# Pre-Trained Language Models for Interactive Decision-Making

**Shuang Li** [1]  **Xavier Puig** [1]  **Chris Paxton** [2]  **Yilun Du** [1]  **Clinton Wang** [1]  **Linxi Fan** [2]  **Tao Chen** [1]  **De-An Huang** [2]
**Ekin Akyürek** [1]  **Anima Anandkumar** [2]  **Jacob Andreas** [1]  **Igor Mordatch** [3]  **Antonio Torralba** [1]  **Yuke Zhu** [2]
\* Junior authors are ordered based on their contributions and senior authors are ordered alphabetically.

## Abstract

Language model (LM) pre-training has proven useful for a wide variety of language processing tasks, but can such pre-training be leveraged for more general machine learning problems? We investigate the effectiveness of language modeling to scaffold learning and generalization in autonomous decision-making. We describe a framework for imitation learning in which goals and observations are represented as a sequence of embeddings, and translated into actions using a policy network initialized with a pre-trained transformer LM. We demonstrate that this framework enables effective combinatorial generalization across different environments, such as VirtualHome and BabyAI. In particular, for test tasks involving novel goals or novel scenes, initializing policies with language models improves task completion rates by 43.6% in VirtualHome. We hypothesize and investigate three possible factors underlying the effectiveness of LM-based policy initialization. We find that sequential representations (vs. fixed-dimensional feature vectors) and the LM objective (not just the transformer architecture) are both important for generalization. Surprisingly, however, the format of the policy inputs encoding (e.g. as a natural language string vs. an arbitrary sequential encoding) has little influence. Together, these results suggest that language modeling induces representations that are useful for modeling not just language, but also goals and plans; these representations can aid learning and generalization even outside of language processing.

## 1. Introduction

In recent years, **language models** (LMs) trained on open-domain text corpora have come to play a central role in machine learning approaches to natural language processing tasks (Devlin et al., 2018). This includes tasks that are not purely linguistic, and additionally require nontrivial planning and reasoning capabilities: for example, vision-language navigation (Majumdar et al., 2020; Fried et al., 2018; Suglia et al., 2021), instruction following (Zhang & Chai, 2021; Hill et al., 2020), and visual question answering (Tsimpoukelli et al., 2021). Indeed, some of these tasks are so remotely connected to language modeling that it is natural to ask whether pre-trained language models can be used as a general framework for different tasks, even for tasks that involve no language at all—and if so, how these capabilities might be accessed in a model trained only to process and generate natural language strings.

In this paper, we study these questions through the lens of **embodied decision-making**, investigating the effectiveness of LM pre-training as a general framework for learning policies across a variety of environments. As shown in Figure 1 (right), we encode the inputs to a policy—including observations, goals, and history—as a sequence of embeddings. These embeddings are passed to a policy network initialized with the parameters of a pre-trained LM, which is fine-tuned to predict actions. This framework is generic, accommodating goals and environment states represented as natural language strings, image patches, or scene graphs.

We find that using pre-trained LMs as policy initializers improves in-domain performance and enables several forms of strong generalization over tasks. For i.i.d. training and evaluation tasks, we find that this approach yields 20% more successful policies than other baseline methods in the VirtualHome environment (Puig et al., 2018; 2020). For combinatorial generalization to out-of-distribution tasks, *i.e.* tasks involving new combinations of goals, states or objects, we find that LM pre-training confers even more benefits: it improves task completion rates by 43.6% for tasks involving novel goals. These results hold for a variety of environment representations: encoding states as natural language strings, when possible, improves the data-efficiency of training, but

---
[1]MIT CSAIL [2]Nvidia [3]Google Brain. Correspondence to: Shuang Li <lishuang@mit.edu>.
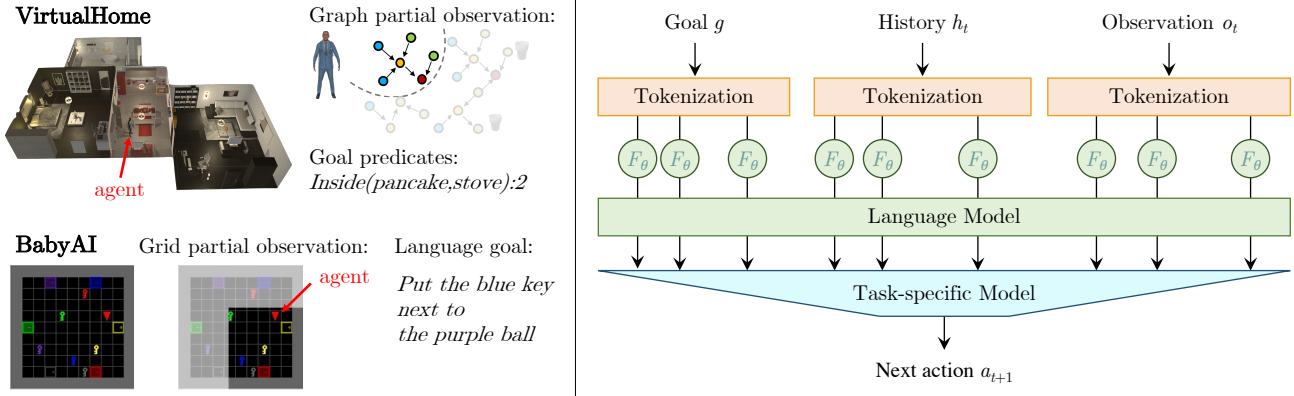
*Figure 1.* **Left: Environments.** Different environments are featured by different types of observations, goals, and histories. For example, VirtualHome has graph partial observations and its goals are specified by a set of goal predicates while BabyAI has grid partial observation and language goals. **Right: Overview of the proposed method.** We propose to use pre-trained language models as a general framework for interactive decision-making by converting all policy inputs into sequential data. Such a framework enables effective combinatorial generalization to novel tasks.

even LMs fine-tuned on random environment encodings generalize combinatorially to new goals and states when trained on large enough datasets.

We hypothesize and investigate three possible factors underlying the effectiveness of language modeling for generalization in policy learning: (1) input encoding scheme; (2) sequential input representations; and (3) parameter pre-training. We investigate (1) by encoding the environment as different types of sequences. Different input encoding schemes have only a negligible impact on model performance: the effectiveness of language modeling is not limited to utilizing natural strings, but in fact extends to arbitrary sequential encodings. We investigate (2) by encoding observations with a single vector embedding, thereby removing its sequential structure. This operation significantly hurts the model's performance on novel tasks. Finally, we investigate (3) by learning the parameters of the policy network from scratch. The success rate on novel tasks after removing the pre-trained LM weights drops by 11.2%.

To summarize, our work has four main contributions. First, we propose to use **pre-trained language models as a general framework** for interactive decision-making across a variety of environments by converting all policy inputs into sequential data. Second, the proposed method demonstrates that **language modeling improves combinatorial generalization in policy learning**: initializing a policy with a pre-trained LM substantially improves out-of-distribution performance on novel tasks. Third, we build a new dataset for evaluating embodied decision-learning that utilizes the VirtualHome environment. The new dataset presents challenging tasks in complex 3D environments, and provides several test sets for evaluating a model's ability to generalize to novel scenes and novel tasks. Finally, we perform several analyses to explain the generalization capabilities of pre-

trained LMs, finding that natural strings are not needed to benefit from LM pre-training, but the sequential input encoding and weight pre-training are important. These results may point to the effectiveness of the proposed framework with pre-trained language models as a general-purpose framework to promote structured generalization in interactive decision-making applications.

## 2. Related Work

In recent years, word and sentence representations from pre-trained LMs (Peters et al., 2018; Devlin et al., 2018; Radford et al., 2018) have become ubiquitous in natural language processing, playing a key role in state-of-the-art models for tasks ranging from machine translation (Zhu et al., 2020) to task-oriented dialog (Platanios et al., 2021). Some of the most successful applications of pre-training lie at the boundary of natural language processing and other domains, as in instruction following (Hill et al., 2020) and language-guided image retrieval (Lu et al., 2019). Building on this past work, our experiments in this paper aim to explain whether these successes result entirely from improved processing of text, or instead from domain-general representational abilities. Below, we briefly survey existing applications of pre-training that motivate the current study.

**Learning representations of language.** From nearly the earliest days of the field, natural language processing researchers observed that representations of words derived from distributional statistics in large text corpora serve as useful features for downstream tasks (Deerwester et al., 1990; Dumais, 2004). The earliest versions of these representation learning schemes focused on isolated word forms (Mikolov et al., 2013; Pennington et al., 2014). However, recent years have seen a number of techniques for training (masked or autoregressive) language models to produce

contextualized word representations (which incorporate information neighboring words in sentences and paragraphs) via a variety of masked-word prediction objectives (Devlin et al., 2018; Yang et al., 2019).

