

Multimodal Embodied Plan Prediction Augmented with Synthetic Embodied Dialogue

Anonymous ACL submission

Abstract

Embodied task completion is a challenge where an agent in a simulated environment must predict environment actions to complete tasks based on natural language instructions and egocentric visual observations. We propose a variant of this problem where the agent predicts actions at a higher level of abstraction called a plan which more directly tests language understanding and reasoning. We show that multimodal transformer models can outperform language-only models for this problem, but fall significantly short of oracle plans. Since collecting human-human dialogues for embodied environments is expensive and time-consuming, we propose a method to synthetically generate such dialogues, which we then use as training data for plan prediction. We demonstrate that multimodal transformer models can attain strong zero shot performance from our synthetic data, outperforming language-only models trained on human-human data.

1 Introduction

Embodied task completion (Shridhar et al., 2020; Padmakumar et al., 2022) is a challenge where an agent in a simulated environment (Kolve et al., 2017; Savva et al., 2019; Chang et al., 2017) is given natural language context in the form of instructions or dialogue and needs to take actions in the environment to complete a desired task, additionally making use of egocentric visual observations. This typically requires the agent to predict actions directly executable in the simulated environment. For example, an action sequence to make coffee could start with actions to move forward 2 steps, turn left and pick up a mug identified by a semantic segmentation mask. In contrast, physical robot systems tend to be more modular with a dedicated component for task planning - composing a sequence of fine-grained motor skills into a more complex task (Chen et al., 2010; Lemaignan et al., 2017; Jiang et al., 2019). In such a system, the

coffee task considered above would likely start by invoking a semantic navigation module to find the mug and a grasping module to pick it up. Some prior work on embodied AI benchmarks suggests that more modular models are capable of outperforming monolithic models (Min et al., 2021; Jia et al., 2022; Zheng et al., 2022; Min et al., 2022). However, these do not evaluate and explore the limitations of individual modules.

In this work, we formulate and explore the problem of task planning for embodied task completion. We improve upon existing plan prediction models and demonstrate the potential for improvement by comparing them with human plans. We adapt the Execution from Dialogue History (EDH) benchmark from the TEACH dataset (Padmakumar et al., 2022) to evaluate models at the level of a plan – a sequence of object interaction actions paired with the type of object on which the action needs to be executed – which are evaluated using task success upon execution with the aid of a heuristic plan execution module. Plan prediction is more challenging in TEACH in comparison to other embodied AI datasets as tasks can be hierarchical and parameterized, and environments are cluttered. Objects may be hidden inside closed receptacles. We evaluate variants of the multimodal Episodic Transformer (E.T.) model, previously used to directly predict low level actions in embodied task completion (Shridhar et al., 2020; Padmakumar et al., 2022) and find that these outperform a fine-tuned language-only baseline.

Data collection for embodied AI tasks involving natural language is expensive and time-consuming to collect (Padmakumar et al., 2022), motivating the need for methods that require less human-human data. We develop a framework for embodied dialogue data generation using agenda-based simulation of semantic dialogue acts (Schatzmann and Young, 2009) paired with a rule-based module for identifying action sequences in the environ-

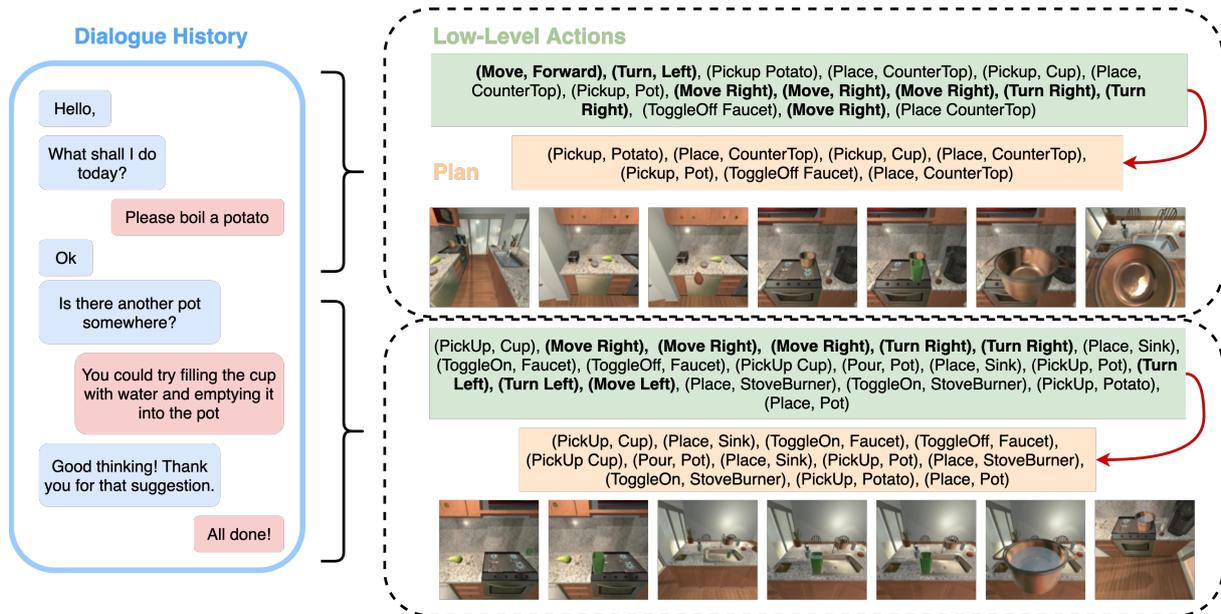


Figure 1: This figure depicts an example EDH instance from the TEACH dataset with modifications converting an action sequence to a plan. Models are trained to predict plans instead of low-level action sequences, and a plan execution module identifies navigation and other adjustment steps to ensure the effective execution of plan actions.

ment to complete tasks. We demonstrate that the E.T.-based models trained only on synthetic data can achieve about 85% of the performance of the same models trained on real data, obtaining a zero-shot success rate of 17.09% which outperforms the full shot success rate of plans generated by the language-only baseline at 10.27%.

To summarize, our contributions include the following:

- We formulate the problem of plan prediction for the TEACH dataset and evaluate a language-only baseline, variants of a multi-modal transformer model (E.T.) and establish oracle performance on this problem.
- We propose a framework for synthesizing embodied dialogues involving both utterances and environment actions to complete a task.
- We demonstrate that the synthetic data generated by our framework results in competitive zero-shot performance in our problem.

2 Task Setup

The TEACH dataset (Padmakumar et al., 2022) is an embodied dialogue dataset consisting of interactions between human annotators role-playing a *Commander* and *Follower* collaborating in a simulated home environment to complete household tasks. Only the *Commander* has access to task

information, and only the *Follower* can take actions in the environment requiring them to communicate to complete the task. An effective *Follower* must engage in dialogue with the *Commander*, obtain relevant information such as details of the task to be completed and locations of objects, and reason about environment actions that can accomplish relevant state changes to make progress in the task. We focus on the EDH benchmark from the TEACH dataset where given some dialogue history, as well as past actions and image observations, the *Follower* must predict subsequent actions in the environment to progress with the task. This is evaluated by comparing state changes in the environment arising from gold and predicted action sequences. We modify the expected prediction from a model to be a **plan**, which we define as a sequence of object interaction actions paired with the object category of the object they are to be executed upon¹. Figure 1 includes an example task of boiling a potato. Note that in plan prediction, the model needs to reason about physical state changes - that the act of boiling requires placing the potato in a container filled with water which is then heated using a stove in the example. Other aspects of execution, such as navigating to required objects and fine-grained position adjustments, can be carried

¹While it is possible to define more abstract plans, we choose this level of abstraction as the training data can be generated automatically from the TEACH EDH instances.

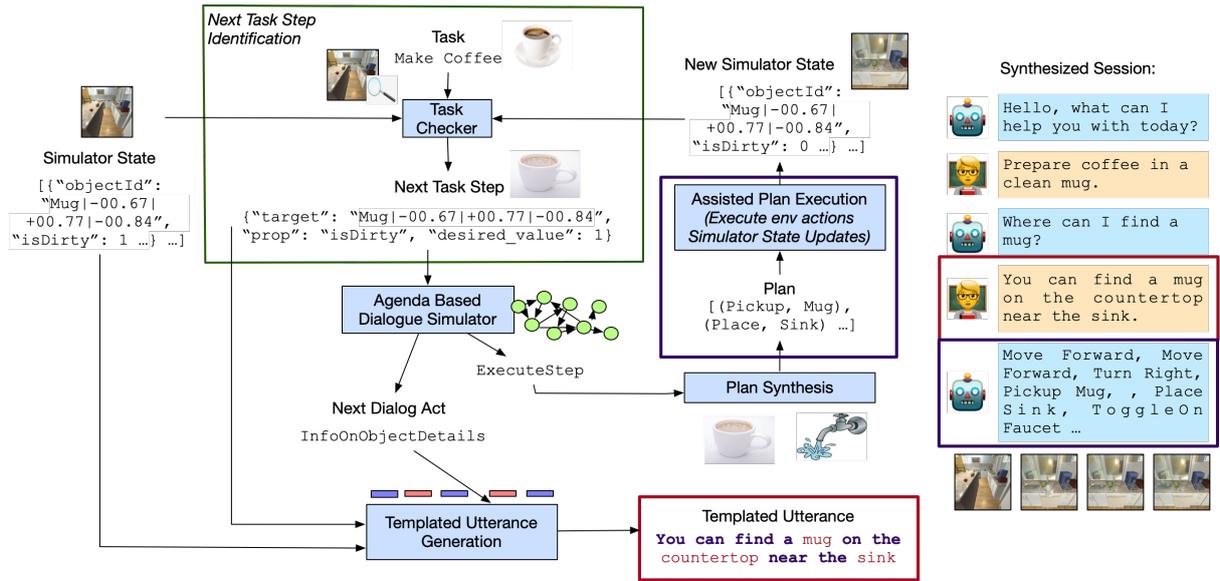


Figure 2: Control Flow for Simulating Synthetic Embodied Dialogues with an example for a task of making coffee.

