

Abstract. Reinforcement learning requires interaction with environments, which can be prohibitively expensive, especially in robotics. This constraint necessitates approaches that work with limited environmental interaction by maximizing the reuse of previous experiences. We propose an approach that maximizes experience reuse while learning to solve a given task by generating and simultaneously learning useful auxiliary tasks. To generate these tasks, we construct an abstract temporal logic representation of the given task and leverage large language models to generate context-aware object embeddings that facilitate object replacements. Counterfactual reasoning and off-policy methods allow us to simultaneously learn these auxiliary tasks while solving the given target task. We combine these insights into a novel framework for multitask reinforcement learning and experimentally show that our generated auxiliary tasks share similar underlying exploration requirements as the given task, thereby maximizing the utility of directed exploration. Our approach allows agents to automatically learn additional useful policies without extra environment interaction.

Motivation

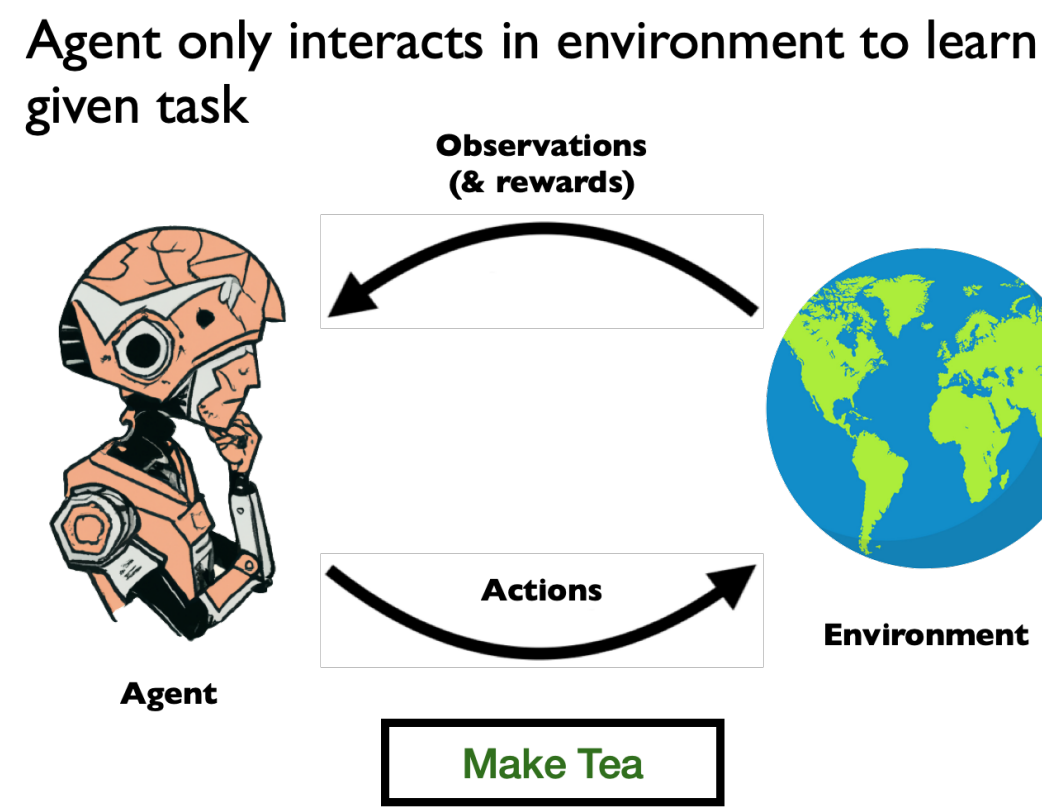
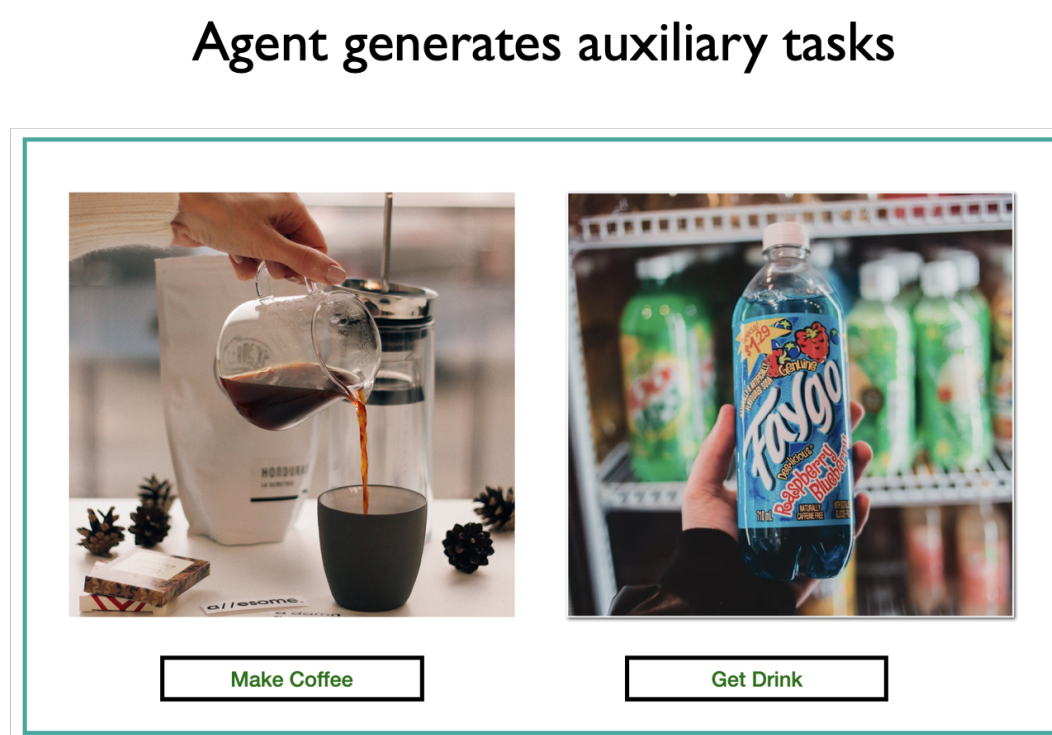
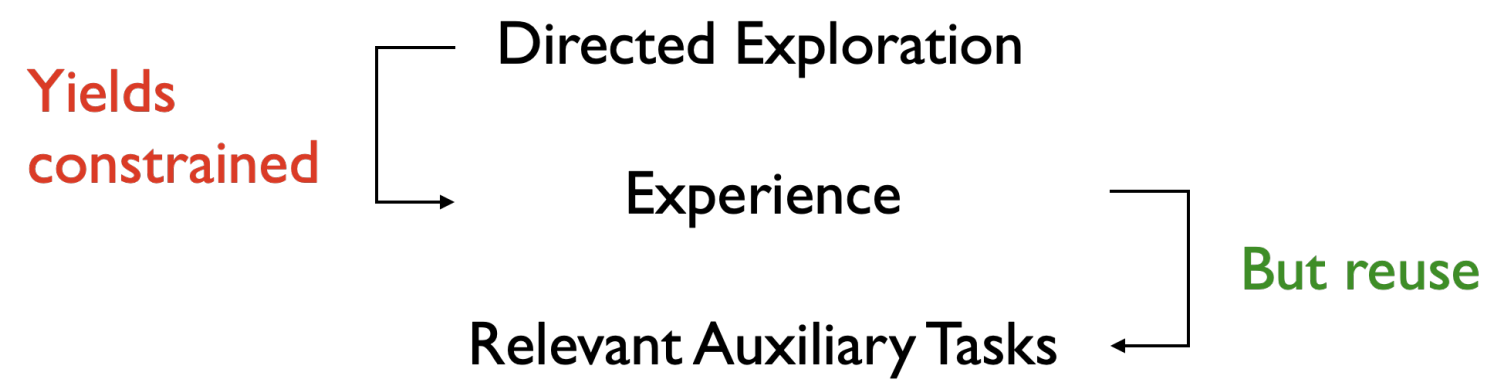
How do we get reinforcement learning agents to maximize the utility of their environment interaction experiences??

Proposal. Generate and learn auxiliary tasks that maximally benefit from the constrained exploration experience of single task curricula.

- Exploit the structure of the given task to obtain a task template.
- Exploit real-world contextual relationships between objects to generate new tasks with that template.