**Applications of pre-trained LMs.** In addition to producing useful representations, these language models can be fine-tuned to perform language processing tasks other than language modeling by casting those tasks as word-prediction problems. Successful uses of representations from pre-trained models include syntactic parsing (Kitaev et al., 2018) and language-to-code translation (Wang et al., 2019); successful adaptations of LM prediction heads include machine translation (Zhu et al., 2020), sentiment classification (Brown et al., 2020) and style transfer (Keskar et al., 2019). Included in these successes are a number of tasks that integrate language and other modalities, including visual question answering and image captioning (Yang et al., 2020). In models that condition on both text and image data, several previous approaches have found that image representations can be injected directly into language models' embedding layers (Tsimpoukelli et al., 2021). A large body of recent work has used LM pre-training to scaffold interactive decision-making in text-based games (Yuan et al., 2018; Narasimhan et al., 2015; Ammanabrolu & Riedl, 2018; Côté et al., 2018; Yao et al., 2020; 2021) in which both observations and actions are encoded as text.

**What do LMs encode?** The possibility that LMs might encode non-linguistic information useful for other downstream tasks is suggested by a number of recent "probing" studies aimed at understanding their predictions and the structure of their internal representations. Pre-trained LMs can answer a non-trivial fraction of queries about both factual and common-sense knowledge (Roberts et al., 2020). Their representations encode information about perceptual relations among concepts, including visual similarity among object classes (Ilharco et al., 2020) and the structure of color spaces (Abdou et al., 2021). Finally, they appear to be capable of basic simulation, modeling changes in entity states and relations described by text (Li et al., 2021).

**LM pre-training beyond language** Several recent papers consider questions related to the ones investigated here. (Brown et al., 2020) show that the GPT-3 model is capable of performing a limited set of arithmetic and string manipulation tasks; (Lu et al., 2021) show that pre-trained LMs require very little fine-tuning to *match* the performance of task-specific models on several image classification and numerical sequence processing tasks. In this paper, we focus on building a general framework for decision-making tasks using pre-trained LMs. In contrast to (Huang et al., 2022), which investigates the existing knowledge in pre-trained LMs by prompting, our work focuses on adapting/fine-tuning pre-trained LMs for decision-making tasks using imitation learning. This enables our model to incorporate

the environment context to make decisions more efficiently, as well as generalize to novel scenes and tasks.

# 3. Language Modeling and Decision-Making

We begin with a brief review of the ingredients of language modeling and policy learning tasks used in our experiments.

## 3.1. Language modeling

Our experiments in this paper focus on **autoregressive**, **transformer-based** language models (Vaswani et al., 2017). These models are trained to fit a distribution $p(\boldsymbol{y})$ over a text sequence $\boldsymbol{y} = \{y_1, y_2, \ldots, y_n\}$ via the chain rule:

$$p(\boldsymbol{y}) = p(y_1) \prod_{i=2}^{n} p(y_i \mid y_1, \ldots, y_{i-1}). \qquad (1)$$

Each term on the right hand side is parameterized by a transformer neural network, which accepts the conditioned tokens as input. Each token passes through a learned embedding layer $F_\theta$, then the full conditioned sequence is fed into the language model, which refines the representations of each token while producing a categorical distribution for the next token. Our experiments utilize a standard language model, GPT-2, that is trained on a large dataset of web text (Radford et al., 2018). We use the implementation provided in the HuggingFace `transformers` library (Wolf et al., 2019). In our work, we use the language model to process the input sequence rather than to predict future tokens.

## 3.2. POMDPs and Policy Learning

Our experiments explore the application of LMs to general sequential decision-making tasks in partially observed environments. These tasks may be formalized as partially observable Markov decision processes (POMDPs). A POMDP is defined by a set of states $\mathcal{S}$, a set of observations $\mathcal{O}$, a set of actions $\mathcal{A}$, and a transition model $\mathcal{T}(s_{t+1}|s_t, a_t)$ that predicts the next state $s_{t+1}$ based on the current state $s_t$ and action $a_t$. Importantly, in a POMDP setting, the observation $o_t$ only captures a portion of the underlying state $s_t$, and an optimal decision-making strategy (a **policy**) must incorporate both the current observation and the previous history of observations and actions. For experiments in this paper, policies are parametric models $\pi_\psi(a_{t+1}|g, h_t, o_t)$ that select actions given the goals $g$, history information $h_t$, and partial observations $o_t$ of the current state $s_t$.

In Figure 1 (right), we show a high-level overview of the proposed policy architecture. We first convert all policy inputs into a sequence and provide them as input to a transformer encoder. Representations from this encoder model are then passed to a task-specific decoder that predicts actions. All our experiments use imitation learning (Santara et al., 2017; Ng et al., 2000; Peng et al., 2018), specifically

behavior cloning (Pomerleau, 1991; 1989; Torabi et al., 2018), to train the policy. We collect a dataset of $N$ expert training trajectories $\mathcal{D} = \{d^1, \cdots, d^N\}$, where each individual trajectory consists of a goal and a sequence of observations and actions: $d^i = \{g^i, o_1^i, a_1^i, \cdots, o_{T_i}^i, a_{T_i}^i\}$, where $T_i$ is the length of the trajectory. We then train the policy to maximize the probability $p(\{\boldsymbol{a}^i\}_{i=1}^N | \mathcal{D})$ of the expert actions $\boldsymbol{a}^i = \{a_1^i, \ldots, a_{T_i}^i\}$ across trajectories using the cross-entropy loss:

$$-\ln p(\{\boldsymbol{a}^i\}_{i=1}^N | \mathcal{D}) = -\sum_{i=1}^N \sum_{t=1}^{T_i} \ln p(a_{t+1}^i | g^i, h_t^i, o_t^i) \quad (2)$$

where $h_t^i$ represents the history information before time $t$.

### 3.3. Language models as policy initializers

Both language modeling and POMDP decision-making are naturally framed as sequence prediction tasks, where successive words or actions/observations are predicted based on a sequence of previous words or actions/observations. This suggests that pre-trained LMs can be used to initialize POMDP policies by fine-tuning them to model high-reward or expert trajectories, as described below.

## 4. Approach

We use the **VirtualHome** (Puig et al., 2018) and **BabyAI** (Hui et al., 2020) environments to evaluate the proposed method. While both environments feature complex goals, the nature of these goals, as well as the state and action sequences that accomplish them, differ substantially across environments (Figure 1 (left)).

### 4.1. Environments

VirtualHome is a 3D realistic environment featuring partial observability, large action spaces, and long time horizons. It provides a set of realistic 3D homes and objects that can be manipulated to perform household organization tasks. Observations in VirtualHome are graphs describing a list of objects and their relations in the current partial observation. Each object has an object name, a state, *e.g. open, close, clean*, and 3D coordinates. The goals are defined as a set of goal predicates, *e.g.* "On(apple, table):2; Inside(milk, fridge):1". The action at each time step consists of a verb and 1 object, *e.g.* "grab" "apple". We build on the original environment to create the VirtualHome-Imitation Learning dataset using regression planning (Korf, 1987), which features more challenging tasks and more diverse scenes and objects. (See **Appendix B**).

BabyAI is a 2D grid world environment designed to evaluate instruction following. Different from VirtualHome, the observation in BabyAI by default is a grid describing a partial and local egocentric view of the state of the environ-
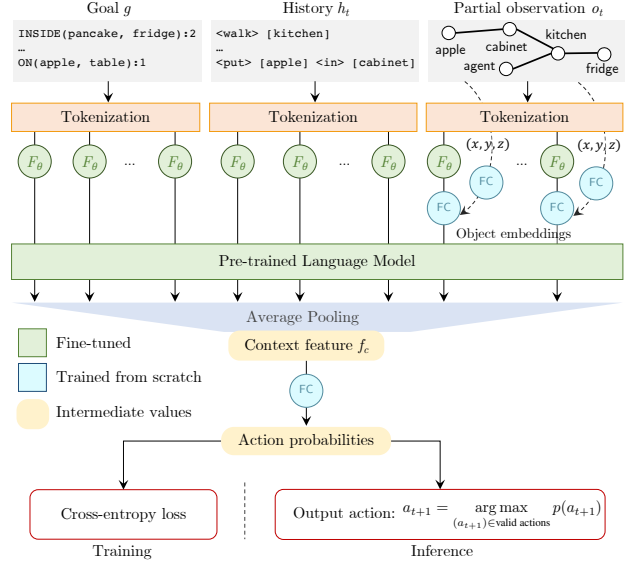


*Figure 2.* **Policy network in VirtualHome.** The observation, goal, and history are first converted into sequences and then passed through an embedding layer $F_\theta$. The combined sequence is passed through a pre-trained LM, and the output tokens are pooled into a context feature vector for action prediction.

ment. This grid shows $7 \times 7$ tiles in the direction the agent is facing. Each tile contains at most one object, encoded using 3 integer values: one storing the type of object, one storing its color, and a state for doors indicating whether it is open, closed or locked. The goals in BabyAI are described using language instructions, *e.g.* "put the blue key next to the purple ball". BabyAI has 7 actions, *e.g.* "turn left", "pick up", and "drop".

VirtualHome and BabyAI have quite different observation spaces, action spaces, and goal spaces; however, we show that embedding policy inputs as sequences and utilizing the pre-trained LM as a policy initializer, enables effective generalization to novel tasks on both environments. We note that the proposed method is not limited to VirtualHome and BabyAI, but is straightforwardly applicable to other embodied environments, such as ALFRED (Shridhar et al., 2020) and iGibson (Shen et al., 2020).