137 out by a separate execution module, which in our
 138 case, is heuristic but can also be learned. This prob-
 139 lem is nontrivial in datasets such as TEACH where
 140 tasks are parameterized and hence highly variable,
 141 and diverse initial states can add or remove task
 142 steps demonstrated in our experiments by the sig-
 143 nificant gap between model generated and human
 144 plans.

145 During inference, at each time step, a plan pre-
 146 diction model is expected to predict one object
 147 interaction action and the category of the object on
 148 which it is to be executed. This is then executed
 149 by one of two possible plan execution modules de-
 150 scribed in section 4. Execution terminates either
 151 when a model predicts a special Stop action, 30
 152 plan steps that result in simulator execution fail-
 153 ures, or reaches a limit of 100 plan steps. A plan
 154 step may fail execution for a variety of reasons. It
 155 may be infeasible (e.g., trying to pick up a cabinet),
 156 a prerequisite step may not be completed (e.g., the
 157 Slice action is only feasible if the agent is hold-
 158 ing a knife), or the agent may be poorly positioned
 159 (e.g., too close to the fridge to open it).

160 3 Plan Prediction Models

161 We adapt the Episodic Transformer (E.T.)
 162 model (Pashevich et al., 2021) for plan prediction.
 163 This is a multimodal transformer model which, in
 164 our case, receives the EDH dialogue history as lan-
 165 guage input and egocentric image observations as
 166 visual input and predicts plan steps consisting of
 167 an object interaction action and the type of object

168 the action is to be taken on². We obtain training
 169 data by filtering EDH action sequences to contain
 170 only object interaction actions and train the model
 171 as in Padmakumar et al. (2022), where the model
 172 receives images and plan steps from the EDH action
 173 history as input and predicts the entire output
 174 plan at once. At inference time, the last plan step
 175 predicted is executed, and the input for the next
 176 time step is updated with an image observation af-
 177 ter executing this plan step. We use 3 variants of
 178 this model.

- 179 • **E.T.** : E.T. model as described above.
- 180 • **E.T. Hierarchical**: E.T. is modified to
 181 pass output from the action prediction head as
 182 input to the object prediction head.
- 183 • **E.T. + Mask**: Uses predefined constraints
 184 to determine whether the predicted action is
 185 feasible to execute on the predicted object, and
 186 if not, backs up to the action with the highest
 187 probability that is feasible.

188 4 Plan Execution

189 While we can compare predicted plans with a
 190 ground truth plan using surface metrics such as edit
 191 distance, we believe a stronger test of predicted
 192 plans is executing them in the environment and
 193 measuring task success. To do this, we pair our
 194 models with heuristic plan execution modules:

- 195 • **Direct Plan Execution**: Given a predicted
 196 object interaction action and object type, we

²See Appendix D for a more detailed explanation and a figure

197 use object coordinates from the simulator to
198 identify the closest object of the type³, use the
199 navigation graph to navigate to it and attempt
200 to execute the action.

- **Assisted Plan Execution:** Direct plan execution can fail for various reasons. For example, if the sink is full and something needs to be placed in it, the sink needs to be cleared first. Although we train our plan prediction model with data that should enable it to predict such intermediate steps, we wish to explore whether models perform better if they abstract out such details. To do this, we analyze common execution failure cases and implement a set of heuristics, detailed in Appendix C, to increase their success.

213 In future work, we hope to replace these with mod-
214 els for plan execution.

215 5 Synthetic Dialogue Generation

216 Collecting embodied dialogue examples is ex-
217 pensive and time-consuming (Padmakumar et al.,
218 2022). However, it would be desirable for embod-
219 ied agents to adapt to new tasks without requiring
220 human interaction data. Ideally, given a task defini-
221 tion, we would like to be able to bootstrap a model
222 for the task, which can then be further refined using
223 techniques such as reinforcement learning. In this
224 work, we propose a method to generate synthetic
225 embodied dialogues to train an initial model with-
226 out human interaction data. Our process for simu-
227 lating synthetic embodied dialogues is outlined in
228 Figure 2 and in the following sections. Additional
229 details are included in Appendix G. We plan to
230 release our synthetic data for future research.

231 5.1 Next Task Step Identification

232 Our dialogue simulation process involves break-
233 ing down a task into task steps corresponding to a
234 single desired object state change, around which di-
235 alogue utterances and environment actions are con-
236 structed. Tasks in the TEACH dataset are defined as
237 sets of object properties that need to be satisfied for
238 the task to be considered complete (Padmakumar
239 et al., 2022). The public TEACH simulation wrap-
240 per also includes a task checker that, when given a
241 task definition and the current state of the simulator
242 can provide pending object state changes that need

³Note that using the closest object is a heuristic and can fail. In our experiments, we evaluate two oracles to quantify the limitations of this.

243 to be accomplished for the task to be considered
244 complete. For example, Figure 2 demonstrates the
245 synthesis of a dialogue session related to making
246 coffee, which requires a mug in the environment to
247 be clean and filled with coffee. As the agent acts in
248 the environment, the task checker examines mugs
249 in the environment and identifies the one closest to
250 completion. The task checker can indicate to the
251 agent the object state changes that still need to be
252 accomplished on this mug; for example, in Figure
253 2, it identifies a mug that needs to be cleaned. This
254 is a single Task Step that can be used as a focus
255 for dialogue utterances and environment actions.
256 Once this Task Step is completed, the simulation
257 process proceeds to the next Task Step of filling
258 coffee, after which the task checker will indicate
259 task completion, ending the dialogue simulation.

260 5.2 Agenda-Based Dialogue Simulation

261 Given the next Task Step, we build a semantic
262 outline for a snippet of synthetic dialogue related
263 to this Task Step that includes dialogue acts ex-
264 changing information related to this Task Step
265 and predicting a special action ExecuteStep that
266 indicates a transition to predicting environment ac-
267 tions that accomplish this Task Step⁴.

268 We do this by building an **agenda-based dia-**
269 **logue simulator** (Schatzmann and Young, 2009)
270 over dialogue acts relevant to the TEACH dataset
271 combined with the ExecuteStep action. We use
272 a subset of the dialogue acts annotated for the
273 TEACH dataset in Gella et al. (2022), focusing
274 on requesting and receiving instructions related to
275 the task and how to complete it, as well as the lo-
276 cations of objects. Our agenda-based simulator
277 consists of 9 dialogue states, each computed as
278 a boolean function of 8 binary dialogue features.
279 We pre-define probabilities for sampling dialogue
280 acts, ExecuteStep, and DoNothing actions in each
281 state. We then generate a dialogue session by sam-
282 pling a sequence of actions, expanding each dia-
283 logue act into an utterance using templates filled
284 in with task and simulator information, and each
285 ExecuteStep into a sequence of environment ac-
286 tions as described in section 5.3. More details of
287 this process are included in Appendix G.

⁴Note that human-human dialogue in TEACH is much more free-form and we hope to achieve more versatility in future work.

5.3 Plan Synthesis

When we predict that the session should transition from dialogue to environment actions, we use a rule-based system to identify an action sequence in the environment that is likely to accomplish the next Task Step. We hard code plans for each possible object state change, detailed in Appendix G. For example, if the object state change requires an object to be cleaned, the plan will involve moving the object to the sink, turning on and off the tap, picking it up, and pouring out water accumulated from cleaning. These hard-coded plans do not account for handling difficulties arising from the current state of the environment, for example, clearing out the sink if it is too full to place the object to be cleaned. Hence, we execute these synthesized plans using Assisted Execution (Section 4) to improve our success rate at completing task steps using these hard-coded plans.

6 Experiments

6.1 Plan Prediction Models

We evaluate our proposed plan prediction models on the EDH task of the TEACH dataset (Padmakumar et al., 2022). We experiment with each of the models in section 3 with each execution method in section 4. Additionally, we evaluate the following baseline and oracle conditions ⁵:

Baseline: Our baseline is a language-only BART model (Lewis et al., 2020), finetuned to take in the EDH dialogue history and predicts the entire plan as a sequence of tokens that are post-processed for validity and executed as in (Gella et al., 2022).