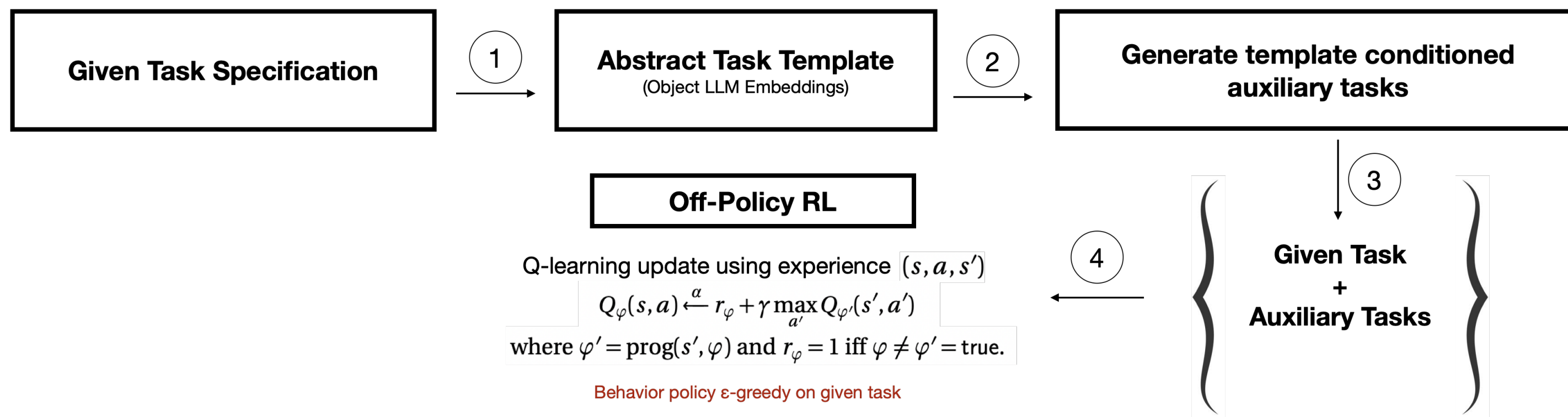
Key takeaway. Policies for auxiliary tasks can be learnt simultaneously without extra environment interaction via **counterfactual experience reasoning** and **off-policy RL**.

Framework for multi-task RL where given a task, an agent can generate auxiliary tasks such that it can perform efficient



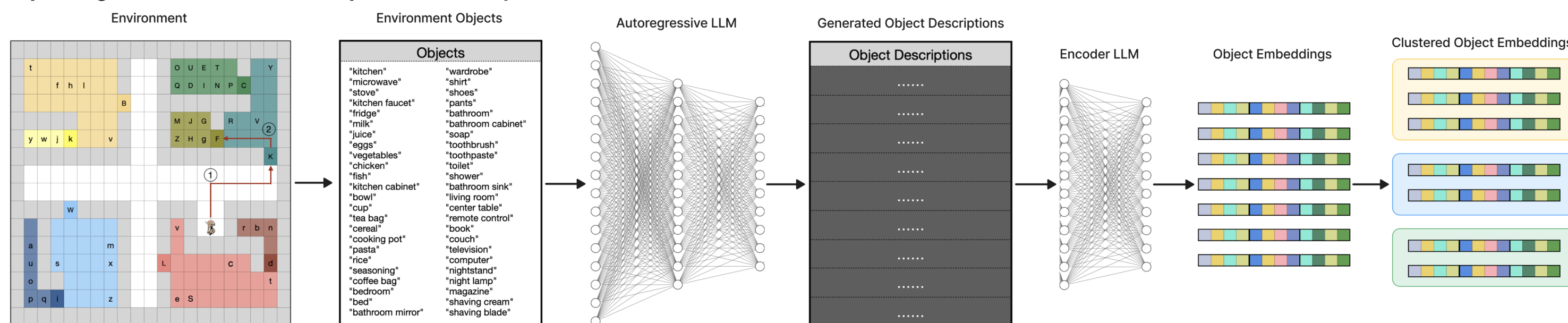
Reuses experience learning given task to simultaneously learn auxiliary tasks

