a list of embeddings of all the current permissible
67 actions (Eq. 1). Then, as shown in Eq. 2, we compute the query Q with a transformer (\mathcal{M}) on the
68 task embedding (H^t), the embedding of current observation (s^t), and the list of action embeddings
69 (H^A). In Eq. 3 the key K_i is computed for each action a_i with the transformer (\mathcal{M}) on the task
70 embedding (H^t), the embedding of current observation (s^t), and the embedding of action (a_i).

71 Finally, we compute the dot-product of query and key as action scores for the policy π (Eq. 4).

$$H^t = \text{avg}_{j \in [1, t-1]} \text{Embed}(s^j), H^A = [\text{Embed}(a_1^t), \dots, \text{Embed}(a_n^t)] \quad (1)$$

$$Q = \mathcal{M}([\text{Embed}(s_{\text{task}}), H^t, \text{Embed}(s^t)], H^A) \quad (2)$$

$$K_i = \mathcal{M}[\text{Embed}(s_{\text{task}}), H^t, \text{Embed}(s^t), \text{Embed}(a_i^t)] \quad (3)$$

$$\pi = \text{softmax}(Q \cdot K) \quad (4)$$

72 **Receptacle/Object Masking** Typical Alfworld scenes can start with around 15 receptacles, each
73 containing up to 15 objects. An agent starting with no knowledge about where to look for objects
74 that are relevant to solving the task at hand can easily get stuck. We make the observation that many
75 receptacles and objects are irrelevant to specific tasks during both training and evaluation, and can be
76 easily filtered with common-sense about the tasks. For example, in Fig 2 the task is to pick up and
77 wash a knife. By removing the irrelevant receptacles like the toaster, fridge, stoveburners, we could
78 significantly shorten our observation.

79 We propose to leverage commonsense knowledge captured by large pre-trained QA models. Note
80 that we do not finetune the pre-trained QA model for our particular task but we use it in a zero-shot
81 manner. We create prompt in the format "Your task is to: <task string>. Where should you go to?"
82 for receptacles and "Your task is to: <task string>. Which objects will be relevant?" for objects. We
83 then obtain scores from the pre-trained QA model representing whether the model believe that the
84 receptacle/object is relevant, and we mask out irrelevant receptacles/objects that have scores below a
85 threshold.

86 4 Experiments and Results

87 **Hyper-parameters.** For the common-sense language model we choose Macaw-11b Tafjord and
88 Clark (2021), which is reported to have common sense QA performance on par with GPT3 Brown
89 et al. (2020) while being order of magnitudes smaller. For embedding of actions and observations, we
90 use pretrained RoBERTa-large Liu et al. (2019) with embedding dimension 1024. Our transformer
91 (\mathcal{M}) is a 12-layer transformer with 12 heads and hidden dimension 768. For receptacle/object
92 masking, we use a score threshold of 0.4 below which the objects are masked out.

93 **Baselines.** We use the BUTLER::BRAIN (**BUTLER+CG**) agent presented in Shridhar et al. (2020),
94 which consists of an encoder, an aggregator, and a decoder. At each time step t , the encoder takes
95 initial observation s^0 , current observation s^t , and task string s_{task} and generates representation r^t .
96 The recurrent aggregator combines r^t with last recurrent state h^{t-1} to produce h^t , which is then
97 decoded into a string a^t representing action. In addition, the BUTLER agent uses beam search to
98 get out of stucked conditions in the event of failed action. Our second baseline (**BUTLER+AC**) is
99 an implementation by Shridhar et al. (2020) to allow BUTLER to directly choose from admissible
100 commands. Both BUTLER agents are trained with an online imitation learning curriculum, DAgger
101 Ross et al. (2011), assisted by a rule-based expert.

102 4.1 Results

103 The results of both DAgger and Behavior Cloning are shown in Table 1. We observe that both the
104 baselines and our models benefit greatly from DAgger training. However, DAgger assumes an expert

Model	Dagger		Behavior Cloning		
	seen	unseen	train	seen	unseen
BUTLER + CG Shridhar et al. (2020)	40	35	9	10	9
BUTLER + AC* Shridhar et al. (2020)	61.7	16.89	-	-	-
Action Attention	90.41 \pm 0.02	33.42 \pm 0.05	30	25	9
Action Attention + Masking	90.53 \pm 0.02	34.92 \pm 0.03	30	25	11

Table 1: Average completion rate with DAgger and Behavior Cloning. * Shridhar et al. (2020) did not provide evaluation for BUTLER+AC, so we report the performance from our own experiment.