Oracle: As an upper bound to the success rate obtainable with each of our plan execution methods, we obtain oracle plans using the ground truth actions present in the EDH instance. We filter these action sequences retaining only object interaction steps and converting object IDs to object types to match the plan representation used by our models.

Oracle with Object IDs (CorefOracle): To further test the limitations of our plan representation combined with the heuristic of selecting the closest object of a particular type, we produce plans from

⁵We do not compare to TEACH EDH models in this paper (Padmakumar et al., 2022) as our execution methods access information that the TEACH baseline models are not allowed to access. We do this to ensure that we are evaluating only the process of plan prediction without additional complications arising from navigation and simulator behavior

human action sequences containing object IDs instead of object types to avoid ambiguity during plan execution.

We additionally include our best zero-shot model and our best model trained on both real and synthetic data. These models use synthetic data generated according to the method outlined in Section 5 using the initial states corresponding to the TEACH train set. The zero-shot model is trained only on synthetic data, and the data-augmented model is trained on a combination of real and synthetic data.

We evaluate models based on the success rate (SR) and goal condition success rate (GC) as defined in the original TEACH paper (Padmakumar et al., 2022). Success rate measures the fraction of EDH instances for which predicted plans produced all expected object state changes, and GC measures the fraction of such object state changes across all instances that were calculated. Since the TEACH test set is not public, we follow the standard protocol proposed in the TEACH codebase ⁶ of using a standardized division of the original validation sets into validation and test sets called the divided validation and divided test sets, each of which are further divided into *Seen* and *Unseen* splits.

We present our results in Table 1. For a subset of these conditions, we train and perform inference with 3 random seeds and perform 2-sided Welch t-tests. Allowing for Bonferroni corrections over 4 tests, we find that E.T. + Mask is trending to be significantly better than the baseline with $p = 0.0381$ on the *divided_val_seen* split and $p = 0.0164$ on the *divided_test_seen* split. We did not find any statistically significant difference between the E.T. Hierarchical and E.T. + Mask models ⁷.

Since oracle plans do not obtain a 100% success rate, we observe the limitations of our plan execution method, which can only handle 78.43% of unseen test instances even with object coreference resolved (CorefOracle). We believe this is due to the difficulty in obtaining perfect positioning and placement even with ground truth simulator information, further supported by the gap between direct and assisted execution of oracle plans. We additionally see that there is considerable scope

⁶<https://github.com/alexa/teach>

⁷We did not perform statistical comparisons across all pairs of conditions as it is expensive and time-consuming to run inference with enough random seeds to allow for Bonferroni corrections as the number of tests grows.

Model	Execution	EDH Plan Divided Val Split				EDH Plan Divided Test Split			
		<i>Seen</i>		<i>Unseen</i>		<i>Seen</i>		<i>Unseen</i>	
		SR	GC	SR	GC	SR	GC	SR	GC
Baseline	Direct	11.26	13.67	7.51	11.03	7.19	9.62	8.87	9.54
	Assisted	11.92	17.27	8.91	12.19	9.80	12.30	10.27	12.31
E.T.	Direct	12.91	16.32	15.58	16.20	15.03	19.52	16.62	15.61
	Assisted	15.89	20.57	18.74	22.36	16.67	19.96	19.98	27.13
E.T. Hierarchical	Direct	14.24	15.67	16.23	17.27	14.71	17.97	17.27	20.30
	Assisted	18.21	20.45	18.09	24.53	17.97	23.67	19.70	25.82
E.T. + Mask	Direct	15.23	22.51	17.81	18.29	16.34	23.84	17.46	18.96
	Assisted	18.87	28.99	19.57	27.64	18.95	26.35	20.07	28.33
Best zero shot	Assisted	18.87	17.52	16.23	19.30	15.36	17.17	17.09	16.82
Best Augmented	Assisted	18.87	26.90	19.48	30.41	17.32	26.52	22.32	34.30
Oracle	Direct	61.92	63.64	55.57	58.48	54.58	53.58	56.77	58.01
	Assisted	68.87	72.13	61.97	63.07	61.44	62.49	63.21	64.87
CorefOracle	Direct	77.81	83.03	70.87	71.87	75.82	79.50	71.90	74.34
	Assisted	80.13	84.50	74.58	77.31	78.43	80.92	76.94	78.30

Table 1: Success rate (SR) and Goal Condition Success Rate (GC) of different models combined with different execution methods on the TEACH EDH task. Oracle performances are separated as upper bounds on the task. Best performance results are bolded for each metric and split in the specific execution method.

for improvement from resolving ambiguity related to which object instance to manipulate, which accounts for an improvement of about 17% between Oracle and CorefOracle.

We observe that while the E.T., E.T. Hierarchical and E.T. + Mask models substantially improve over the baseline, there is also a considerable gap between them and Oracle which uses the same plan representation, which demonstrates that there is considerable scope for improvement in understanding the details of the task to be completed from dialogue, and reasoning about actions to take to achieve the corresponding state changes. Qualitatively, we find that the main benefits of multimodal input are in breaking down complex tasks such as watering a plant, for which detailed steps are rarely provided in the dialogue, and in identifying how much of the task has already been completed. Failures of the E.T. models mainly arise from not learning when to stop, which is a limitation of the current inference procedure. Other causes of failure include performing unrelated manipulations on easily visible objects, or ignoring small objects in favor of larger ones.

On comparing models trained on real data with synthetic data, we find that the zero-shot models perform surprisingly well, outperforming the base-

line trained on real data, and approaching the performance of the models which have the same architecture but are trained on real data. This suggests that for this application when generalizing to new tasks, it might be reasonable to train a model purely on synthetic data and expect reasonable performance from interaction with human users.

6.2 Zero Shot Training Ablation

We perform further ablations to identify how zero-shot model performance varies according to data size in Table 2. While we see a trend towards improvement in performance with increasing data size, and the best results are obtained at higher data sizes, the improvement is not perfectly consistent with data size. We believe this is due to fundamental limitations of the E.T. model capacity and in future work, we hope to experiment with stronger models.

Note that we did not perform another hyperparameter tuning for models in Table 2 besides the changes in training data. The best condition, E.T. Hierarchical with a synthetic training set of size 3x as large as the human-human training set from Table 2 has been reported in Table 1 for comparison with models trained on human-human data. We find that this model trained purely on

Model	Train Set Size	EDH Plan Divided Val Split				EDH Plan Divided Test Split			
		<i>Seen</i>		<i>Unseen</i>		<i>Seen</i>		<i>Unseen</i>	
		SR	GC	SR	GC	SR	GC	SR	GC
E.T.	1x	12.25	9.50	13.17	9.33	14.05	11.04	13.82	9.36
E.T.	2x	15.89	13.95	14.47	13.58	14.38	12.03	15.97	11.85
E.T.	3x	17.88	11.94	14.84	10.61	13.07	12.14	16.06	11.22
E.T.	4x	16.89	11.28	13.45	13.07	14.05	11.26	14.94	13.54
E.T. Hierarchical	1x	15.56	10.15	13.36	10.77	11.76	9.73	12.32	9.74
E.T. Hierarchical	2x	15.89	11.16	14.19	16.68	11.44	13.50	14.47	18.53
E.T. Hierarchical	3x	18.87	17.52	16.23	19.30	15.36	17.17	17.09	16.82
E.T. Hierarchical	4x	18.21	10.45	14.75	10.14	14.05	11.97	15.69	11.21
E.T. + Mask	1x	13.58	9.98	13.82	10.90	13.73	12.30	13.35	9.83
E.T. + Mask	2x	14.24	16.39	15.49	16.08	14.38	12.41	15.59	14.94
E.T. + Mask	3x	16.89	12.05	14.10	12.45	12.09	11.21	15.03	12.09
E.T. + Mask	4x	15.23	14.37	14.66	15.20	14.38	13.45	14.29	14.98

Table 2: We explore how zero shot model performance varies with the size of the synthetic training set (as a proportion of the size of the human-human training data). This table reports success rate (SR) and goal condition success rate (GC) on the EDH divided test split with assisted execution.