### 4.2. Policy Network

We first examine whether pre-trained LMs provide effective initializers when states and action histories are represented as natural language strings. We encode the inputs to the policy—including observations, goals, and action histories—as sequences of words. These word sequences are passed to the LM (using its pre-trained word embedding layer $F_\theta$) and used to obtain contextualized token representations. These token representations are averaged, and used to predict actions. We design a policy network following the general policy framework proposed in Figure 1. We use

VirtualHome as an example to show the policy network with detailed input encoding and task-specific model in Figure 2.

**Environment encodings in VirtualHome.** In Virtual-Home, each goal consists of a sequence of predicates and multiplicities, and is translated into a templated English sentence (*e.g.* "`Inside(apple, fridge):2`" becomes "put two apples inside the fridge"). To encode the agent's partial observation, we extract a list of currently visible objects, their states (*e.g.* "open, clean"), and 3D world coordinates. We use a fully-connected layer to encode the 3D information and generate a feature representation of each object in the observation (See **Appendix A.1**). To encode history, we store information about all previous actions and convert them into templated English sentences (*e.g.* "I have put the plate on the kitchen table and the apple inside the fridge").

**Environment encodings in BabyAI.** The observation by default is a $7 \times 7$ grid. We convert the observation into $7 \times 7$ text descriptions, *e.g.* "purple ball", "grey wall", "open door", and combine them into a long sentence. We then convert the history actions into text descriptions, *e.g.* "turn left" and "go forward". We combine the language instruction (without modification) with the observation and history text descriptions, and feed them to the pre-trained LM.

We note that the policy architecture described above does not strictly require that these encodings take the form of natural language strings—other encodings of the environment as a sequence also work. We explore several alternative encoding schemes in Section 7. This framework could be also generalized to support pixel-based observations using discretization schemes like the one employed in the Vision Transformer (Dosovitskiy et al., 2020).

**Action prediction.** To produce action probabilities, we pool the outputs of the pre-trained LM and pass this "context feature" through a fully-connected layer as shown in Figure 2. In training, we maximize the probabilities of the expert actions. In inference, we select the valid action with the highest probability.

## 5. Experiment Setup

We evaluate the effectiveness of the proposed method on VirtualHome and BabyAI.

### 5.1. VirtualHome

We build three test sets that evaluate policies' ability from three aspects: (1) performance on in-distribution tasks; (2) generalization to novel scenes; and (3) generalization to novel tasks. (See **Appendix C** for more details.)

**In-Distribution.** The predicate types and their counts in the goal are randomly sampled from the same distribution as the

training data. There are $2 \sim 10$ predicates in each task. The objects are initially placed in the environment according to common-sense layouts; (*e.g.* plates appear inside the kitchen cabinets rather than the bathtub). While goal predicates are drawn from the same distribution as the training data, the initial states in the test set are different from the training set.

**Novel Scenes.** The objects are placed in random positions in the initial environment without common-sense constraints (*e.g.* apples may appear inside the dishwasher). There are $2 \sim 10$ goal predicates in each task.

**Novel Tasks.** The components of all goal predicates are never seen together during training (*e.g.* both plates and fridge appear in training goals, but `Inside(plate, fridge)` only appears in the test set. There are $2 \sim 8$ goal predicates in each task.

We evaluate the success rates of different methods on each test set. A given episode is scored as successful if the policy completes its entire goal within $T$ steps, where $T = 70$ is the maximum allowed steps of the environment. On each of the 3 test subsets, we use 5 different random seeds and test 100 tasks under each seed. Thus there are 1500 examples used to evaluate each model.

### 5.1.1. BASELINE MODELS

We compare the proposed method with a variety of baselines using different policy architectures.

**LM-Text (Ours)** is the proposed method that converts all environments inputs into text descriptions. The pre-trained LM is fine-tuned for decision-making (conditioned on goals, observations, and histories) as described in Section 4.2.

**Recurrent Network**. We compare our method with a recurrent baseline using an LSTM (Hochreiter & Schmidhuber, 1997) to encode the history information. The hidden representation from the last time step together with the goal, and the current observation are used to predict the next action.

**MLP** and **MLP-1**. We perform additional comparisons with baselines that do not use the recurrent network and pre-trained language models. *MLP* and *MLP-1* take the goal, history actions, and the current observation as input and send them to the multilayer perceptron neural network (MLP) to predict actions. *MLP-1* has three more average-pooling layers than *MLP* that average the features of tokens in the goal, history actions, and the current observation, respectively, before sending them to the MLP layer.

### 5.2. BabyAI

BabyAI contains programmatically generated natural language instructions, e.g. "put the green ball next to the box on your left" and requires the agent to navigate the world and move objects to target locations. BabyAI provides demos
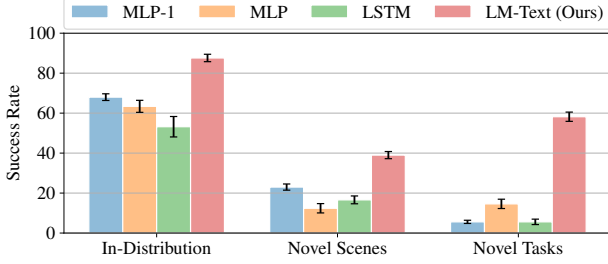
*Figure 3.* **Comparisons of the proposed method and baselines on different testing subsets on VirtualHome.** *MLP-1*, *MLP*, and *LSTM* are baselines without using the pre-trained LM. The proposed method, *LM-Text (Ours)*, outperforms all baselines.

for 19 tasks. In our experiments, we compare the proposed method with baselines on four representative tasks: *Go-ToRedBall*, *GoToLocal*, *PickupLoc*, and *PutNextLocal*. In BabyAI, performing well on the test set requires the model to generalize to new environment layouts and goals, resulting in new combinations of tasks not seen in training.

## 6. Experiments: LMs as Policy Initializers

### 6.1. VirtualHome

Each model is trained on $20K$ demos from the VirtualHome-Imitation Learning dataset, and then evaluated on the three test subsets described in Section 5.1: **In-Distribution**, **Novel Scenes**, and **Novel Tasks**. The results are shown in Figure 3. We find that *LM-Text (Ours)*, which initializes the policy with a pre-trained language model, has higher success rates than the baseline models.

This difference is most pronounced in the **Novel Tasks** setting, where test tasks require combinatorial generalization across goals that are never seen during training. Here, *LM-Text (Ours)* dramatically (**43.6%**) improves upon all baseline models. Such combinatorial generalization is necessary to construct general purpose agents, but is often difficult for existing deep learning approaches. Our results suggest that pre-trained LMs can serve as a computational backbone for combinatorial generalization. We analyze the source of such combinatorial generalization in Section 7.

### 6.2. BabyAI

In Table 1, we compare the results of the proposed method and baselines. We use the standard training and test data provided by (Hui et al., 2020). (Hui et al., 2020) is the method used in the original paper. **LM-Text (Ours)** is the proposed method that converts policy inputs into a text sequence. In BabyAI, performing well on the unseen test tasks with new environment layouts and goals requires approaches that have combinatorial reasoning capabilities.

We report the success rate given different number of training demos in Table 1. Given enough training data, *i.e.* 5K and

| Tasks | Methods | Number of Demos | | | | |
|---|---|---|---|---|---|---|
| | | **100** | **500** | **1K** | **5K** | **10K** |
| **GoToRedBall** | (Hui et al., 2020) | 81.0 | 96.0 | 99.0 | 99.5 | 99.9 |
| | LM-Text (Ours) | **93.9** | **99.4** | **99.7** | **100.0** | **100.0** |
| **GoToLocal** | (Hui et al., 2020) | 55.9 | 84.3 | 98.6 | **99.9** | 99.8 |
| | LM-Text (Ours) | **64.6** | **97.9** | **99.0** | 99.5 | 99.5 |
| **PickupLoc** | (Hui et al., 2020) | 28.0 | 58.0 | 93.3 | 97.9 | **99.8** |
| | LM-Text (Ours) | **28.7** | **73.4** | **99.0** | **99.6** | 99.8 |
| **PutNextLocal** | (Hui et al., 2020) | **14.3** | 16.8 | 43.4 | 81.2 | 97.7 |
| | LM-Text (Ours) | 11.1 | **93.0** | **93.2** | **98.9** | **99.9** |

*Table 1.* **Success rates of the proposed method vs. baseline on four BabyAI tasks.** *LM-Text (Ours)* outperforms (Hui et al., 2020), the method used in the original paper.

10K demos, both (Hui et al., 2020) and *LM-Text (Ours)* achieve high success rates on four tasks, but *LM-Text (Ours)* outperforms (Hui et al., 2020) given fewer training data. For example, *LM-Text (Ours)* is 12.9% higher than (Hui et al., 2020) on the *GoToRedBall* task given 100 training demos, indicating the proposed method can generalize well to novel language tasks using fewer training data.

## 7. Analysis: Understanding the Sources of Generalization

*LM-Text (Ours)* exhibits effective combinatorial generalization to novel tasks in VirtualHome and BabyAI. Is this simply because, as past work has observed, LMs are effective models of relations between natural language descriptions of states and actions (Ammanabrolu & Riedl, 2018; Li et al., 2021), or because they provide a more general framework for combinatorial generalization in decision-making? In this part, we selectively ablate various components of the model described above to understand the sources of its success.

### 7.1. Input Encoding Scheme

We first explore the role of *natural language* in these results by investigating three alternative ways of encoding policy inputs to our model architecture without using natural language strings: two in VirtualHome, and one in BabyAI.

