Sub-goal Distillation: A Method to Improve Small Language Agents

Anonymous ACL submission

Abstract

While Large Language Models (LLMs) have shown great promise as agents in interactive tasks, their high computational costs limit their utility, especially for long-horizon tasks. We propose a method for transferring the performance of an LLM with billions of parameters to a much smaller language model (770M parameters). Specifically, we develop a hierarchical agent composed of a planning module that learns via Knowledge Distillation from an LLM to generate sub-goals and an execution module that learns to achieve sub-goals with elementary actions. Because neither module relies on online access to an LLM at inference, our method has a fixed cost of LLM interactions all happening during training. In Science-World - a challenging interactive text environment - our approach outperforms standard imitation learning on elementary actions alone by 16.7% (absolute). Our analysis underscores our method's efficiency with respect to other LLMbased methods. We release our code and data for distillation at anon_url.com.

1 Introduction

011

017

018

019

037

041

Recently, Large Language Models (LLMs) have found applications in various fields, including answering questions, summarizing documents, translating languages, completing sentences, and serving as search assistants. They showcase a remarkable ability to make predictions based on input, enabling their use in generative AI applications to produce content based on input prompts (Devlin et al., 2018; Brown et al., 2020; Rae et al., 2021; Chowdhery et al., 2023; Scao et al., 2022; Patel and Pavlick, 2021; Han et al., 2021; Bommasani et al., 2021).

The promising advantage of LLMs is attributed to their training on extensive text datasets, resulting in impressive capabilities. This prior knowledge can be leveraged for action planning to solve tasks in robotics and reinforcement learning (Huang et al., 2022b; Brohan et al., 2023; Liang et al., 2023). Recent works have utilized in-context learning with LLMs to provide actions in autonomous decision-making agents and interactive environments (Mahowald et al., 2023; Yao et al., 2022; Schick et al., 2023; Shen et al., 2023; Nakano et al., 2021; Park et al., 2023; Lin et al., 2023; Brohan et al., 2023). 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

078

079

081

However, the extreme size of LLMs makes them computationally unaffordable for many applications. Moreover, closed-source models like ChatGPT (OpenAI, 2022) and GPT-4 (OpenAI, 2023) limit accessibility and reproducibility. Consequently, there is an increasing demand to find approaches that are less computationally intensive while still capitalizing on the knowledge embedded in LLMs. One prevalent technique is the use of Knowledge Distillation (KD) (Buciluă et al., 2006; Hinton et al., 2015), wherein a smaller model is trained with guidance from a larger model. Through this approach, we can leverage the knowledge in an LLM to train a more compact model with a reduced number of parameters.

Distilling knowledge from LLMs offers significant advantages, allowing for the training of specialized local models rather than depending on an LLM as a general model. This approach not only enhances privacy, particularly for systems with security-sensitive considerations like co-pilot models, but also provides greater flexibility in tailoring models for specific tasks. Additionally, the use of a smaller model offers the advantage of versatility across diverse applications without size constraints, including device models and mobile apps. Another challenge with LLMs is their susceptibility to hallucinations. This tendency poses a hindrance to their effective execution of long-tail planning, especially in interactive decision-making scenarios.

In our research, we leverage the knowledge of LLMs to train an autonomous agent for effective decision-making in complex interactive text environments, utilizing simple language models as our policy. Knowledge Distillation facilitates the training of smaller policies, allowing seamless integration of LLM knowledge. To address the challenges at hand, adopting a two-level planning approach proves beneficial for reducing hallucination – one for high-level reasoning to formulate sub-goals and another for low-level action planning to execute each sub-goal.

084

100

101

103

104

105

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

Fig. 1 illustrates this concept in the task of freezing water from ScienceWorld (Wang et al., 2022a). The agent's subtasks involve navigating to the kitchen, finding a thermometer and a metal pot, pouring water into the pot, placing it in the freezer, and continuously monitoring its temperature until frozen. These constitute sub-goals generated by a high-level model, with each sub-goal subsequently executed by a low-level model. The generation of sub-goals empowers an autonomous agent to expedite learning for the current task and reuse similar sub-goals in various tasks to have more generalization.

The contributions in this work are:

- We employ Knowledge Distillation from an LLM to train a high-level policy capable of generating sub-goals without making assumptions about the specific set of sub-goals. Notably, these sub-goals remain flexible, accommodating various sequences of actions.
- We demonstrate that employing Knowledge Distillation with hierarchical policies surpasses the performance achieved by both standalone imitation learning and its combination with in-context learning.
- We illustrate that this approach is more costeffective in terms of the number of calls to an LLM compared to other methods utilizing in-context learning.