Model	Synthetic Data Size	EDH Plan Divided Val Split				EDH Plan Divided Test Split			
		<i>Seen</i>		<i>Unseen</i>		<i>Seen</i>		<i>Unseen</i>	
		SR	GC	SR	GC	SR	GC	SR	GC
E.T.	1x	18.21	17.52	17.90	29.12	17.65	23.67	21.38	31.63
E.T.	2x	17.55	25.00	19.02	24.99	18.63	23.51	19.70	24.60
E.T.	4x	17.22	26.84	18.37	26.12	17.97	25.26	19.79	26.01
E.T. Hierarchical	1x	16.89	24.70	18.74	25.73	18.30	23.46	20.45	25.65
E.T. Hierarchical	2x	18.54	19.48	18.74	23.55	20.92	28.05	19.61	27.91
E.T. Hierarchical	4x	14.90	20.31	17.90	24.22	17.97	26.95	20.17	24.30
E.T. + Mask	1x	18.87	26.90	19.48	30.41	17.32	26.52	22.32	34.30
E.T. + Mask	2x	16.89	32.84	19.02	30.66	18.63	27.99	20.45	29.98
E.T. + Mask	4x	18.21	27.38	18.65	28.01	20.26	27.88	19.79	27.28

Table 3: We explore how synthetic data size impacts model performance. Models here are trained on a combination of the human-human TEACH training data and synthetic data of the size indicated in column Synthetic Data Size (as a proportion of the size of the human-human training data). This table reports success rate (SR) and goal condition success rate (GC) on the EDH divided test split with assisted execution.

synthetic data obtains a success rate of 17.09% on the divided unseen test split, outperforming the language-only baseline trained on real data at 8.87%, and approaches close to the performance of E.T. Hierarchical trained on human-human

data at 19.70%.

6.3 Data Augmentation Training Ablation

In Table 3, we ablate different sizes of synthetic data when used in data augmentation. We find that when both real and synthetic data are included,

431
432
433
434
435

436
437
438
439
440

larger sizes of the synthetic training set are less beneficial than when trained only on synthetic data. The best condition E.T. + Mask with a combination of human-human and synthetic data of equal size has been included in Table 1 as "Best Augmented". We find that at 22.32% on the divided unseen test split, this slightly outperforms the same model condition trained on human-human data at 20.07%, and is much stronger than the best condition using only synthetic data at 17.09%.

7 Related Work

Task Planning: Interactive systems on physical robots typically have a modular structure in which task planning plays a significant role (Chen et al., 2010; Khandelwal et al., 2017; Peshkin et al., 2001). In simulated environments, Logeswaran et al. (2022) propose a language-only finetuned GPT-2 model for task planning on ALFRED. Some end-to-end ALFRED models also have task planning as a component (Min et al., 2021; Jia et al., 2022). However, this is a much simpler dataset where task planning can be cast as a 7-way classification problem. Prior work has also explored language-only task planning using finetuned BART models in TEACH (Gella et al., 2022; Zheng et al., 2022), which we compare to as a baseline.

Dialogue Simulation: User simulation in the dialogue community originally consisted of rule-based systems designed using linguistic knowledge to enable fine-tuning dialogue systems to individual user preferences and subsequently evolved into trainable probabilistic models that can be used to bootstrap a dialogue system in the initial development phase and further finetune it through reinforcement learning (Schatzmann et al., 2006; Young et al., 2013). A common method for building user simulators is agenda-based simulation (Schatzmann and Young, 2009) which uses a predefined set of transition probabilities between dialogue acts in combination with goal information to sample subsequent dialogue acts. This has been used to bootstrap a range of dialogue models ranging from probabilistic POMDP models (Schatzmann et al., 2007) and text-to-SQL models (Liu et al., 2022) to hierarchical deep reinforcement learning methods (Peng et al., 2017). In this work, we augment a standard agenda-based simulator with an additional intent to determine transitions to acting in the environment in order to generate situated dialogues. Another popular paradigm for dialogue simulation is to de-

velop two models - one for the user and one for the agent side and train them simultaneously using reinforcement learning (Liu and Lane, 2017; Shah et al., 2018). This has also been used for some multimodal dialogue domains (Das et al., 2017) but we choose not to adapt it in our domain as the time to generate a single dialogue is higher in situated applications due to the latency from executing environment actions in the simulator.

8 Conclusions and Future Work

We develop a model for multi-modal plan prediction for the TEACH dataset using the Episodic Transformer architecture and evaluate end-to-end performance on the TEACH EDH task in conjunction with heuristic plan execution modules. We additionally experiment with training this model using only synthetic data generated using an agenda-based dialogue simulator combined with an environment action generator that uses a combination of rules and simulator information to identify sequences of actions in the environment that can make progress with the task. We find that our E.T. plan prediction models outperform a BART baseline, even when BART is finetuned on human-human embodied dialogue data but E.T. is finetuned only on synthetic data, suggesting that our dialogue simulation approach is a viable alternative to expensive data collection for bootstrapping an embodied task completion model on new tasks. We additionally find that there remains a considerable gap in performance between models and humans in plan prediction, that cannot be easily closed by techniques such as data augmentation. In future work, we plan to explore other approaches to plan prediction such as exploration of the environment with reinforcement learning and making use of very large language models.

9 Limitations

In this paper, we explore the problem of plan prediction for embodied task completion, and conduct our experiments in the TEACH dataset which is set in the AI2-THOR simulator. While we hypothesize that models developed for plan prediction will transfer better to physical robots as they align better with levels of abstraction at which robotics systems are currently implemented, further experimentation is needed to evaluate such transferability. In the short term, such experiments will likely need to work on problems with a simpler action space

as some of the *Follower* actions supported in AI2-THOR such as slicing are not supported in most robots available currently. Additionally, it would be beneficial to test plan prediction models and our dialogue simulation method on similar tasks set in other simulators. This is a direction we plan to pursue in future work. We believe this is beyond the scope of this publication due to the considerable engineering effort involved in adapting models across different embodied AI simulators and task representations.

In this work, we describe a method to generate synthetic dialogues for embodied task completion by augmenting an agenda-based dialogue simulator with modules that break up an embodied AI task into individual object state changes, and rule-based methods to identify actions that complete them. Additionally, we currently manually define transition probabilities between dialogue acts for simulation, which would also need to be modified for a new dataset. While our approach is effective at bootstrapping plan prediction models without any human-human interaction data. The approach requires simulator-specific engineering, particularly when defining heuristics for assisted plan execution. Additionally, for both plan prediction inference and dialogue simulation, every action needs to be executed in the simulator. This results in considerable computing time, as discussed in the appendix. While dialogue simulation does not require a GPU, the plan prediction models do - both for training and inference. We have additionally found that the Episodic Transformer models used in the paper show a noticeable variance in performance when trained with the same hyperparameters but with different random seeds. We attempt to strengthen our conclusions by training models with multiple seeds for statistical analysis where it is beneficial, but we would also like to highlight the development of models with less variance as an important direction for future research.

Finally, our work is currently limited to English as we are not familiar with datasets in other languages that provide language instructions for tasks that require complex reasoning over multi-step action sequences.

10 Ethics Statement

This work is part of a broader tradition of building natural language interfaces to control various devices. Natural language interfaces such as

language-based search, and intelligent personal assistants provide convenience and have the potential to make various forms of technology ranging from mobile phones and computers, as well as robots or other machines such as ATMs or self-checkout counters more accessible and less intimidating to users who are unfamiliar or uncomfortable with other interfaces on such devices such as command shells, button-based interfaces or changing visual user interfaces. Spoken language interfaces can also be used to make such devices more accessible for the visually impaired or users who have difficulty with fine motor control.

User trust in such interfaces is important. Depending on the circumstances, therefore, some considerations to keep in mind are: (1) whether the collection of personal data benefits the user; (2) whether the collection of personal data is transparent to the user; (3) whether the user understands whether and how they can control the collection of personal data; and (4) whether the user has the ability to elect whether to use such interfaces, opt into or out of the collection of certain personal data, access and update certain personal data, and delete certain personal data. Additionally, safeguards may be required as embodied agents become capable of interacting with arbitrary objects in the world to reduce the likelihood of accidents or malicious misuse.

In this work, we also experiment with the use of simulated embodied dialogue data to train models. On the one hand, the use of simulated data can limit the collection of personal data. On the other hand, simulators may not be designed to represent the full range of user behavior, and hence may perform better for some groups of users than others.