**Index encoding in VirtualHome.** Different from *LM-Text (Ours)* that convert policy inputs into natural language strings, *LM-Index (Ours)* converts policy inputs into indexes. *LM-Index (Ours)* retains the discrete, serial format of the goal, history, and observation representation, but replaces each word with an index, and the embedding layer from the pre-trained language model with a new embedding layer trained from scratch. For example, *grab apple* is mapped to *"5 3"* based on the positions of *grab* and *apple* in the vocabulary set.

**Unnatural string encoding in VirtualHome.** *LM-Unnatural (Ours)* replaces the *natural language* tokens

| Methods | Number of Demos | | | | | |
|---|---|---|---|---|---|---|
| | 100 | 500 | 1K | 5K | 10K | 20K |
| LM-Text (Ours) | **8.8** | **22.2** | 26.8 | 46.0 | **58.2** | 58.2 |
| LM-Index (Ours) | 6.4 | 18.0 | 18.8 | 45.5 | 54.6 | 57.8 |
| LM-Unnatural (Ours) | 6.8 | 18.6 | **27.0** | **47.2** | 55.8 | **58.8** |

*Table 2.* **Success rates of policies trained with different input encodings in the *Novel Tasks* setting of VirtualHome.** The text encoding is most sample-efficient, but all models converge to similar performance given sufficient training data.

| Tasks | Methods | Number of Demos | | | | |
|---|---|---|---|---|---|---|
| | | 100 | 500 | 1K | 5K | 10K |
| GoToRedBall | LM-Text (Ours) | 93.9 | 99.4 | 99.7 | **100.0** | 100.0 |
| | LM-Conv (Ours) | 92.5 | 98.8 | **100.0** | 100.0 | 100.0 |
| GoToLocal | LM-Text (Ours) | 64.6 | **97.9** | 99.0 | 99.5 | 99.5 |
| | LM-Conv (Ours) | **69.5** | 86.0 | 98.2 | **99.9** | **99.9** |
| PickupLoc | LM-Text (Ours) | 28.7 | **73.4** | 99.0 | 99.6 | 99.8 |
| | LM-Conv (Ours) | 25.0 | 58.8 | 95.1 | **99.6** | **100.0** |
| PutNextLocal | LM-Text (Ours) | 11.1 | **93.0** | 93.2 | 98.9 | **99.9** |
| | LM-Conv (Ours) | **17.9** | 53.6 | 91.3 | 97.7 | 99.5 |

*Table 3.* **Success rate of policies trained with *text encoding* vs. *convolutional encoding* in BabyAI.** The text encoding is more sample-efficient, but both models converge to near perfect performance given sufficient training data.

(e.g. converting the goal "On(fork, table):1" as *put one fork on the table*) with random ones (e.g. converting On(fork, table) as *brought wise character trees fine yet*). This is done by randomly permuting the entire vocabulary, mapping each token to a new token. Such a permutation breaks the semantic information in natural strings.

*LM-Index (Ours)* and *LM-Unnatural (Ours)* have the same policy network as *LM-Text (Ours)*. All of them are fine-tuned on our expert data. Each of them is trained on different numbers of demos. The averaged results using 5 different random seeds on the Novel Tasks setting are reported in Table 2. Given few training data, e.g. 100 demos, all the models perform poorly, with success rates lower than 10%. *LM-Text (Ours)* achieves higher success rates than *LM-Index (Ours)* and *LM-Unnatural (Ours)* when the training data is increased, *e.g. LM-Text (Ours)* is around 4% higher than *LM-Index (Ours)* and *LM-Unnatural (Ours)* when there are 500 training demos. However, when the training dataset is further enlarged, *e.g.* 20K demos, success rates of all approaches reach similar performance. Such a result indicates that the effectiveness of pre-trained LMs in compositional generalization stems is not unique to natural language strings, but can be leveraged from arbitrary encodings, although adapting the model to arbitrary encodings may require more training data.

**Convolutional encoding in BabyAI.** We test a new model, **LM-Conv (Ours)**, that converts environment inputs into
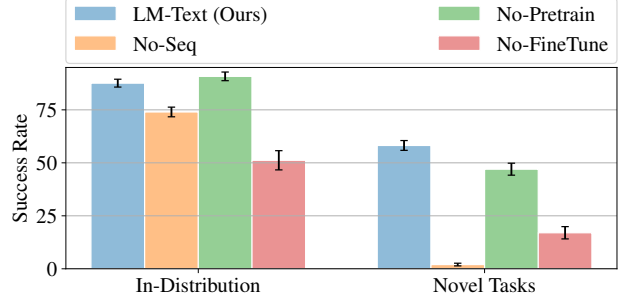


*Figure 4.* **Experiments on sequential representation and parameter pre-training.** *ML-Text (Ours)* refines a pre-trained LM while *No-Pretrain* learns it from scratch. *No-FineTune* freezes the pre-trained weights. *No-Seq* uses non-sequential inputs. Fine-tuning the pre-trained weights and the usage of sequential encoding are important for combinatorial generalization.

*convolutional embeddings*. We pass the $7 \times 7 \times 3$ grid observation in BabyAI to convolutional layers and obtain a $7 \times 7 \times d$ feature map, where $d$ is the feature dimension. We flatten the feature map and get a sequence of features to describe the observation. The rest parts are the same as *LM-Text (Ours)*. Table 3 shows the results of policies using the *text encoding* and *convolutional encoding*. *LM-Text (Ours)* and *LM-Conv (Ours)* have similar results given enough training data, but *LM-Text (Ours)* is slightly better when there are fewer training data. This conclusion is coincident with the results on VirtualHome.

Different input encoding schemes have only a negligible impact on model performance: the effectiveness of pre-training is not limited to utilizing natural strings, but in fact extends to arbitrary sequential encodings.

### 7.2. Sequential Representation

We want to further investigate whether the LM pre-trained policy will still be effective when the input encoding is not sequential. We thus compare the proposed method with a baseline, **No-Seq**, that uses the non-sequential inputs. *No-Seq* uses three features to represent the goal, history, and observation, respectively. It encodes the goal as a single vector by averaging the features of all goal embeddings. The history and observation features are obtained in the same way. All the features are then sent to the pre-trained LM to predict actions. As shown in Figure 4, removing the sequential structure significantly hurts the performance on the *Novel Tasks* setting. *No-Seq* can achieve good performance on test tasks that are closer to training tasks, but cannot generalize well to more challenging unseen tasks. Thus to enable the combinatorial generalization of pre-trained LMs, it is important to convert environment inputs into sequential embeddings.
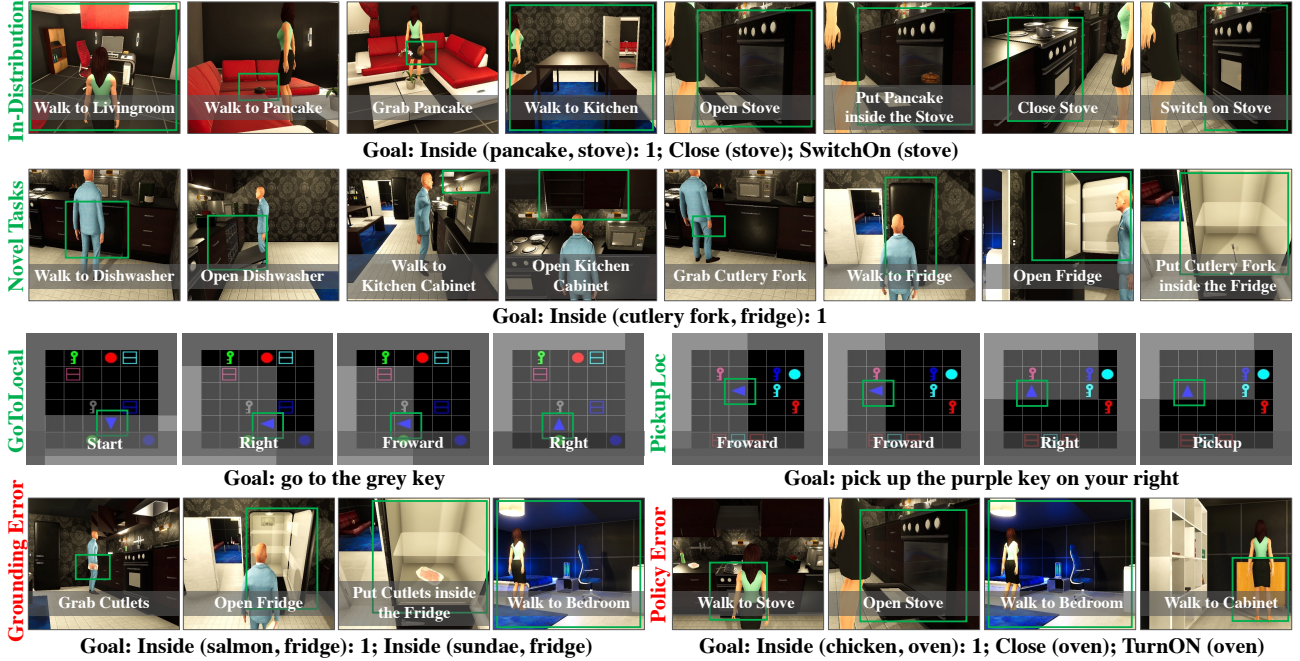
*Figure 5.* **Qualitative results of our model.** We show two examples of successes in VirtualHome, two examples of successes in BabyAI, and two failure cases caused by the grounding error and policy error. We only show a sub-trajectory in each example and omit most exploration actions to save space. The interacted objects are labelled by green bounding boxes.