2 Related Work

Using LLMs for Action Planning Recent works have demonstrated the ability of LLMs to perform 123 action planning for interactive decision making pro-124 cess without any additional training (Huang et al., 126 2022a). ReAct (Yao et al., 2022) proposes a way of prompting an LLM with interleave reasoning 127 step and action taking step. That led the resolu-128 tion of a variety of language-based reasoning and 129 decision-making tasks. This approach empowers 130

Task Description:

Your task is to change the state of matter of water. First, focus on the substance. Then, take actions that will cause it to change its state of matter.

Annotated Trainatory

Autotated Hajectory
Navigate_to(kitchen)
open door to kitchen
go to kitchen
Pick_up(thermometer)
pick up thermometer
Find(metal pot)
open cupboard
pick up metal pot
Fill(metal pot, water)
move metal pot to sink
activate sink
deactivate sink
pick up metal pot
Focus_on(substance in metal pot
focus on substance in metal pot
Freeze(water, metal pot)
pour metal pot into metal pot
pick up metal pot
open freezer
move metal pot to freezer
Monitor_temperature(metal pot, freeze)
examine substance in metal pot

Figure 1: Example of annotating an expert trajectory with sub-goals for a particular variation of task 1-4 (*change-the-state-of-matter-of*). Looking only at the original trajectory (i.e., ignoring the red rows), we gather the expert ended up changing the state of *water* to be frozen. The expert had to navigate to the kitchen, find a thermometer and a metal pot, pour water into the pot, place it in the freezer, and continually monitor its temperature until frozen. Each of those milestones (highlighted in red) can be considered a sub-goal, encompassing a sequence of actions. Sub-goals can be shared across different tasks, facilitating generalization. We opted for a format that looks like function calls to encourage reusability (e.g., *fill(metal pot, water)*).

the model to construct high-level plans for effective action. Reflexion (Shinn et al., 2023) draws inspiration from reinforcement learning, employing a framework to reinforce language agents through linguistic feedback. At the end of each trial, it uses self-reflection to determine what went wrong with the task and keeps it in a memory. Then it leverages this information for the next trial.

136

137

138

Some works use a programmatic LLM prompt 139 structure with available actions and objects in an 140 environment to translate natural language com-141 mands into robot policy code via few-shot exam-142 ples (Liang et al., 2023; Singh et al., 2023). Khot 143 et al. (2022) introduced a decomposed prompt-144 ing approach wherein a task is broken down into 145 simpler sub-tasks, allowing for recursive handling. 146 Subsequently, these sub-tasks are assigned to sub-147 task-specific LLMs, with both the decomposer 148 and the sub-task LLMs with their own few-shot 149 prompts. Sun et al. (2023) uses three steps, action 150 mining, plan formulation, and plan execution to 151 decompose a question into a sequence of actions 152 by few-shot prompting of LLMs. In Prasad et al. 153 (2023) tasks are decomposed explicitly by a separate LLM through prompting when an executor is unable to execute a given sub-task. 156

Imitation learning Some works employ imitation learning to train a language model as the 158 agent's policy, as seen in offline decision trans-159 160 formers (Torabi et al., 2018). The inputs consist of states, actions, and reward-to-go values, which are fed into a transformer. This transformer then 162 predicts actions in an autoregressive manner, utilizing a causal self-attention mask (Chen et al., 2021). 164 165 Contextual Action Language Model (CALM) (Yao et al., 2020) is another work which uses a fine-166 tuned language model with oracle data to generate 167 a set of candidate actions which are then passed to 168 a policy network to select the best one. In Nakano 169 et al. (2021), the authors fine-tune GPT-3 to ad-170 dress long-form questions within a web-browsing 171 context. Human feedback is employed as a direct 172 optimization measure for enhancing the quality of 173 answers generated by the model. 174

Knowledge Distillation: Knowledge Distillation 175 (KD) typically falls into two categories: black-box 176 KD and white-box KD. In black-box KD, only 177 the teacher's predictions are available for guid-178 ance, while in white-box KD, we have access to 179 the teacher's parameters (Gou et al., 2021). Recently, black-box KD has gained widespread use 181 for fine-tuning original models using self-instruct techniques, as proposed by Wang et al. (2022b), em-183 ploying self-generated instructions or for smaller 185 models (Taori et al., 2023; Chiang et al., 2023; Peng et al., 2023) through the utilization of promptresponse pairs generated by LLMs APIs. West 187 et al. (2021) introduces symbolic KD from text rather than logits. This process involves the trans-189

fer of knowledge from a large, general model to a more compact commonsense model, facilitated by a commonsense corpus, yielding a commonsense knowledge graph and model.