References

- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, et al. 2022. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*.
- Ahmed Al-Moadhen, Renxi Qiu, Michael Packianather, Ze Ji, and Rossi Setchi. 2013. [Integrating robot task planner with common-sense knowledge base to improve the efficiency of planning](#). *Procedia Computer Science*, 22:211–220. 17th International Conference in Knowledge Based and Intelligent Information and Engineering Systems - KES2013.
- Peter Anderson, Angel Chang, Devendra Singh Chap-

641	lot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, et al. 2018a. On evaluation of embodied navigation agents. <i>arXiv preprint arXiv:1807.06757</i> .	698
642		699
643		700
644		701
645		702
646	Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton van den Hengel. 2018b. Vision-and-Language Navigation: Interpreting visually-grounded navigation instructions in real environments. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i> .	703
647		
648		704
649		705
650		706
651		707
652		
653	Sugato Bagchi, Gautam Biswas, and Kazuhiko Kawamura. 1996. Interactive task planning under uncertainty and goal changes. <i>Robotics and Autonomous Systems</i> , 18(1-2):157–167.	
654		
655		
656		
657	Dhruv Batra, Aaron Gokaslan, Aniruddha Kembhavi, Oleksandr Maksymets, Roozbeh Mottaghi, Manolis Savva, Alexander Toshev, and Erik Wijmans. 2020. Objectnav revisited: On evaluation of embodied agents navigating to objects. <i>arXiv preprint arXiv:2006.13171</i> .	
658		
659		
660		
661		
662		
663	Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. 2017. Matterport3d: Learning from rgb-d data in indoor environments. <i>International Conference on 3D Vision (3DV)</i> .	
664		
665		
666		
667		
668		
669	Howard Chen, Alane Suhr, Dipendra Misra, Noah Snaveley, and Yoav Artzi. 2019. Touchdown: Natural language navigation and spatial reasoning in visual street environments. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 12538–12547.	
670		
671		
672		
673		
674		
675	Xiaoping Chen, Jianmin Ji, Jiehui Jiang, Guoqiang Jin, Feng Wang, and Jiongkun Xie. 2010. Developing high-level cognitive functions for service robots. In <i>AAMAS</i> , volume 10, pages 989–996.	
676		
677		
678		
679	Abhishek Das, Satwik Kottur, José MF Moura, Stefan Lee, and Dhruv Batra. 2017. Learning cooperative visual dialog agents with deep reinforcement learning. In <i>Proceedings of the IEEE international conference on computer vision</i> , pages 2951–2960.	
680		
681		
682		
683		
684	Richard E Fikes and Nils J Nilsson. 1971. Strips: A new approach to the application of theorem proving to problem solving. <i>Artificial intelligence</i> , 2(3-4):189–208.	
685		
686		
687		
688	Chuang Gan, Jeremy Schwartz, Seth Alter, Martin Schrimpf, James Traer, Julian De Freitas, Jonas Kubilius, Abhishek Bhandwaldar, Nick Haber, Megumi Sano, et al. 2020. Threedworld: A platform for interactive multi-modal physical simulation. <i>arXiv preprint arXiv:2007.04954</i> .	
689		
690		
691		
692		
693		
694	Michael Gelfond and Yulia Kahl. 2014. <i>Knowledge representation, reasoning, and the design of intelligent agents: The answer-set programming approach</i> . Cambridge University Press.	
695		
696		
697		
	Spandana Gella, Aishwarya Padmakumar, Patrick Lange, and Dilek Hakkani-Tur. 2022. Dialog Acts for Task-Driven Embodied Agents. In <i>Proceedings of the 23rd Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDial)</i> , pages 111–123.	698
		699
		700
		701
		702
		703
	N Gopalan, E Rosen, GD Konidaris, and S Tellex. 2020. Simultaneously learning transferable symbols and language groundings from perceptual data for instruction following. <i>Robotics: Science and Systems XVI</i> .	704
		705
		706
		707
	Peter E Hart, Nils J Nilsson, and Bertram Raphael. 1968. A formal basis for the heuristic determination of minimum cost paths. <i>IEEE transactions on Systems Science and Cybernetics</i> , 4(2):100–107.	708
		709
		710
		711
	Malte Helmert. 2004. A planning heuristic based on causal graph analysis. In <i>ICAPS</i> , volume 16, pages 161–170.	712
		713
		714
	Zhiwei Jia, Kaixiang Lin, Yizhou Zhao, Qiaozhi Gao, Govind Thattai, and Gaurav Sukhatme. 2022. Learning to act with affordance-aware multimodal neural slam. <i>arXiv preprint arXiv:2201.09862</i> .	715
		716
		717
		718
	Yuqian Jiang, Nick Walker, Justin Hart, and Peter Stone. 2019. Open-world reasoning for service robots. In <i>Proceedings of the International Conference on Automated Planning and Scheduling</i> , volume 29, pages 725–733.	719
		720
		721
		722
		723
	Yash Kant, Arun Ramachandran, Sriram Yenamandra, Igor Gilitschenski, Dhruv Batra, Andrew Szot, and Harsh Agrawal. 2022. Housekeep: Tidying virtual households using commonsense reasoning. <i>arXiv preprint arXiv:2205.10712</i> .	724
		725
		726
		727
		728
	Piyush Khandelwal, Shiqi Zhang, Jivko Sinapov, Matteo Leonetti, Jesse Thomason, Fangkai Yang, Ilaria Gori, Maxwell Svetlik, Priyanka Khante, Vladimir Lifschitz, et al. 2017. Bwibots: A platform for bridging the gap between ai and human–robot interaction research. <i>The International Journal of Robotics Research</i> , 36(5-7):635–659.	729
		730
		731
		732
		733
		734
		735
	Hyoungun Kim, Abhaysinh Zala, Graham Burri, Hao Tan, and Mohit Bansal. 2020. Arramon: A joint navigation-assembly instruction interpretation task in dynamic environments. In <i>Findings of EMNLP</i> .	736
		737
		738
		739
	Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. 2017. Ai2-thor: An interactive 3d environment for visual ai. <i>arXiv preprint arXiv:1712.05474</i> .	740
		741
		742
		743
		744
	George Konidaris, Leslie Pack Kaelbling, and Tomas Lozano-Perez. 2018. From skills to symbols: Learning symbolic representations for abstract high-level planning. <i>J. Artif. Int. Res.</i> , 61(1):215–289.	745
		746
		747
		748
	Séverin Lemaignan, Mathieu Warnier, E Akin Sisbot, Aurélie Clodic, and Rachid Alami. 2017. Artificial cognition for social human–robot interaction: An implementation. <i>Artificial Intelligence</i> , 247:45–69.	749
		750
		751
		752

753	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880.	809
754		810
755		811
756		812
757		813
758		
759	Nir Lipovetzky. 2014. <i>Structure and inference in classical planning</i> . Lulu. com.	
760		
761		
762		
763	Bing Liu and Ian Lane. 2017. Iterative policy learning in end-to-end trainable task-oriented neural dialog models. In <i>2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 482–489. IEEE.	
764		
765		
766		
767		
768	Qi Liu, Zihuiwen Ye, Tao Yu, Phil Blunsom, and Linfeng Song. 2022. Augmenting multi-turn text-to-sql datasets with self-play. <i>arXiv preprint arXiv:2210.12096</i> .	
769		
770		
771		
772	Lajanugen Logeswaran, Yao Fu, Moontae Lee, and Honglak Lee. 2022. Few-shot subgoal planning with language models. <i>arXiv preprint arXiv:2205.14288</i> .	
773		
774		
775	Drew V McDermott. 1996. A heuristic estimator for means-ends analysis in planning. In <i>AIPS</i> , volume 96, pages 142–149.	
776		
777		
778	So Yeon Min, Devendra Singh Chaplot, Pradeep Kumar Ravikumar, Yonatan Bisk, and Ruslan Salakhutdinov. 2021. Film: Following instructions in language with modular methods . In <i>International Conference on Learning Representations</i> .	
779		
780		
781		
782		
783	So Yeon Min, Hao Zhu, Ruslan Salakhutdinov, and Yonatan Bisk. 2022. Don’t copy the teacher: Data and model challenges in embodied dialogue. <i>arXiv preprint arXiv:2210.04443</i> .	
784		
785		
786		
787	Anjali Narayan-Chen, Prashant Jayannavar, and Julia Hockenmaier. 2019. Collaborative dialogue in minecraft. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 5405–5415.	
788		
789		
790		
791		
792	Aishwarya Padmakumar, Jesse Thomason, Ayush Shrivastava, Patrick Lange, Anjali Narayan-Chen, Spanandana Gella, Robinson Piramuthu, Gokhan Tur, and Dilek Hakkani-Tur. 2022. Teach: Task-driven embodied agents that chat . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pages 2017–2025.	
793		
794		
795		
796		
797		
798		
799	Alexander Pashevich, Cordelia Schmid, and Chen Sun. 2021. Episodic transformer for vision-and-language navigation. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 15942–15952.	
800		
801		
802		
803		
804	Baolin Peng, Xiujun Li, Lihong Li, Jianfeng Gao, Asli Celikyilmaz, Sungjin Lee, and Kam-Fai Wong. 2017. Composite task-completion dialogue policy learning via hierarchical deep reinforcement learning. <i>arXiv preprint arXiv:1704.03084</i> .	
805		
806		
807		
808		
	Michael A Peshkin, J Edward Colgate, Wit Wannasuphoprasit, Carl A Moore, R Brent Gillespie, and Prasad Akella. 2001. Cobot architecture. <i>IEEE Transactions on Robotics and Automation</i> , 17(4):377–390.	814
		815
		816
		817
		818
		819
	Caroline Ponzoni Carvalho Chanel, Alexandre Albore, Jorrit T’hoft, Charles Lesire, and Florent Teichteil-Königsbuch. 2019. Ample: an anytime planning and execution framework for dynamic and uncertain problems in robotics. <i>Autonomous Robots</i> , 43(1):37–62.	
	Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba. 2018. Virtualhome: Simulating household activities via programs. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 8494–8502.	820
		821
		822
		823
		824
		825
	Silvia Richter and Matthias Westphal. 2010. The lama planner: Guiding cost-based anytime planning with landmarks. <i>Journal of Artificial Intelligence Research</i> , 39:127–177.	826
		827
		828
		829
	Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. 2019. Habitat: A platform for embodied ai research. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 9339–9347.	830
		831
		832
		833
		834
		835
		836
	Jost Schatzmann, Blaise Thomson, Karl Weilhammer, Hui Ye, and Steve Young. 2007. Agenda-based user simulation for bootstrapping a POMDP dialogue system . In <i>Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Companion Volume, Short Papers</i> , pages 149–152, Rochester, New York. Association for Computational Linguistics.	837
		838
		839
		840
		841
		842
		843
		844
		845
	Jost Schatzmann, Karl Weilhammer, Matt Stuttle, and Steve Young. 2006. A survey of statistical user simulation techniques for reinforcement-learning of dialogue management strategies. <i>The knowledge engineering review</i> , 21(2):97–126.	846
		847
		848
		849
		850
	Jost Schatzmann and Steve Young. 2009. The hidden agenda user simulation model. <i>IEEE transactions on audio, speech, and language processing</i> , 17(4):733–747.	851
		852
		853
		854
	Pararth Shah, Dilek Hakkani-Tür, Gokhan Tür, Abhinav Rastogi, Ankur Bapna, Neha Nayak, and Larry Heck. 2018. Building a conversational agent overnight with dialogue self-play. <i>arXiv preprint arXiv:1801.04871</i> .	855
		856
		857
		858
	Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020. Alfred: A benchmark for interpreting grounded instructions for everyday tasks . In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 10740–10749.	859
		860
		861
		862
		863
		864
		865