## 7.3. Parameter Pre-training

Finally, we verify whether the pre-trained LM weights contribute to the effective generalization of the proposed model. We first compare the proposed model with a baseline model, **No-Pretrain**, that does not initialize the policy using the pre-trained LM, but instead training the policy on our expert data from scratch. In Figure 4, we find that removing the pre-trained weights can fit the in-domain data and thus performs well on the *In-Distribution* setting. However, its success rate is 11.2% lower than the proposed model on the *Novel Tasks* setting, indicating the pre-trained weights are important for effective generalization.

We further test a baseline, **No-FineTune**, that keeps the pre-trained weights of the language model but freezes them while training the rest model on our expert data. Freezing the pre-trained weights significantly hurts the performance on both the *In-Distribution* and *Novel Tasks* settings, suggesting that fine-tuning of the transformer weights is essential for effective combinatorial generalization to novel tasks.

Together, these results suggest that sequential representations (vs. fixed-dimensional feature vectors) and the LM objective (not just the transformer architecture) are both important for generalization, however, the input encoding schemes (e.g. as a natural language string vs. an arbitrary encoding scheme) has little influence. They point to the potential broader applicability of pre-trained LMs as a computational backbone for compositional embodied decision

making, where arbitrary inputs, such as language, images, or grids, may be converted to sequential encodings.

## 8. Qualitative Results

In Figure 5, we show examples of *LM-Text (Ours)* completing tasks in VirtualHome and BabyAI. Given a VirtualHome task described by goal predicates, *e.g.* `Inside(cutlery fork, fridge):1`, the agent first explores the environment until it finds the *cutlery fork*. Then the agent grabs the *cutlery fork*, walks to the *fridge*, and puts the *cutlery fork* inside the *fridge*. Then the agent moves to the next subtask until it finishes all the subtasks. In Figure 5, we show two examples of successes from the *In-Distribution* and *Novel Tasks* settings. We show another two successful examples from BabyAI on solving the *GoToLocal* and *PickupLoc* tasks. We only show short trajectories or extract a sub-trajectory in each example. Most exploration actions are omitted for saving space.

**Failure case analysis.** We observed two main types of failure cases: grounding error and policy error. For failures caused by the grounding error, the agent interacts with a wrong object that is not related to the given goal, *e.g.* the agent puts *cutlets* instead of the *salmon* inside the fridge. For failures caused by the policy error, the agent cannot find the target objects or does not interact with them. The proposed method that converts policy inputs into sequential encodings and feeds them to the general LM framework can accomplish decision-making tasks efficiently, however,

there are still challenging tasks that the policy fails to accomplish. Larger LMs, *e.g*. GPT-3 (Brown et al., 2020), may improve the success rate of those challenging tasks.

## 9. Conclusion

In this paper, we propose to use LM pre-training as a general framework for decision-making by converting goals and observations into sequences and sending them to a policy network initialized with a pre-trained LM. The proposed method enables effective combinatorial generalization across environments with substantially different states, actions, and goals. We find that sequential representations and LM pre-training both contribute to the effective generalization while the format of the input encoding has little influence. These results show great promise for pre-trained language models as a general-purpose framework in decision-making tasks.

## References

Abdou, M., Kulmizev, A., Hershcovich, D., Frank, S., Pavlick, E., and Søgaard, A. Can language models encode perceptual structure without grounding? a case study in color. *arXiv preprint arXiv:2109.06129*, 2021.

Ammanabrolu, P. and Riedl, M. O. Playing text-adventure games with graph-based deep reinforcement learning. *arXiv preprint arXiv:1812.01628*, 2018.

Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.

Côté, M.-A., Kádár, A., Yuan, X., Kybartas, B., Barnes, T., Fine, E., Moore, J., Hausknecht, M., El Asri, L., Adada, M., et al. Textworld: A learning environment for text-based games. In *Workshop on Computer Games*, pp. 41–75. Springer, 2018.

Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., and Harshman, R. Indexing by latent semantic analysis. *Journal of the American society for information science*, 41(6):391–407, 1990.

Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

Dumais, S. T. Latent semantic analysis. *Annual review of information science and technology*, 38(1):188–230, 2004.

Fried, D., Hu, R., Cirik, V., Rohrbach, A., Andreas, J., Morency, L.-P., Berg-Kirkpatrick, T., Saenko, K., Klein, D., and Darrell, T. Speaker-follower models for vision-and-language navigation. *arXiv preprint arXiv:1806.02724*, 2018.

Hill, F., Mokra, S., Wong, N., and Harley, T. Human instruction-following with deep reinforcement learning via transfer-learning from text. *arXiv preprint arXiv:2005.09382*, 2020.

Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

Huang, W., Abbeel, P., Pathak, D., and Mordatch, I. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. *arXiv preprint arXiv:2201.07207*, 2022.

Hui, D. Y.-T., Chevalier-Boisvert, M., Bahdanau, D., and Bengio, Y. Babyai 1.1, 2020.

Ilharco, G., Zellers, R., Farhadi, A., and Hajishirzi, H. Probing text models for common ground with visual representations. *arXiv e-prints*, pp. arXiv–2005, 2020.

Keskar, N. S., McCann, B., Varshney, L. R., Xiong, C., and Socher, R. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*, 2019.

Kitaev, N., Cao, S., and Klein, D. Multilingual constituency parsing with self-attention and pre-training. *arXiv preprint arXiv:1812.11760*, 2018.

Korf, R. E. Planning as search: A quantitative approach. *Artificial intelligence*, 33(1):65–88, 1987.

Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.

Li, B. Z., Nye, M., and Andreas, J. Implicit representations of meaning in neural language models. *arXiv preprint arXiv:2106.00737*, 2021.

Lu, J., Batra, D., Parikh, D., and Lee, S. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *arXiv preprint arXiv:1908.02265*, 2019.

Lu, K., Grover, A., Abbeel, P., and Mordatch, I. Pretrained transformers as universal computation engines. *arXiv preprint arXiv:2103.05247*, 2021.

Majumdar, A., Shrivastava, A., Lee, S., Anderson, P., Parikh, D., and Batra, D. Improving vision-and-language navigation with image-text pairs from the web. In *European Conference on Computer Vision*, pp. 259–274. Springer, 2020.

Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pp. 3111–3119, 2013.

Narasimhan, K., Kulkarni, T., and Barzilay, R. Language understanding for text-based games using deep reinforcement learning. *arXiv preprint arXiv:1506.08941*, 2015.

Ng, A. Y., Russell, S. J., et al. Algorithms for inverse reinforcement learning. In *Icml*, volume 1, pp. 2, 2000.

Peng, X. B., Kanazawa, A., Toyer, S., Abbeel, P., and Levine, S. Variational discriminator bottleneck: Improving imitation learning, inverse rl, and gans by constraining information flow. *arXiv preprint arXiv:1810.00821*, 2018.

Pennington, J., Socher, R., and Manning, C. D. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 1532–1543, 2014.

Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. Deep contextualized word representations. *arXiv preprint arXiv:1802.05365*, 2018.

Platanios, E. A., Pauls, A., Roy, S., Zhang, Y., Kyte, A., Guo, A., Thomson, S., Krishnamurthy, J., Wolfe, J., Andreas, J., et al. Value-agnostic conversational semantic parsing. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 3666–3681, 2021.

Pomerleau, D. A. Alvinn: An autonomous land vehicle in a neural network. Technical report, CARNEGIE-MELLON UNIV PITTSBURGH PA ARTIFICIAL INTELLIGENCE AND PSYCHOLOGY . . . , 1989.

Pomerleau, D. A. Efficient training of artificial neural networks for autonomous navigation. *Neural computation*, 3(1):88–97, 1991.

Puig, X., Ra, K., Boben, M., Li, J., Wang, T., Fidler, S., and Torralba, A. Virtualhome: Simulating household activities via programs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8494–8502, 2018.

Puig, X., Shu, T., Li, S., Wang, Z., Tenenbaum, J. B., Fidler, S., and Torralba, A. Watch-and-help: A challenge for social perception and human-ai collaboration. *arXiv preprint arXiv:2010.09890*, 2020.

Radford, A., Narasimhan, K., Salimans, T., and Sutskever, I. Improving language understanding by generative pre-training. 2018.

Roberts, A., Raffel, C., and Shazeer, N. How much knowledge can you pack into the parameters of a language model? *arXiv preprint arXiv:2002.08910*, 2020.

Santara, A., Naik, A., Ravindran, B., Das, D., Mudigere, D., Avancha, S., and Kaul, B. Rail: Risk-averse imitation learning. *arXiv preprint arXiv:1707.06658*, 2017.

Shen, B., Xia, F., Li, C., Martín-Martín, R., Fan, L., Wang, G., Buch, S., D'Arpino, C., Srivastava, S., Tchapmi, L. P., et al. igibson, a simulation environment for interactive tasks in large realisticscenes. *arXiv preprint arXiv:2012.02924*, 2020.