190

191

192

193

194

195

196

197

198

200

201

203

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

Complex interactive text environments In textbased games, an agent interacts with the environment by reading and writing text while aiming towards the end game or solving a given task. Out of the recent frameworks that deals with generating and interfacing text-based games (Côté et al., 2018; Hausknecht et al., 2019; Shridhar et al., 2021; Murugesan et al., 2021), we use ScienceWorld (Wang et al., 2022a) which is very complicated by having a large set of objects, actions, and tasks.

3 Model

In this paper, we propose to train a hierarchical policy by combining KD from an LLM and imitation learning from expert trajectories. This section describes both modules in detail and we refer the reader to Fig. 2 for a schematic view. We first formulate the problem as a POMDP (Section 3.1). Next, we describe what knowledge we are distilling from an LLM to guide the agent in accomplishing tasks (Section 3.2). Then, we detail how both the high-level and low-level policies of the hierarchical policy are trained (Section 3.3).

3.1 Problem Formulation

ScienceWorld (Wang et al., 2022a) can be defined as a partially observable Markov decision process (POMDP), where observations provide information solely on environmental changes induced by the current action. ScienceWorld is an interactive text environment meaning all task instructions, observations and actions are expressed in textual form. Importantly, both observations and rewards in this environment are conditioned by the ongoing task.

Given a language vocabulary V and an arbitrary maximum number of tokens N, an observation is defined such as $o \in \Omega \subset V^N$, a reward such as $r \in \mathbb{R}$ and an action as $a \in A \subset V^N$. Finally, a task or goal description is shown by $g \in G \subset V^N$.

We formalize the problem as a goal-augmented POMDP $M = (S, V, A, \Omega, G, T, R, O, \gamma)$ with S the state space, $A \subset V^N$ the action space, $\Omega \subset V^N$ the observation space, $G \subset V^N$ the goal space, $T : S \times A \times G \to S$ the goal-conditioned transition function, $R : S \times A \times G \to R$ the goalconditioned reward function, $O : S \to V^N$ an (unknown) observation function mapping a state to



Figure 2: On the left, a schematic view of our approach is shown. There are two modules: the sub-goal generator and action generator. The sub-goal generator provides a sub-goal for the action generator, which predicts the next action given the current sub-goal and history. On the right, the inputs and outputs of both modules are illustrated. The input comprises different parts including task description, completed sub-goal, current sub-goal, a history of recent actions-observations, and more, each highlighted in a distinct color.

a textual description and γ the discount factor. We assume $\gamma = 1$ in our experiments.

3.2 Distilling Knowledge from an LLM

The initial step in training our policies is creating a dataset. This dataset should include sub-goals along with their corresponding aligned sequences of actions for each task.

To generate this dataset, we assume access to a list of expert trajectories which we will augment. Specifically, each trajectory will be segmented and each segment will be annotated with a textual subgoal. To achieve this, we prompt an LLM with two in-context examples. Each example is composed of a task description and an expert trajectory for a similar task as the one we wish to annotate. The example also contains the expected response which corresponds to a set of sub-goals, accompanied by sequences of actions linked to each sub-goal.

Given the two in-context examples, the LLM is then instructed to generate a response for a new task description and its expert trajectory. It is important to note that these responses are collected only for the training tasks for which we assume having access to expert trajectories.

Fig. 3 illustrates the prompt examples for task 1-1 which is boiling a given substance. To ensure more uniform sub-goals that can generalize across tasks, we opted for a format that looks like function calls. Since that format was shown in the in-context examples, the LLM-generated sub-goals mimic this format as well making them easier to parse.

Given expert trajectories for some tasks can be long (+100 actions), the generated sub-sequence of actions corresponding to each sub-goal may not align exactly with the expert trajectory. Sometimes, it might miss certain actions, while in other instances, it might include additional actions, especially when there are repeated actions in the trajectory. To address this challenge, we use dynamic programming to determine the minimal number of actions that need to be added to or removed from the generated trajectory. So, if we need to remove some actions, we simply take them out of the generated trajectory. However, when some actions should be added, we ask the LLM to provide sub-goals for the actions missed in the generated trajectory. Then, we add them to the generated trajectory. By that, we make sure the generated trajectories remain in line with the expert data (see Appendix A.3). 275

276

277

278

279

281

282

283

285

286

287

289

290

291

293

294

295

296

297

299

300

301

302

303

304

305

306

307

308

In the resulting annotated dataset, each data point follows the same format as used by Lin et al. (2023) but with the added mention of completed sub-goals and the current sub-goal. Precisely, it corresponds to:

- **Input**: task description, number of steps, current score, completed sub-goal, current sub-goal, a history of 10 recent actions-observations, current items in the room, inventory, and the visited rooms.
- Target: next action, next sub-goal.