Approach

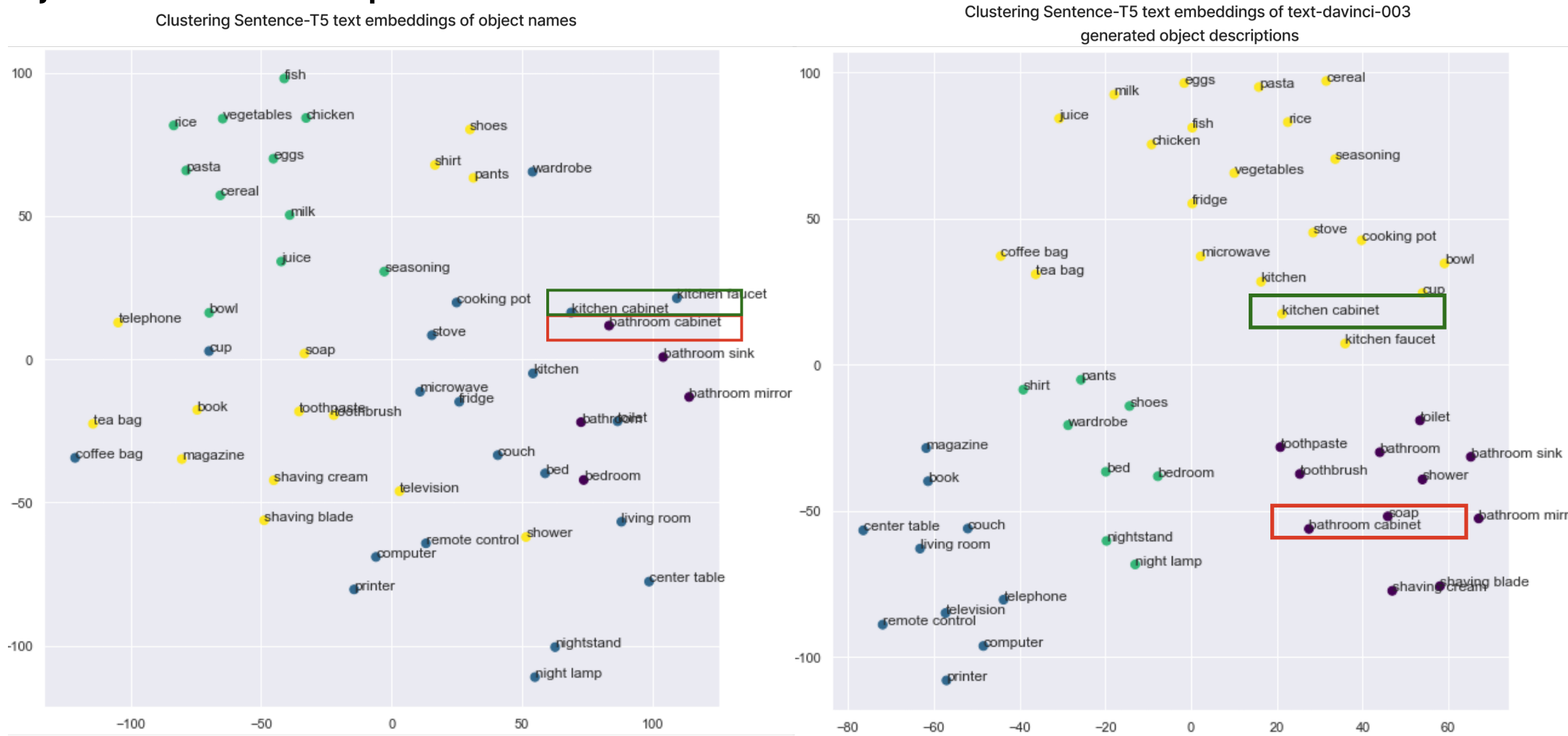


Generating Auxiliary Tasks

Exploiting Contextual Relationships between Objects

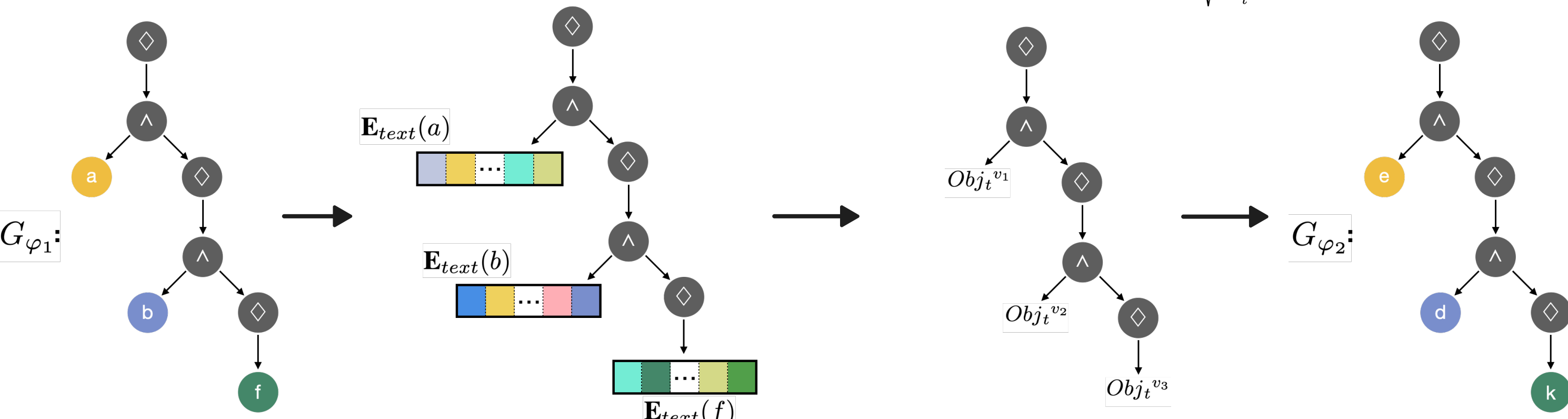


Visualizing Emergent Object Classes in Latent Space



Exploiting Task Structure

$$\varphi_1 : \diamond(a \wedge \diamond(b \wedge \diamond f)) \quad \diamond(\mathbf{E}_{text}(a) \wedge \diamond(\mathbf{E}_{text}(b) \wedge \diamond \mathbf{E}_{text}(f))) \quad Obj_i^{v_i} = \max_{obj_j} [\text{sim}(\mathbf{E}_{text}(G_{\varphi_1}^{v_i}), \mathbf{E}_{text}(obj_j)) + c \sqrt{\frac{\log N_t}{N_t^{obj_j}}}] \quad \varphi_2 : \diamond(e \wedge \diamond(d \wedge \diamond k))$$