Shridhar, M., Thomason, J., Gordon, D., Bisk, Y., Han, W., Mottaghi, R., Zettlemoyer, L., and Fox, D. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10740–10749, 2020.

Suglia, A., Gao, Q., Thomason, J., Thattai, G., and Sukhatme, G. Embodied bert: A transformer model for embodied, language-guided visual task completion. *arXiv preprint arXiv:2108.04927*, 2021.

Torabi, F., Warnell, G., and Stone, P. Behavioral cloning from observation. *arXiv preprint arXiv:1805.01954*, 2018.

Tsimpoukelli, M., Menick, J., Cabi, S., Eslami, S., Vinyals, O., and Hill, F. Multimodal few-shot learning with frozen language models. *arXiv preprint arXiv:2106.13884*, 2021.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.

Wang, B., Shin, R., Liu, X., Polozov, O., and Richardson, M. Rat-sql: Relation-aware schema encoding and linking for text-to-sql parsers. *arXiv preprint arXiv:1911.04942*, 2019.

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., et al. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.

Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R. R., and Le, Q. V. Xlnet: Generalized autoregressive pretraining for language understanding. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings. neurips.cc/paper/2019/file/ dc6a7e655d7e5840e66733e9ee67cc69-Paper. pdf.

Yang, Z., Garcia, N., Chu, C., Otani, M., Nakashima, Y., and Takemura, H. Bert representations for video question answering. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1556–1565, 2020.

Yao, S., Rao, R., Hausknecht, M., and Narasimhan, K. Keep calm and explore: Language models for action generation in text-based games. *arXiv preprint arXiv:2010.02903*, 2020.

Yao, S., Narasimhan, K., and Hausknecht, M. Reading and acting while blindfolded: The need for semantics in text game agents. *arXiv preprint arXiv:2103.13552*, 2021.

Yuan, X., Côté, M.-A., Sordoni, A., Laroche, R., Combes, R. T. d., Hausknecht, M., and Trischler, A. Counting to explore and generalize in text-based games. *arXiv preprint arXiv:1806.11525*, 2018.

Zhang, Y. and Chai, J. Hierarchical task learning from language instructions with unified transformers and self-monitoring. *arXiv preprint arXiv:2106.03427*, 2021.

Zhu, J., Xia, Y., Wu, L., He, D., Qin, T., Zhou, W., Li, H., and Liu, T.-Y. Incorporating bert into neural machine translation. *arXiv preprint arXiv:2002.06823*, 2020.

# Appendix

In this appendix, we first provide more implementation details of the proposed model in Appendix A. We then introduce the VirtualHome-Imitation Learning dataset in Appendix B. Appendix C lists the goal predicates in different test subsets. Finally, we visualize the attention weights in language models in Appendix D.

## A. More implementation Details of the Proposed Method in VirtualHome

In Appendix A.1, we provide more details of the model architecture used in the main paper Section 4.2. We then introduce the training detail in Appendix A.2.

### A.1. Model architecture details

In the main paper Section 4.2, we introduced the model architecture used for training an interactive policy. Our model consists of three parts, *i.e.* inputs, the pre-trained language model, and outputs. We take the goal $g$, history $h_t$, and the current partial observation $o_t$ as inputs and send them to the LM-pretrained policy. The output action $a_{t+1}$ in VirtualHome consists of a verb and an object. For brevity, we will omit the time subscript $t$ from now on.

In VirtualHome, the partial observation $o$ of the environment state can be represented as a list of objects in the agent's view. We represent each object by its name, *e.g.* "oven", a state description, *e.g.* "open, clean", and position both in the world and relative to the agent. In this part, we provide more details of how **LM-Text (Ours)** encodes the name, state, and position of each object in the current observation. Figure 6 shows the model architecture we used to encode the observation.

**Name encoding.** For each object node, we serialize its object name as an English phrase $s^o$. We extract its tokens and features using the tokenizer and the embedding layer of the pre-trained language model, respectively. Since one object name might generate several English tokens using the tokenizer from the pre-trained language model, *e.g.* the tokens of "kitchencabinet" is $[15813, 6607, 16212, 500]$, we take the averaged features of all the tokens in the object name and obtain a "name" feature $f_i^{o,\text{name}}$ for each object node as shown in Figure 6.

**State encoding.** Some objects have a state description, *e.g.* "oven: open, clean". There are six types of object states in the training dataset: "clean", "closed", "off", "on", "open", and "none". For each object node, we use a binary vector to represent its state. Taking the "oven" as an example, if the oven is open and clean, its state vector would be $[1, 0, 0, 0, 1, 0]$. This state vector is then passed through a fully-connected layer to generate a state feature $f_i^{o,\text{state}}$ of object $o_i$.

**Position encoding.** To encode the position information of each object $o_i$, we take their world coordinates $\{o_{i,x}, o_{i,y}, o_{i,z}\}$ and their spatial distance to the agent $\{a_x, a_y, a_z\}$ to generate a position vector $[o_{i,x}, o_{i,y}, o_{i,z}, o_{i,x} - a_x, o_{i,y} - a_y, o_{i,z} - a_z]$. This position vector is then passed through two fully-connected layers with a ReLU layer in the middle to generate a position feature $f_i^{o,\text{position}}$ of object $o_i$.

The final feature $f_i^o$ of each object node is obtained by passing the concatenation of its name feature $f_i^{o,\text{name}}$, state feature $f_i^{o,\text{state}}$, and position feature $f_i^{o,\text{position}}$ through a fully-connect layer. The observation at a single step can be written as a set of features $\{f_1^o, \cdots, f_N^o\}$, where $N$ is the number of objects in the current observation $o$.

### A.2. Training details

Our proposed approach and baselines are trained on Tesla 32GB GPUs. We train every single model on 1 Tesla 32GB GPU. All experiments used the AdamW optimizer with the learning rate of $10^{-5}$. We utilize a standard pre-trained language model, GPT-2 that is trained on Webtext dataset (Radford et al., 2018), in our experiments by using the Huggingface library (Wolf et al., 2019).

## B. Expert Data Collection in VirtualHome

In this section, we provide more details of collecting expert data in VirtualHome for training the decision-making policies. We first introduce the VirtualHome environment and then the VirtualHome-Imitation Learning dataset.
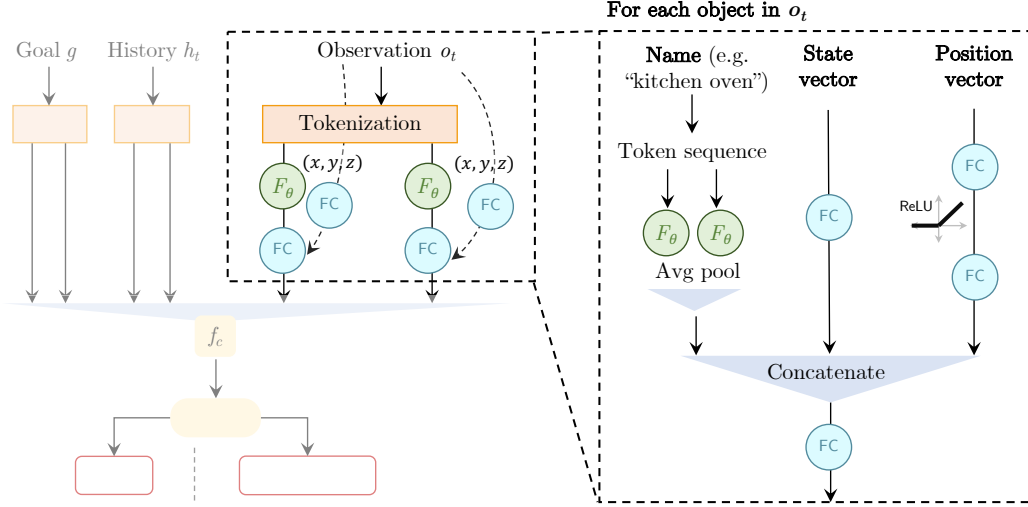
*Figure 6.* **Object encoding.** In VirtualHome, the partial observation of the environment state can be represented as a list of objects in the agent's view. Each object is represented by a name, a state vector, and position vector. **Object name encoding:** each object's name is an English phrase. We tokenize the phrase, embed the tokens, and average the embeddings. **Object state encoding:** each object is assigned one of six states: "clean", "closed", "off", "on", "open", or "none". This state is represented as a 6-dimensional binary vector and passed through a fully-connected layer. **Object position encoding:** an object's position vector is a 6-dimensional vector containing its world coordinates alongside its displacement to the agent (*i.e.* the signed difference in their world coordinates). This position vector is passed through two fully-connected layers. These three features are concatenated and passed through a fully-connected layer to obtain the object representation.

### B.1. VirtualHome

VirtualHome is a 3D realistic environment featuring partial observability, large action spaces, and long time horizons. It provides a set of realistic 3D homes and objects that can be manipulated to perform household organization tasks.

**Goal Space.** For each task, we define the goal as a set of predicates and multiplicities. For example, `Inside(apple, fridge):2; Inside(pancake, fridge):1;` means "put two apples and one pancake inside the fridge". In each task, the initial environment (including initial object locations), the goal predicates, and their orders and multiplicities are randomly sampled. There are 59 different types of predicates in total. Each predicate can appear in a goal with a multiplicity between 0 and 3. See Appendix C for the full list of goal predicates.