105 that is well-defined at all observation spaces, which is infeasible in most practical scenarios. We also
 106 observe that training is 100x slower with DAgger compared to behavior cloning (3 weeks for DAgger
 107 v.s. 6 hours for Behavior Cloning).

108 In the DAgger training scenario, our action attention agent greatly exceeds baseline performance in
 109 **seen** evaluation (we observe a 50.41% absolute improvement), and receptacle/object masking further
 110 improves the performance on unseen evaluation.

111 In the behavior cloning scenario, where there is not enough training data, we observe that Recepta-
 112 cle/Object Masking is more effective in the behavior cloning setting (we observe a 22.2% relative
 113 improvement).

114 **Quality of Pre-trained QA for receptacle/object masking** We evaluate the zero-shot recepta-
 115 cle/object masking performance of Macaw on the three splits of AlfWorld. In Fig 3, we plot the AUC
 116 curve of the relevance-score that the model assigns to the objects v.s. objects that the rule-based
 117 expert interacted with when completing each task. In practice decision threshold of 0.4 retains around
 118 80% relevant objects, 70% relevant receptacles and reduces the length of observations by 50% on
 119 average. In addition, the zero-shot QA model demonstrates consistent masking performance on all
 three splits of the environment, even on the unseen split.

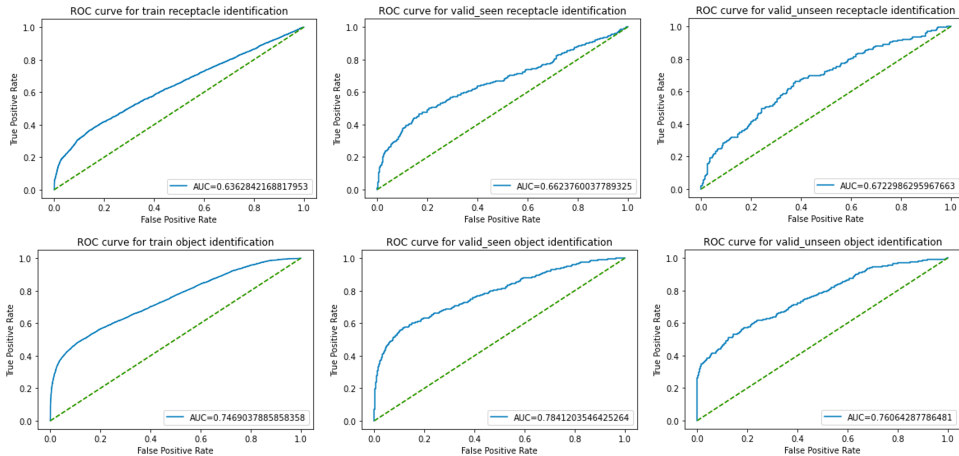


Figure 3: Plot of AUC scores of zero-shot relevance identification across all tasks in the Alfworld-Thor environment, with the Macaw-11b model. The ground truth is obtained as receptacles/objects accessed by the rule-based expert. **Top:** Receptacle relevance identification. **Bottom:** Object relevance identification. The QA model achieves average AUC-ROC score of 65 for receptacles, and 76 on objects.

120

121 5 Conclusion

122 In this work, we present (1) a novel question answering framework to simplify observation and (2) an
 123 action attention framework to handle large and variable size action space. Future works can focus on
 124 adding ways from which LMs can assist learning of the policy, such as providing high-level plans.