Assumptions

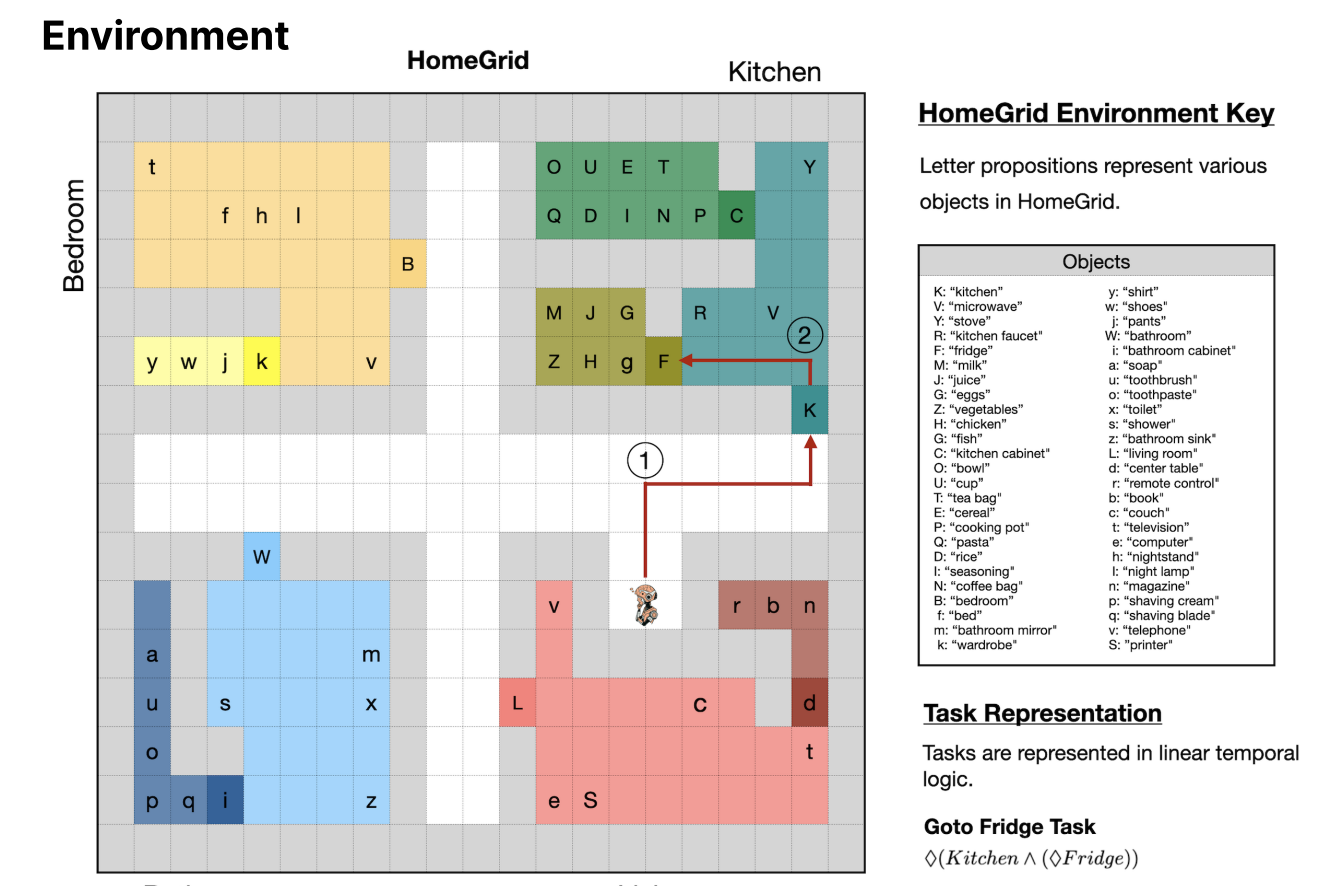
Tasks are represented in linear temporal logic (LTL), an expressive grammar for specifying temporal behavior composed of atomic propositions, logical connectives and the temporal operators:

Next (\circ) Until (U) Always (\square) Eventually (\diamond)

LTL formulae can be progressed given a sequence of proposition truth assignments to determine which parts of the formula have been satisfied by prior states and which parts remain—a useful property for tracking non-Markovian objectives in the reinforcement learning setting.

Example. A sequential task of visiting a Kitchen then a Fridge, where Kitchen and Fridge are Boolean propositions that can be observed is represented as:

$$\diamond(Kitchen \wedge (\diamond Fridge))$$



Experiments & Results

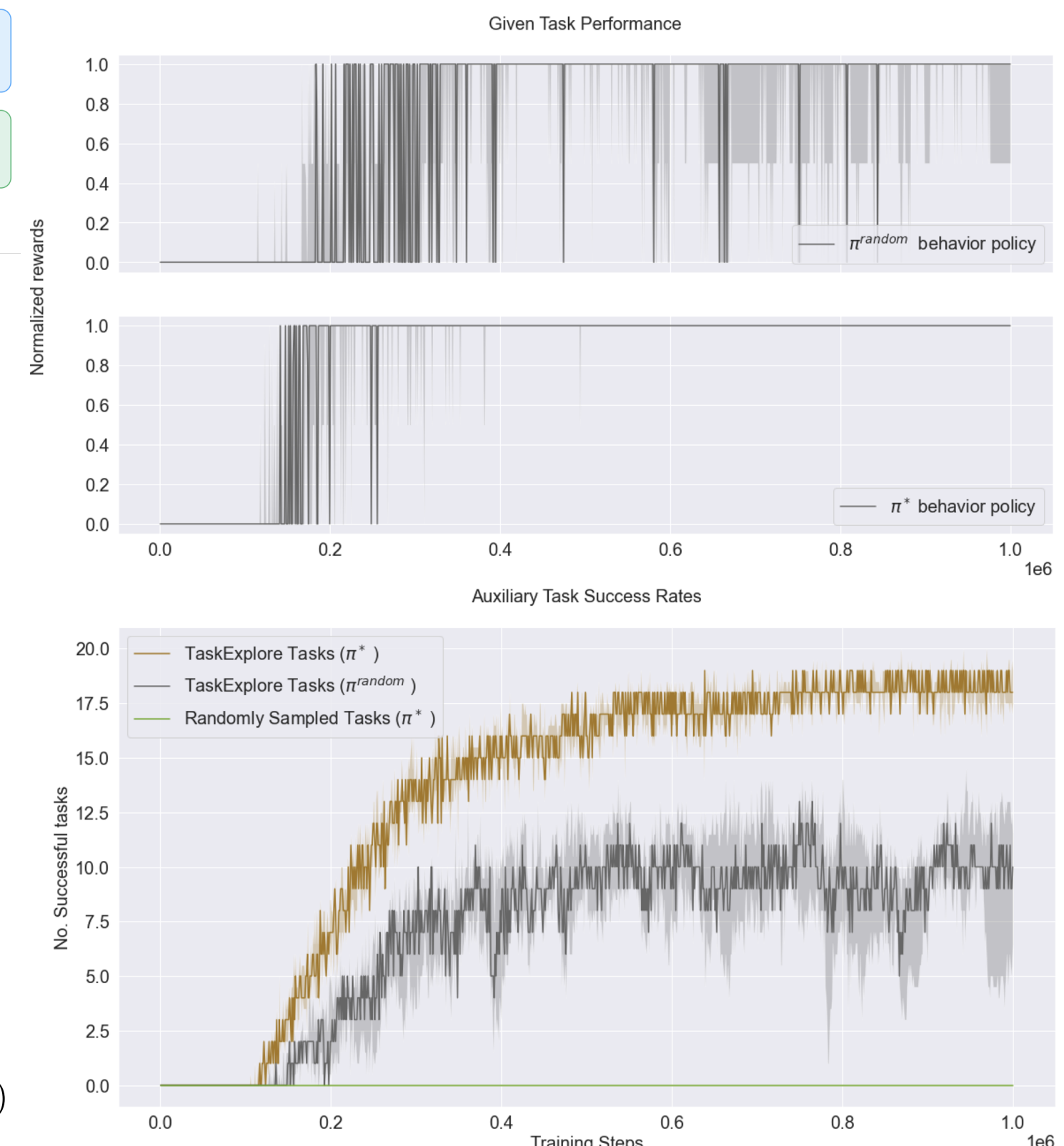
The task specification used in our experiments was a food preparation task where the agent had to go to the **kitchen cabinet**, obtain a **cooking pot**, obtain **seasoning**, then go to the **fridge**, obtain **chicken**, and finally go to the **stove**. In HomeGrid, this task corresponds to visiting the right cells in the correct order. The LTL formula below represents this task using the atomic propositions that represent each of the relevant objects:

$$\diamond(C \wedge \diamond(P \wedge \diamond(I \wedge \diamond(F \wedge \diamond(H \wedge \diamond Y))))))$$

Ours: Given our food preparation task φ_1 we generate 20 auxiliary tasks following our approach. We learn these tasks simultaneously with φ_1 with a behavior policy epsilon greedy on φ_1 , directing exploration towards more relevant experiences for φ_1 .

Baseline 1: To demonstrate that tasks generated by TaskExplore uniquely leverage the directed experience of a single-task curriculum, we repeat the approach described above, replacing the behavior policy with a random one that explores more widely.

Baseline 2: To show that tasks generated by TaskExplore more relevantly benefit from directed exploration experience than the general set of possible tasks, we generate 20 auxiliary tasks by randomly sampling equal length sequential tasks from the set of propositions in our environment, as typically done in prior works. We learn these tasks simultaneously with φ_1 , using a behavior policy epsilon greedy on φ_1 .



Conclusion

This work introduced an approach that allows agents to generate relevant auxiliary tasks that maximally benefits from the directed exploration experience of a single task curriculum.

This approach to auxiliary task generation is particularly valuable in the lifelong learning setting, as agents can generate and solve new tasks from constrained experience datasets.

Modern vision-language models (VLMs) that detect open vocabulary objects in real world environments can be employed with this approach to relax the dependence on predefined sets of object propositions from which tasks can be expressed and labelling functions that map states to proposition truth values.