**Observation Space.** The observation in VirtualHome by default is a graph describing a list of objects and their relations in the current partial observation.

**Action Space.** Agents can navigate in the environment and interact with objects. To interact with an object, the agent must predict an action name and the index of the interested object, *e.g.* `Open(5)` to opening the object with index (5). The agent can only interact with objects that are in the current observation or execute the navigation actions, such as `Walk(bathroom)`. For some actions, such as `open`, the agent must be close to the object. There are also strict preconditions for actions, *e.g.* the agent must `grab` an object before it can `put` the object on a target position. As a result of these constraints, the subset of actions available to the agent changes at every timestep.

### B.2. VirtualHome-Imitation Learning Dataset

To train the models, we collect a set of expert trajectories in VirtualHome using regression planning (RP) (Korf, 1987). We follow the implementation of the regression planner used in (Puig et al., 2020). Given a task described by goal predicates, the planner generates an action sequence to accomplish this task. As shown in Figure 7, the agent has a belief about the environment, *i.e.* an imagined distribution of object locations. As the agent explores the environment, its belief of the world becomes closer to the real world. At every step, the agent updates its belief based on the latest observation (see (Puig et al., 2020)), finds a new plan using the regression planner, and executes the first action of the plan. If the subtask (described by the goal predicate) has been finished, the agent will select a new unfinished subtask, otherwise, the agent will keep doing this subtask until it finishes.

Similarly to previous work (Shridhar et al., 2020; Shen et al., 2020; Puig et al., 2020), we generate training data using a planner that has access to privileged information, such as full observation of the environment and information about the pre-conditions and effects of each action. The planner allows an agent to robustly perform tasks in partially observable environments and generate expert trajectories for training and evaluation. We generate $20,000$ trajectories for training and $3,000$ trajectories for validation. Each trajectory has a goal, an action sequence, and the corresponding observations after executing each action.
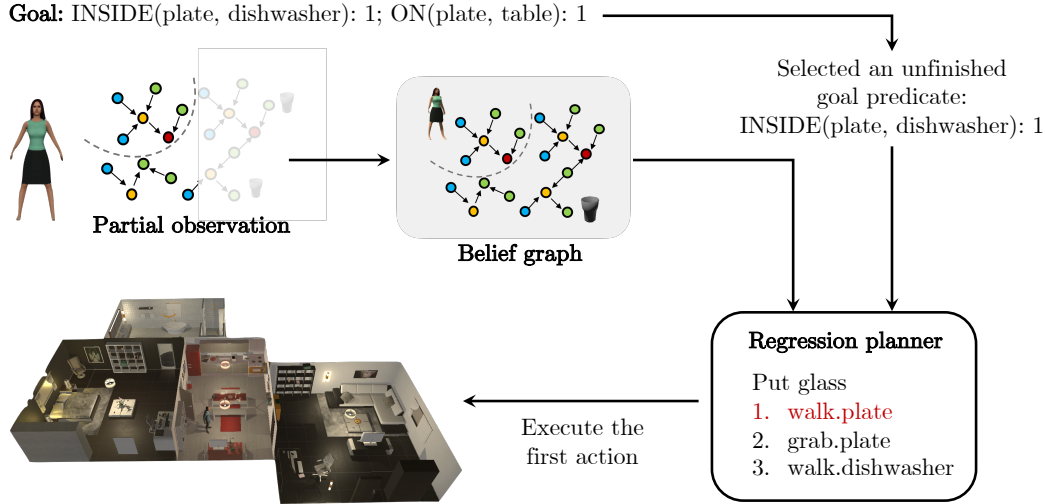


*Figure 7.* **Regression planner**. Given a task described by goal predicates, the planner generates an action sequence to accomplish this task. The agent has a belief about the environment, *i.e.* an imagined distribution of object locations. As the agent explores the environment, its belief of the world becomes closer to the real world. At every step, the agent updates its belief based on the latest observation, finds a new plan using the regression planner, and executes the first action of the plan. If the subtask (described by the goal predicate) has been finished, the agent will select a new unfinished subtask, otherwise, the agent will keep doing this subtask until finish it.

## C. Test subsets

In this section, we provide more details of each test subset. In Table 4, we provide a detailed description of each subset, including the count of goal predicate types and the number of goal predicates in each task. The **In-Distribution** setting has 37 goal predicates in total and each task has $2 \sim 10$ goal predicates. The tasks are drawn from the same distribution as the training tasks. The **Novel Scenes** setting also has 37 goal predicates and each task has $2 \sim 10$ goal predicates. The objects are randomly placed in the initial environment. The **Novel Tasks** setting has 22 goal predicates in total and each task has $2 \sim 8$ goal predicates. The tasks are never seen during training. In Figure 5 and Figure 6, we list the full goal predicates used in each subset. The goal predicates of a task are randomly sampled from the corresponding predicates pool.

## D. Visualization of Attention Weights

To better understand how does LM pre-trained policies make decisions, we visualize the attention weights from the self-attention layers of GPT-2 (Vaswani et al., 2017) in Figure 8 and Figure 9. We show the attention weights from the input to the output of **LM-Text (Ours)**. The order of tokens in the input and ouput is observation, goal, and history.

Figure 8 illustrates the attention weights of a layer named "Head 3 Layer 2". We show attention weights on two different tasks. We find that "Head 3 Layer 2" can capture objects in the goal predicates, such as "wineglass" and "cutleryfork" in the left figure, and "pancake" and "chicken" in the right figure (the figures are cropped for visualization).

Figure 9 illustrates the attention weights of layers named "Head 1 Layer 2" (left) and "Head 4 Layer 11" (right). Given the goal predicates, history, and the current observation, "LM (ft)" predicts the next action, *i.e.* "grab milk". We find that "Head 1 Layer 2" is able to capture objects in the goal predicates, such as "milk", "pancake", and "chicken" while "Head 4 Layer 11" focuses on the interacted object in the predicted action, such as "milk".

The attention weights from different self-attention layers are significantly different—some self-attention layers assign high

attention weight to objects in the goal predicates while some layers focus on the interacted object. There are also some layers that do not have interpretable meanings. The attention weights just provide us an intuition of how does the internal language model works, more quantified results are reported in the main paper.

*Table 4.* **Summary of test subsets.** We show the count of goal predicate types and the number of goal predicates used in each task.

| Test Sets | Pred. Types | #Pred. Per Task | Compared with the training set |
|---|---|---|---|
| **In-Distribution** | 37 | $2 \sim 10$ | Tasks are drawn from the same distribution as training tasks. |
| **Novel Scenes** | 37 | $2 \sim 10$ | The objects are randomly placed in the initial environment. |
| **Novel Tasks** | 22 | $2 \sim 8$ | Tasks are unseen during training. |

*Table 5.* **Goal predicates used in the *In-Distribution* and *Novel Scenes* settings.**

```
ON(cutleryfork, kitchentable):0 ~ 3
ON(plate, kitchentable):0 ~ 3
ON(waterglass, kitchentable):0 ~ 3
ON(wineglass, kitchentable):0 ~ 3
INSIDE(cutleryfork, dishwasher):0 ~ 3
INSIDE(plate, dishwasher):0 ~ 3
INSIDE(waterglass, dishwasher):0 ~ 3
INSIDE(wineglass, dishwasher):0 ~ 3
INSIDE(cutleryfork, sink):0 ~ 3
INSIDE(plate, sink):0 ~ 3
INSIDE(waterglass, sink):0 ~ 3
INSIDE(wineglass, sink):0 ~ 3
INSIDE(milk, fridge):0 ~ 3
INSIDE(chicken, fridge):0 ~ 3
INSIDE(cupcake, fridge):0 ~ 3
INSIDE(pancake, fridge):0 ~ 3
INSIDE(poundcake, fridge):0 ~ 3
ON(milk, kitchentable):0 ~ 3
ON(chicken, kitchentable):0 ~ 3
ON(cupcake, kitchentable):0 ~ 3
ON(pancake, kitchentable):0 ~ 3
ON(poundcake, kitchentable):0 ~ 3
INSIDE(chicken, microwave):0 ~ 3
INSIDE(cupcake, microwave):0 ~ 3
INSIDE(pancake, microwave):0 ~ 3
INSIDE(poundcake, microwave):0 ~ 3
INSIDE(chicken, oven):0 ~ 3
INSIDE(cupcake, oven):0 ~ 3
INSIDE(pancake, oven):0 ~ 3
INSIDE(poundcake, oven):0 ~ 3
CLOSE(oven):0 ~ 1
CLOSE(dishwasher):0 ~ 1
CLOSE(microwave):0 ~ 1
CLOSE(fridge):0 ~ 1
TurnON(oven):0 ~ 1
TurnON(dishwasher):0 ~ 1
TurnON(microwave):0 ~ 1
```