866	Alane Suhr, Claudia Yan, Jacob Schluger, Stanley Yu, Hadi Khader, Marwa Mouallem, Iris Zhang, and Yoav Artzi. 2019. Executing instructions in situated collaborative interactions. <i>arXiv preprint arXiv:1910.03655</i> .	believe will be consistent with AI2-THOR and TEACH license terms. We believe our work is consistent with the intended usage of AI2-THOR, which is to further embodied AI research in the household domain.	916 917 918 919 920
871	Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer. 2020. Vision-and-dialog navigation. In <i>Conference on Robot Learning</i> , pages 394–406. PMLR.	• BART model: We finetune BART (Lewis et al., 2020) as a baseline for our task released under the Apache 2.0 license. We believe our usage is consistent with the license and do not intend to redistribute it. We believe that our usage is consistent with the intended usage of BART as a general purpose sequence to sequence language model.	921 922 923 924 925 926 927 928
875	Marc Toussaint and Christian Goerick. 2007. Probabilistic inference for structured planning in robotics. In <i>2007 IEEE/RSJ International Conference on Intelligent Robots and Systems</i> , pages 3068–3073. IEEE.	• Episodic Transformer Model: We use the Episodic Transformer model (Pashevich et al., 2021), which is released under an MIT License, for plan prediction with some architectural modifications. We believe our usage is consistent with the intent of this model, which was designed for a very similar embodied AI application in a similar dataset.	929 930 931 932 933 934 935 936
879	Claudia Yan, Dipendra Misra, Andrew Bennett, Aaron Walsman, Yonatan Bisk, and Yoav Artzi. 2018. <i>Chalet: Cornell house agent learning environment</i> .	The TEACH dataset was manually inspected by the original authors to remove identifying information and offensive utterances (Padmakumar et al., 2022). Since our generated data is templated, we do not believe it is possible for this to contain personally identifying information or offensive utterances. Our synthetic dialogues are set in the TEACH environment and only cover the tasks listed in the TEACH dataset. The dialogues are in English and are mainly intended to cover requesting and informing of task steps (in the form of object state changes), and locations of objects. They support a limited breakage of strict turn-taking by allowing an agent to not generate an utterance at some time steps in the generation procedure.	937 938 939 940 941 942 943 944 945 946 947 948 949 950 951
882	Steve Young, Milica Gašić, Blaise Thomson, and Jason D Williams. 2013. Pomdp-based statistical spoken dialog systems: A review. <i>Proceedings of the IEEE</i> , 101(5):1160–1179.		
886	Kaizhi Zheng, Kaiwen Zhou, Jing Gu, Yue Fan, Jialu Wang, Zonglin Li, Xuehai He, and Xin Eric Wang. 2022. Jarvis: A neuro-symbolic commonsense reasoning framework for conversational embodied agents. <i>arXiv preprint arXiv:2208.13266</i> .		
891	A Licensing and Responsible Use		
892	In this section, we discuss the licensing and usage of existing scientific artifacts in this paper:		
894	• TEACH dataset: The TEACH dataset is released under a CDLA-Sharing V 1.0 license, with images released under Apache 2.0 and code under MIT license. We believe our usage does not violate any of these license terms. We do not redistribute any of these as part of our work as they are publicly available. The TEACH dataset was created to study models for translating natural language instructions and egocentric visual observations into action sequences. Our use case is a subtask of this intended use case. We also believe using the TEACH code to generate synthetic dialogue sessions is consistent with this goal.		
908	• AI2-THOR simulator: The AI2-THOR simulator (Kolve et al., 2017) is necessary for the use of the TEACH dataset and for our process of generating synthetic dialogue sessions. We believe our usage of AI2-THOR does not violate the license. We plan to release our data under CDLA-Sharing V 1.0 license, with images released under Apache 2.0, which we		
915		B Additional Related Work	952
		Task Planning on Physical Robots Task planning has long been a standard component of physical robot architectures (Chen et al., 2010), particularly with general purpose service robots (Khandelwal et al., 2017; Peshkin et al., 2001). Classical task planners include a symbolic representation of the state of the world, a goal, and skills the robot can execute. They are expected to find a sequence of skills that, when executed, will transform the world into the goal state, typically using heuristic search algorithms (Lipovetzky, 2014). Over the years,	953 954 955 956 957 958 959 960 961 962 963

964 research in planning has improved the symbolic
 965 representations used in planners (Fikes and Nils-
 966 son, 1971; McDermott, 1996; Gelfond and Kahl,
 967 2014; Konidaris et al., 2018; Gopalan et al., 2020),
 968 search algorithms (Hart et al., 1968; Helmert, 2004;
 969 Richter and Westphal, 2010) and handling uncer-
 970 tainty via probabilistic methods (Toussaint and
 971 Goerick, 2007; Bagchi et al., 1996; Ponzoni Car-
 972 valho Chanel et al., 2019). More recent work
 973 has focused on expanding beyond fully defined
 974 world representations by expanding to use com-
 975 mon sense (Al-Moadhen et al., 2013) and open
 976 worlds (Jiang et al., 2019). Some interesting efforts
 977 in this direction use Large Language Models to
 978 perform planning (Ahn et al., 2022).

979 **Embodied AI Tasks in Simulation** Simulated
 980 environments (Kolve et al., 2017; Savva et al.,
 981 2019; Puig et al., 2018; Chang et al., 2017; Yan
 982 et al., 2018) have been used over recent years to
 983 explore the efficacy of deep learning methods that
 984 directly use egocentric visual observations instead
 985 of data from expensive sensors. While there is
 986 a challenge in transferring results from simulated
 987 to real environments, simulated environments are
 988 more accessible, less expensive, and allow for the
 989 testing of technologies that may not be sufficiently
 990 safe for use in the real world (Savva et al., 2019).
 991 Additionally, while simulated environments can be
 992 used for tasks that do not require the use of lan-
 993 guage (Anderson et al., 2018a; Batra et al., 2020;
 994 Gan et al., 2020; Kant et al., 2022), they play a par-
 995 ticularly valuable role in developing language un-
 996 derstanding and reasoning capabilities over actions
 997 that are currently difficult for physical robots to
 998 complete, but we hope will become a reality in the
 999 future (Kolve et al., 2017). Much of the work in lan-
 1000 guage understanding for embodied AI happens us-
 1001 ing vision and language navigation, where an agent
 1002 must learn to navigate through a previously un-
 1003 seen environment purely based on natural language
 1004 route instructions (Anderson et al., 2018b; Chen
 1005 et al., 2019; Thomason et al., 2020). Embodied task
 1006 completion additionally requires performing object
 1007 manipulation actions (Shridhar et al., 2020; Pad-
 1008 makumar et al., 2022; Suhr et al., 2019; Kim et al.,
 1009 2020; Narayan-Chen et al., 2019). In this work,
 1010 we specifically focus on the TEACH dataset (Pad-
 1011 makumar et al., 2022) as it involves more complex
 1012 tasks which require nontrivial task planning.