125 **References**

- 126 Ahn, M., Brohan, A., Brown, N., Chebotar, Y., Cortes, O., David, B., Finn, C., Fu, C., Gopalakrishnan,
127 K., Hausman, K., Herzog, A., Ho, D., Hsu, J., Ibarz, J., Ichter, B., Irpan, A., Jang, E., Ruano, R. J.,
128 Jeffrey, K., Jesmonth, S., Joshi, N. J., Julian, R., Kalashnikov, D., Kuang, Y., Lee, K.-H., Levine,
129 S., Lu, Y., Luu, L., Parada, C., Pastor, P., Quiambao, J., Rao, K., Rettinghouse, J., Reyes, D.,
130 Sermanet, P., Sievers, N., Tan, C., Toshev, A., Vanhoucke, V., Xia, F., Xiao, T., Xu, P., Xu, S., Yan,
131 M., and Zeng, A. (2022). Do as i can, not as i say: Grounding language in robotic affordances.
- 132 Blukis, V., Paxton, C., Fox, D., Garg, A., and Artzi, Y. (2021). A persistent spatial semantic
133 representation for high-level natural language instruction execution.
- 134 Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S.,
135 Bohg, J., Bosselut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji,
136 N., Chen, A., Creel, K., Davis, J. Q., Demszky, D., Donahue, C., Doumbouya, M., Durmus, E.,
137 Ermon, S., Etchemendy, J., Ethayarajh, K., Fei-Fei, L., Finn, C., Gale, T., Gillespie, L., Goel, K.,
138 Goodman, N., Grossman, S., Guha, N., Hashimoto, T., Henderson, P., Hewitt, J., Ho, D. E., Hong,
139 J., Hsu, K., Huang, J., Icard, T., Jain, S., Jurafsky, D., Kalluri, P., Karamcheti, S., Keeling, G.,
140 Khani, F., Khattab, O., Koh, P. W., Krass, M., Krishna, R., Kuditipudi, R., Kumar, A., Ladhak,
141 F., Lee, M., Lee, T., Leskovec, J., Levent, I., Li, X. L., Li, X., Ma, T., Malik, A., Manning, C. D.,
142 Mirchandani, S., Mitchell, E., Muniyikwa, Z., Nair, S., Narayan, A., Narayanan, D., Newman, B.,
143 Nie, A., Niebles, J. C., Nilforoshan, H., Nyarko, J., Ogut, G., Orr, L., Papadimitriou, I., Park, J. S.,
144 Piech, C., Portelance, E., Potts, C., Raghunathan, A., Reich, R., Ren, H., Rong, F., Roohani, Y.,
145 Ruiz, C., Ryan, J., Ré, C., Sadigh, D., Sagawa, S., Santhanam, K., Shih, A., Srinivasan, K., Tamkin,
146 A., Taori, R., Thomas, A. W., Tramèr, F., Wang, R. E., Wang, W., Wu, B., Wu, J., Wu, Y., Xie,
147 S. M., Yasunaga, M., You, J., Zaharia, M., Zhang, M., Zhang, T., Zhang, X., Zhang, Y., Zheng, L.,
148 Zhou, K., and Liang, P. (2021). On the opportunities and risks of foundation models.
- 149 Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam,
150 P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *Advances in neural
151 information processing systems*, 33:1877–1901.
- 152 Chaplot, D. S., Gandhi, D., Gupta, A., and Salakhutdinov, R. (2020). Object goal navigation using
153 goal-oriented semantic exploration.
- 154 Côté, M.-A., Kádár, A., Yuan, X., Kybartas, B., Barnes, T., Fine, E., Moore, J., Hausknecht, M.,
155 Asri, L. E., Adada, M., et al. (2018). Textworld: A learning environment for text-based games. In
156 *Workshop on Computer Games*, pages 41–75. Springer.
- 157 Huang, W., Abbeel, P., Pathak, D., and Mordatch, I. (2022a). Language models as zero-shot planners:
158 Extracting actionable knowledge for embodied agents.
- 159 Huang, W., Xia, F., Xiao, T., Chan, H., Liang, J., Florence, P., Zeng, A., Tompson, J., Mordatch, I.,
160 Chebotar, Y., Sermanet, P., Brown, N., Jackson, T., Luu, L., Levine, S., Hausman, K., and Ichter, B.
161 (2022b). Inner monologue: Embodied reasoning through planning with language models.
- 162 Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and
163 Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach. *arXiv preprint
164 arXiv:1907.11692*.
- 165 Min, S. Y., Chaplot, D. S., Ravikumar, P., Bisk, Y., and Salakhutdinov, R. (2021). Film: Following
166 instructions in language with modular methods.
- 167 Ross, S., Gordon, G., and Bagnell, D. (2011). A reduction of imitation learning and structured
168 prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on
169 artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings.
- 170 Shridhar, M., Yuan, X., Côté, M.-A., Bisk, Y., Trischler, A., and Hausknecht, M. (2020). Alf-
171 world: Aligning text and embodied environments for interactive learning. *arXiv preprint
172 arXiv:2010.03768*.
- 173 Tafjord, O. and Clark, P. (2021). General-purpose question-answering with macaw. *arXiv preprint
174 arXiv:2109.02593*.

175 Yao, S., Rao, R., Hausknecht, M., and Narasimhan, K. (2020). Keep calm and explore: Language
176 models for action generation in text-based games.