*Table 6.* **Goal predicates used in the *Novel Tasks* setting.**

```
INSIDE(milk, dishwasher):0 ~ 3
INSIDE(chicken, dishwasher):0 ~ 3
INSIDE(cupcake, dishwasher):0 ~ 3
INSIDE(pancake, dishwasher):0 ~ 3
ON(milk, sink):0 ~ 3
ON(chicken, sink):0 ~ 3
ON(cupcake, sink):0 ~ 3
ON(pancake, sink):0 ~ 3
INSIDE(cutleryfork, fridge):0 ~ 3
INSIDE(plate, fridge):0 ~ 3
INSIDE(waterglass, fridge):0 ~ 3
INSIDE(wineglass, fridge):0 ~ 3
INSIDE(cutleryfork, microwave):0 ~ 3
INSIDE(plate, microwave):0 ~ 3
INSIDE(waterglass, microwave):0 ~ 3
INSIDE(wineglass, microwave):0 ~ 3
INSIDE(milk, microwave):0 ~ 3
INSIDE(cutleryfork, oven):0 ~ 3
INSIDE(plate, oven):0 ~ 3
INSIDE(waterglass, oven):0 ~ 3
INSIDE(wineglass, oven):0 ~ 3
INSIDE(milk, oven):0 ~ 3
```
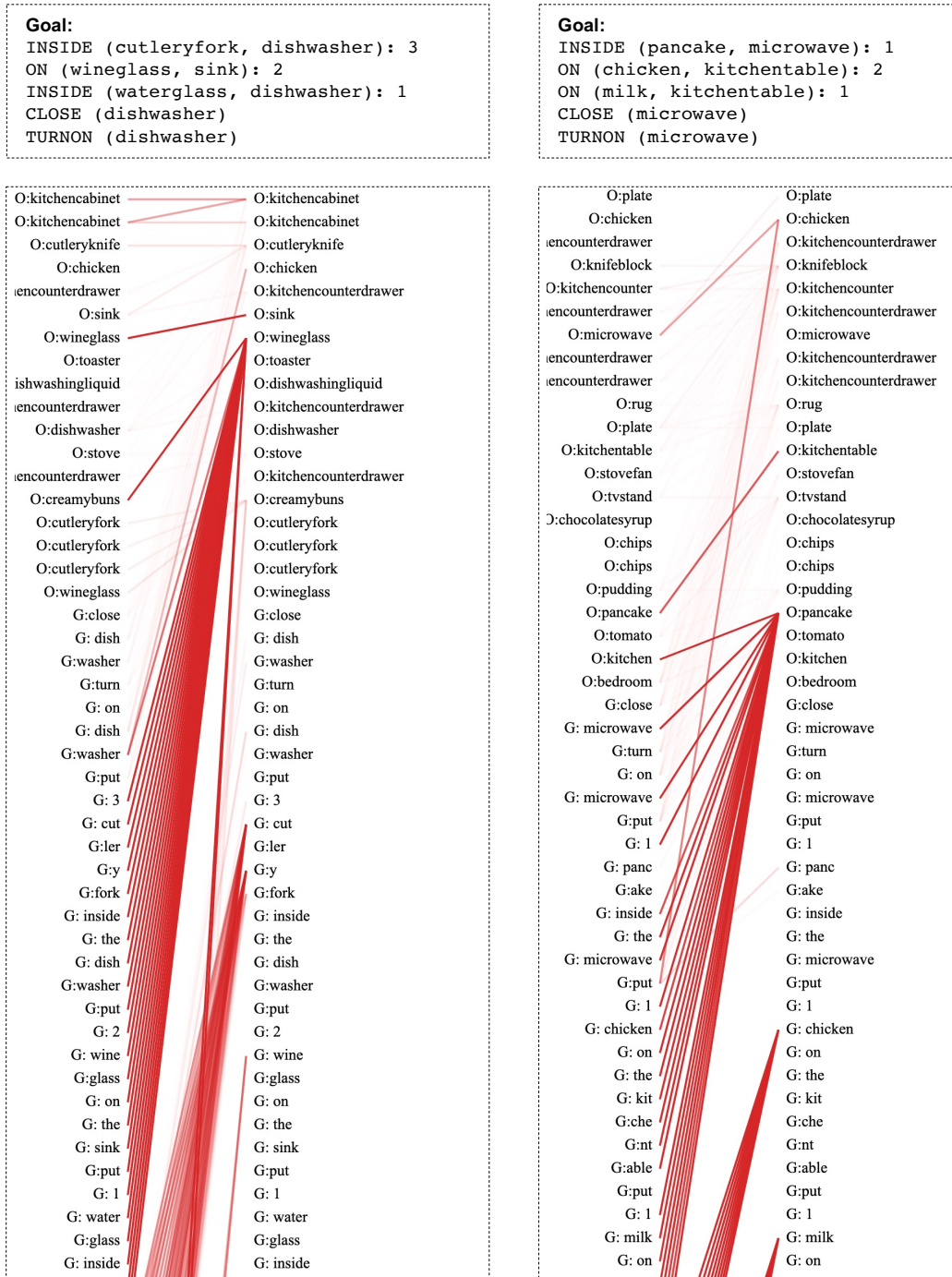
*Figure 8.* **Attention weights of a layer named "Head 3 Layer 2".** We show attention weights on two different tasks. We find that "Head 3 Layer 2" is able to capture objects in the goal predicates, such as "wineglass" and "cutleryfork" in the left figure, and "pancake" and "chicken" in the right figure (the figures are cropped for visualization).
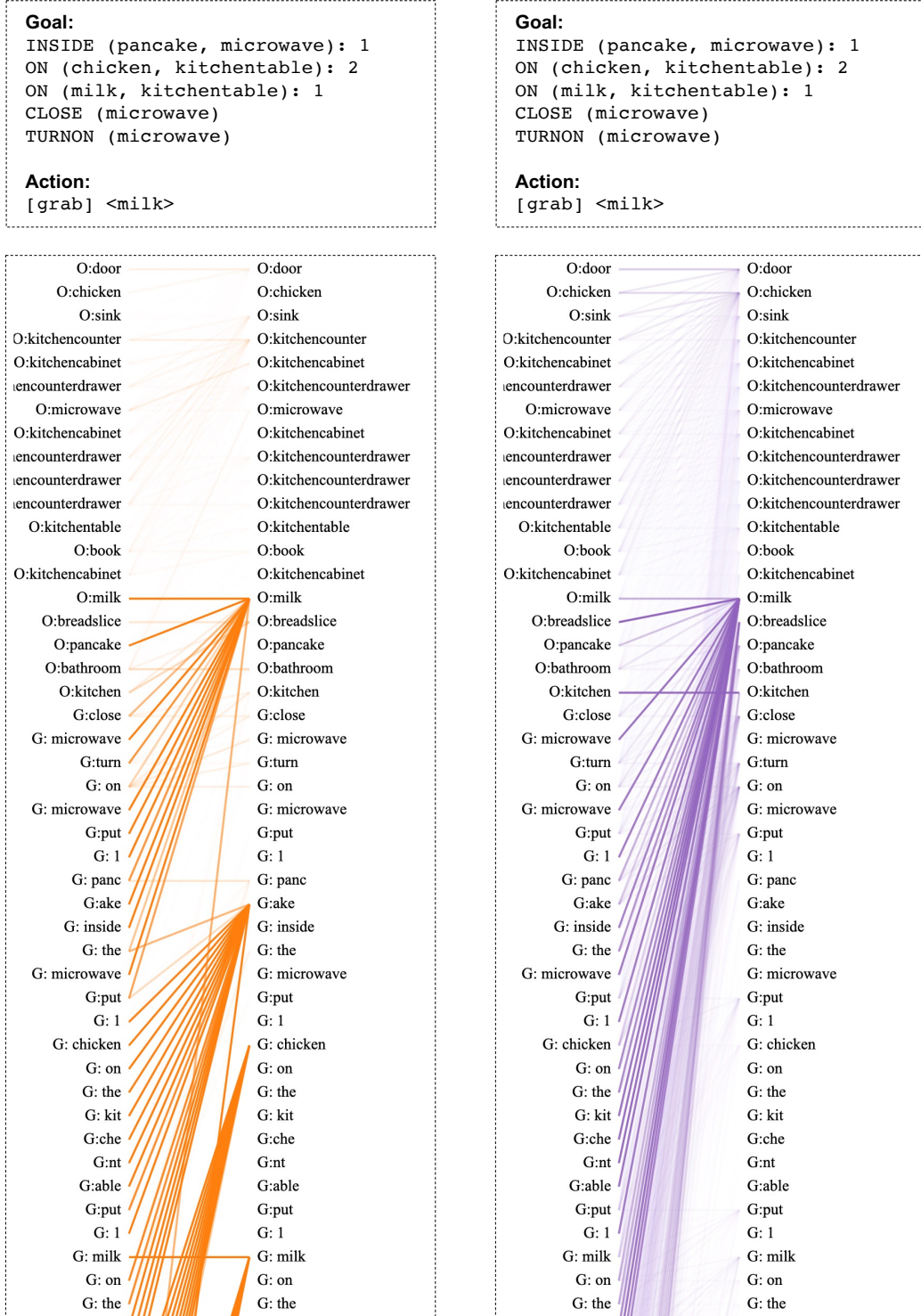
*Figure 9.* **Attention weights of layers named "Head 1 Layer 2" (left) and "Head 4 Layer 11" (right).** Given the goal predicates, history, and the current observation, the policy model predicts the next action, *i.e.* "grab milk". We find that "Head 1 Layer 2" can capture objects in the goal predicates, such as "milk", "pancake", and "chicken" while "Head 4 Layer 11" focuses on the interacted object in the predicted action, such as "milk".