C Assisted Plan Execution 1013

This section outlines the full set of heuristics in-
 1014 volved in assisted plan execution: 1015

- For all actions, if the target object property
 1016 change is already complete, do nothing to
 1017 avoid an execution failure. 1018
- *Pickup*: If the object is inside a receptacle
 1019 (container), open the receptacle. After pickup,
 1020 if a receptacle was opened, close it. 1021
- *Place*: If the target receptacle is in a recep-
 1022 tacle, take it out and place it on the counter
 1023 first (for example, if we need to place some-
 1024 thing on a plate that is inside a drawer). If the
 1025 target receptacle needs to be opened, open it
 1026 and close it after placement (for example, a
 1027 drawer or microwave needs to be opened to
 1028 place something inside). If a placement at-
 1029 tempt fails, try removing the existing contents
 1030 of the receptacle one by one to make more
 1031 space. 1032
- *Open, Close*: Toggle off the target object if
 1033 relevant (for example, microwaves need to be
 1034 turned off to open them). 1035
- *ToggleOn, ToggleOff*: If the target is open,
 1036 close it first (for example, microwaves need
 1037 to be closed to turn them on). 1038
- *Slice*: If the target is in a receptacle, move it
 1039 to the counter first. 1040

1041 Additionally, we also attempt position adjustments
 1042 to increase the chance of success.

D E.T. model 1043

This section discusses the Episodic Transformer
 1044 (E.T.) Model (Pashevich et al., 2021) along with
 1045 our modifications in greater detail. For conve-
 1046 nience, we include a diagram of the model in Fig-
 1047 ure 3. The E.T. model receives language (in our
 1048 case EDH dialogue history) and egocentric image
 1049 observations of size 900 x 900 as input. Visual
 1050 observations are first resized to 224 x 244 and then
 1051 encoded using a visual encoder that is based on a
 1052 Faster R-CNN model trained on 325K frames of
 1053 expert demonstrations from the ALFRED train fold
 1054 (which comprises seen splits of TEACH) and not
 1055 finetuned in any of our experiments. This encoder
 1056 average-pools ResNet features 4 times and adds a
 1057

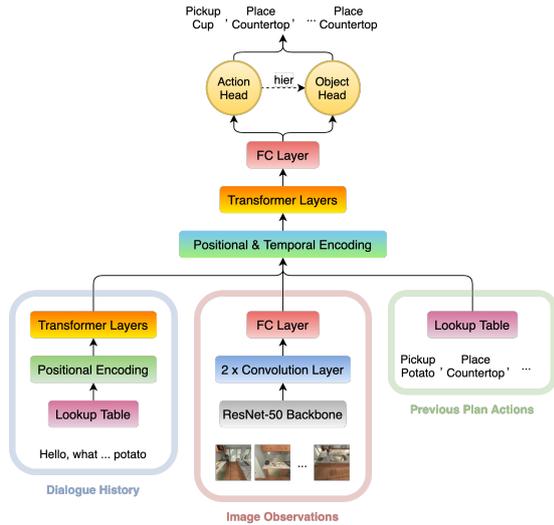


Figure 3: This depicts the architecture for the E.T.-based models. The basic E.T. model does not have a connection between the action and object heads, but E.T. Hierarchical does. There are three main input processing components (dialogue history, image observation, and previous plan actions). These are then input into the positional and temporal encodings and transformer layers to come up with the next feasible set of high-level plan actions.

dropout of 0.3 to obtain feature maps of $512 \times 7 \times 7$. These are then fed into 2 convolutional layers with 256 and 64 filters of size 1×1 , respectively, and mapped using a fully connected layer to size 768.

The language input is tokenized using revtok⁸ and encoded using two transformer layers with 12 attention heads and an embedding size of 768, which are trained from scratch. The language and visual encodings are then concatenated and passed through two multimodal transformer layers with 12 attention heads trained from scratch in our experiments.

To predict actions and objects, the final embedding corresponding to each image from the input is projected using a fully connected layer and then passed to two independent softmax prediction heads over the action and object space, respectively. At inference time, the last action, which is not a padding token, is identified and used for prediction, along with the object, at the same time step. The E.T. Hierarchical model connects the two prediction heads by passing the output of the action head as an additional input to the object head. The E.T. + Mask model examines the action-object

⁸<https://pypi.org/project/revtok/>

pair predicted at inference time, and if they are found to be incompatible, replaces the action with the action of highest probability with the predicted object. It is assumed that the object is the more reliable prediction as it is more likely to be directly visible.

E Plan Prediction Model Hyperparameters

For training the E.T., E.T. + Mask and E.T. Hierarchical methods, we retained hyperparameters from the original TEACH paper (Padmakumar et al., 2022), except the batch size, without further hyperparameter tuning, and used the largest batch size that could fit in a single GPU of a p3.8xlarge AWS EC2 instance. Note that the following hyperparameters were kept constant for all experiments reported in this paper, and results in different tables arise only from changes in training data, model choice (between E.T., E.T. Hierarchical and E.T. + Mask), and plan execution method. We used the AdamW optimizer with 0.33 weight decay with a learning rate of $1e-4$ for the first 10 epochs and $1e-5$ for the last 10 epochs. We trained all models for 20 epochs with a batch size of 3, and report results from the final epoch. We replace sampling with rotation permutations of our training dataset per epoch, ensuring that every training example is seen exactly once in our dataset. For the language decoder in the transformer, we use a drop-out of 0.1, and for the encoder, we use a dropout of 0.1. The different E.T. models required 4 hours for preprocessing (extracting image features using the ResNet-50 backbone), and about 2 hours per model for training using 4 GPUs of a p3.8xlarge AWS EC2 instance. At inference time we could use a maximum of 3 GPUs for inference as one GPU was required by the simulator. When using 3 GPUs of a p3.8xlarge AWS EC2 instance, E.T. models took about 11 hours to complete inference jointly on the divided_val_seen and divided_test_seen splits and about 35 hours to complete inference jointly on the divided_val_unseen and divided_test_unseen splits. The time difference is due to the size of the various splits.

For the baseline BART model, we retain hyperparameters from the model presented in (Gella et al., 2022). We take the pretrained BART-base model from the Huggingface library⁹ and fine-tune for 20 epochs using a batch size of 2 per

⁹<https://huggingface.co/>

GPU. The training was done using gradient accumulation across 4 GPUs of a p3.8xlarge AWS EC2 instance. We use the AdamW optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.99$, $\epsilon = 1e08$ and weight decay of 0.01. We use a learning rate of $5e05$ with a linear warmup of over 500 steps. The BART model can be finetuned in under an hour using all 4 p3.8xlarge AWS EC2 instance GPUs. We first performed inference on the BART model and saved the predicted plans to file before separately executing them in the AI2-THOR simulator. This process can also be completed in under an hour. Executing stored plans either in the case of the BART model or the oracle conditions took about 2.5 hours using 3 GPUs of a p3.8xlarge AWS EC2 instance for the combined `divided_val_seen` and `divided_test_seen` splits and about 8 hours for the combined `divided_val_unseen` and `divided_test_unseen` splits.

F Data Statistics

The number of games and EDH instances in the TEACH data splits and batches of synthetic data used in this paper are included in Table 4.

Split	Number of games	Number of EDH instances
TEACH train	3121	3895
TEACH <code>divided_val_seen</code>	83	302
TEACH <code>divided_val_unseen</code>	309	1078
TEACH <code>divided_test_seen</code>	98	306
TEACH <code>divided_test_unseen</code>	303	1071
Synthetic 1x	1587	2034
Synthetic 2x	3272	3875
Synthetic 3x	4952	7201
Synthetic 4x	6360	9284

Table 4: Data statistics in TEACH and synthetic data splits.

G dialogue Simulation Details

This section provides a more detailed description of the dialogue simulation process. dialogue simulation begins by sampling which agent starts the interaction - the *Commander* or the *Follower*, each

with 50% probability - and the task and initial state of the environment in which to start the dialogue. While it is possible to create new initial states, we iterate over each initial state in the train split of the TEACH dataset, each of which is associated with a task to be completed in that state. For each agent, the *Commander* and the *Follower* we maintain a state for the agent that is factored into binary state features. We then use predefined probabilities for sampling different dialogue acts in each state and alternate taking turns between the two agents. In this implementation, we use the following dialogue acts defined in Gella et al. (2022):

- RequestForInstruction 1173
- Instruction 1174
- RequestForObjectLocationAndOtherDetails 1175
- InfoOnObjectLocationAndOtherDetails 1176
- Acknowledge 1177
- FeedbackPositive 1178
- FeedbackNegative 1179

However, we use both `FeedbackPositive` and `FeedbackNegative` only to end the dialogue either as a success or failure respectively. We additionally divide the `Instruction` dialogue act into two sub-types for convenience:

- `Instruction`: For communicating the task and its parameters 1185
- `Step`: For communicating a single desired object state change that would result in progress towards completion of the task, for example cleaning a mug to eventually fill it with coffee. 1188

We correspondingly also create a special `RequestStep` action. Besides this, we have two non-dialogue actions that an agent can also choose to perform:

- `ExecuteStep`: This is a cue to transition to actions in the environment. This is only performed by the *Follower* which identifies a rule-based plan to accomplish the desired state change and executes it with the aid of assisted plan execution described in section 4. 1195
- `DoNothing`: This allows an agent to skip a turn and hence avoids rigid turn taking in the 1202

1203	resultant dialogue, and introduces variability	Given the value of the state features, the next	1247
1204	in the amount of information communicated	state of the agent is computed as boolean functions	1248
1205	in the dialogue to better mimic real dialogues.	over state features included in Table 5. Given the	1249
1206	The state features used are:	dialogue acts we then sample dialogue acts accord-	1250
1207	• <i>dialogue_started</i> : Agents start with this	ing to Table 6 for the <i>Commander</i> and Table 7 for	1251
1208	feature set to <i>False</i> indicating that the agent	the <i>Follower</i> .	1252
1209	is in the initial state before any dialogue has	Given a dialogue act, we use the following tem-	1253
1210	taken place and must initiate dialogue. It gets	plates to get utterances:	1254
1211	set to <i>True</i> once an initial utterance has been	• <i>RequestForInstruction</i> : Hello, what	1255
1212	exchanged.	can I help you with today?	1256
1213	• <i>goal_communicated</i> : This feature is <i>False</i>	• <i>Instruction</i> : This is filled in with the	1257
1214	initially and set to <i>True</i> after an <i>Instruction</i>	description field from the task definition,	1258
1215	utterance has been sent from the <i>Commander</i>	for example <i>Make coffee</i> or <i>Put all Forks</i>	1259
1216	to the <i>Follower</i> communicating the high level	in any Sink.	1260
1217	task.	• <i>RequestStep</i> : Done. What should I do	1261
1218	• <i>cur_step_requested</i> : This feature is <i>False</i>	next?	1262
1219	initially, gets set to <i>True</i> when the <i>Follower</i>	• <i>Step</i> : This is filled it from the	1263
1220	sends a <i>RequestStep</i> action to the <i>Comman-</i>	<i>condition_failure_desc</i> field from	1264
1221	<i>der</i> and reset to <i>False</i> when environment ac-	the <i>TEACH</i> task definition. For example <i>The</i>	1265
1222	tions are executed.	<i>Mug</i> does not have coffee. or <i>The Fork</i>	1266
1223	• <i>cur_step_sent</i> : This feature is <i>False</i> ini-	must be placed in a Sink.	1267
1224	tially, gets set to <i>True</i> when the <i>Commander</i>	• <i>RequestForObjectLocationAndOtherDetails</i> :	1268
1225	sends the current step to the <i>Follower</i> through	Where can I find a/an <i><object></i> ? where	1269
1226	a <i>Step</i> dialogue act and reset to <i>False</i> when	<i><object></i> is identified from the task step.	1270
1227	environment actions are executed.	• <i>InfoOnObjectDetails</i> : You can find a/	1271
1228	• <i>cur_step_obj_requested</i> : This feature is	and <i><object></i> in/on a/an <i><object></i> near	1272
1229	<i>False</i> initially, gets set to <i>True</i> when the <i>Fol-</i>	a/an <i><object></i> where the reference objects	1273
1230	<i>lower</i> requests the location of an object and	are identified using location information in the	1274
1231	reset to <i>False</i> when environment actions are	simulator.	1275
1232	executed.	• <i>FeedbackPositive</i> : Great! We're all	1276
1233	• <i>cur_step_obj_sent</i> : This feature is <i>False</i>	done.	1277
1234	initially, gets set to <i>True</i> when the <i>Comman-</i>	If an <i>ExecuteStep</i> action is sampled, we then	1278
1235	<i>der</i> sends the location of an object and reset	synthesize the plan for the property to be changed	1279
1236	to <i>False</i> when environment actions are exe-	as part of the current task step using Table 8. This	1280
1237	cuted.	is then executed using assisted execution (Sec-	1281
1238	• <i>task_complete</i> : This feature is <i>False</i> ini-	tion 4, Appendix C) to obtain the final sequence	1282
1239	tially and gets set to <i>True</i> when the task is	of environment actions. If this produces the de-	1283
1240	completed successfully.	sired object state change, dialogue simulation con-	1284
1241	• <i>follower_stuck</i> : This feature is used to	tinues, otherwise it is terminated by entering the	1285
1242	identify failed dialogues. It is <i>False</i> initially	<i>follower_stuck</i> state. Our final implementation	1286
1243	and set to <i>true</i> if the <i>Follower</i> attempts envi-	of dialogue simulation is able to successfully simu-	1287
1244	ronment actions but is unable to accomplish	late a dialogue that completes the task in 35.4% of	1288
1245	the intended object state change. When this is	the initial states in the <i>TEACH</i> dataset.	1289
1246	identified, the dialogue is terminated early.		

dialogue_not_started	\neg dialogue_started
goal_or_step_start	dialogue_started \wedge \neg goal_communicated \wedge \neg task_complete \wedge \neg follower_stuck \wedge \neg cur_step_requested \wedge \neg cur_step_sent \wedge \neg cur_step_obj_requested \wedge \neg cur_step_obj_sent
step_start	dialogue_started \wedge goal_communicated \wedge \neg task_complete \wedge \neg follower_stuck \wedge \neg cur_step_requested \wedge \neg cur_step_sent \wedge \neg cur_step_obj_requested \wedge \neg cur_step_obj_sent
step_requested	dialogue_started \wedge \neg task_complete \wedge \neg follower_stuck \wedge cur_step_requested \wedge \neg cur_step_sent
step_sent	dialogue_started \wedge \neg task_complete \wedge \neg follower_stuck \wedge cur_step_sent \wedge \neg cur_step_obj_requested \wedge \neg cur_step_obj_sent
obj_loc_requested	dialogue_started \wedge \neg task_complete \wedge \neg follower_stuck \wedge cur_step_obj_requested \wedge \neg cur_step_obj_sent
obj_loc_sent	dialogue_started \wedge \neg task_complete \wedge \neg follower_stuck \wedge cur_step_obj_sent
task_complete	dialogue_started \wedge task_complete \wedge \neg follower_stuck
follower_stuck	dialogue_started \wedge \neg task_complete \wedge follower_stuck

Table 5: Boolean functions over dialogue state features to compute current dialogue state (\neg represents NOT and \wedge represents AND).

dialogue_not_started	Instruction	0.8
	Step	0.1
	DoNothing	0.1
goal_or_step_start	Instruction	0.8
	Step	0.1
	DoNothing	0.1
step_requested	Step	0.9
	DoNothing	0.1
step_start	Step	0.9
	DoNothing	0.1
obj_loc_requested	InfoOnObject	0.9
	Details	0.1
	DoNothing	0.1
follower_stuck	Feedback Negative	1.0
task_complete	Feedback Positive	1.0

Table 6: Probabilities for sampling various dialogue acts for the *Commander* given the dialogue state.

dialogue_not_started	RequestFor	1.0
	Instruction	1.0
goal_or_step_start	RequestStep	0.8
	ReqForObjLoc	0.1
	AndOD	0.1
step_start	RequestStep	0.8
	ReqForObjLoc	0.1
	AndOD	0.1
step_requested	ExecuteStep	1.0
step_sent	ReqForObjLoc	0.9
	AndOD	0.1
	ExecuteStep	0.1
obj_loc_requested	ExecuteStep	1.0
obj_loc_sent	ExecuteStep	1.0

Table 7: Probabilities for sampling various dialogue acts for the *Follower* given the dialogue state.

Desired Outcome	Plan
Cleaning	(Pickup, <TARGET>), (Place, Sink), (ToggleOn, Faucet), (ToggleOff, Faucet), (Pickup, <TARGET>), (Pour, Sink)
Making Coffee	(Pickup, <TARGET>), (Place, CoffeeMachine), (ToggleOn, CoffeeMachine)
Slicing	(Pickup, Knife), (Slice, <TARGET>), (Place, CounterTop)
Toasting Bread	(Pickup, <TARGET>), (Place, Toaster), (ToggleOn, Toaster), (Pickup, <TARGET>), (Place, CounterTop)
Cooking	(Pickup, <TARGET>), (ToggleOff, Microwave), (Open, Microwave), (Place, Microwave), (Close, Microwave), (ToggleOn, Microwave), (ToggleOff, Microwave), (Open, Microwave), (Pickup, <TARGET>), (Close, Microwave), (Place, CounterTop)
Placing	(Pickup, <TARGET>), (Place, <DESIRED_VALUE>)
Boiling	(Pickup, <TARGET>), (Place, Bowl), (Pickup, Bowl), (Place, Sink), (ToggleOn, Faucet), (ToggleOff, Faucet), (Pickup, Bowl), (ToggleOff, Microwave), (Open, Microwave), (Place, Microwave), (Close, Microwave), (ToggleOn, Microwave), (ToggleOff, Microwave), (Open, Microwave), (Pickup, <TARGET>), (Close, Microwave), (Place, CounterTop)
Fill With Water	(Pickup, Cup), (Place, Sink), (ToggleOn, Faucet), (ToggleOff, Faucet), (Pickup, Cup), (Pour, <TARGET>)
Turn On	(ToggleOn, <TARGET>)
Open	(Open, <TARGET>)

Table 8: Plans to accomplish object state changes