3.3 Hierarchical Imitation Learning

With the dataset obtained from distilling knowledge from an LLM, we can now focus on training the policies.

Low-level policy: The low-level policy is a language model (LM) which is trained through imitation learning using the annotated dataset. The goal is to have a model much smaller than an LLM so it can fit on a single machine and run faster, ideally

274

239



Figure 3: The figure demonstrates KD to generate sub-goals using an LLM. The LLM is presented with a prompt containing two in-context examples. Each example is composed of a task description in green and an expert trajectory detailing the steps to accomplish that task in blue. It also includes the expected set of sub-goals with their corresponding sequences of actions in red. Following this, we provide a new task description and trajectory, and we let the LLM generate the associated sub-goals and segmented actions.

below a billion of parameters. This policy learns
to predict the next action given the current task description, the 10 previous observation-action pairs,
the previous completed sub-goals, and the current
sub-goal. We refer to this policy as the action
generator.

High-level policy: The high-level policy is another LM with a reasonable size which is trained to
generate the next sub-goal. This policy conditions
on the same input information as for the action generator. We call this policy the sub-goal generator.

Hierarchical policy: During inference, we first leverage the high-level policy to generate a subgoal. This generated sub-goal is then fed into the action generator, allowing it to produce the next action aligned with the provided sub-goal. This sequential approach serves as a guiding cue for the action generator, particularly when the trajectory to achieve the goal is complex or long. Moreover, it serves to prevent the action generator from generating actions that might deviate the agent from the correct path, thereby improving the precision and relevance of the actions being generated.

4 Experiments

332

334

337

4.1 Environment

We chose ScienceWorld (Wang et al., 2022a) as the environment due to its complexity and the diverse range of tasks it encompasses. This environment is an interactive text-based game where the agent conducts elementary science experiments in a simulated environment. Each experiment is designed as a separate task. For example, "Your task is to boil water. For compounds without a boiling point, combusting the substance is also acceptable. First, focus on the substance. Then, take actions that will cause it to change its state of matter". To complete a task, the agent must perform multiple actions and receives the result of each action as an observation and a score. The observations and actions are in text format. An observation describes the changes in the environment, and the score is a numerical value ranging from 0% to 100%, indicating the degree of completion of the current task through the current action. 338

339

340

341

342

343

344

347

348

349

350

351

352

356

357

358

360

361

362

363

364

365

367

The environment includes 30 task types covering 10 different science topics (Appendix A.5). It has 10 different locations, more than 200 object types, and 25 action templates which makes the search space very larger for the agent. Each type of task has different variations in which the task objects, the agent's initial location, and random contents of each room are altered.

4.2 Experimental Setup

The environment has separate sets of variations for train and test. In the test variations, the combinations of objects and conditions are not seen in the train set. Following the experimental setup in (Lin et al., 2023), if the number of variations is more than 10, we consider only the first 10 variations.

Our base models for both policies is a pre-trained FLAN-T5-LARGE (Chung et al., 2022) with 700M parameters. For the both polices, we used greedy decoding at inference. We also conduct an ablation study over different model sizes (Fig. 4). For finetuning the policies, we use all the training tasks and their variations (3600 games in total) from ScienceWorld. We vary the number of training epochs in function of the size of the models (see Appendix A.4).

4.3 Baseline Agents

371

374

377

381

384

391

396

400

401

402

403

404

405

We compare our approach with other works that leverage LLMs. Some rely only on prompting such as SayCan, ReAct, and Reflexion, but SwiftSage also do imitation learning. Here is a brief description of each method.

SayCan: the LLM initially offers a set of actions along with their respective ranks. Then, a valuebased method is employed to re-rank these actions in order to determine the most rewarding action for execution (Brohan et al., 2023).

ReAct: the LLM generates actions by incorporating the provided prompt and the history of generated texts. It employs reasoning traces as intermediate thought steps during the action generation to refine a plan for the upcoming steps (Yao et al., 2022).

Reflexion: the language agent reflects the task feedback at each trial in the form of text and retains this information within an episodic memory. During the subsequent trial, it leverages the stored memory text to enhance its decision-making process (Shinn et al., 2023).

SwiftSage: it comprises two components: Swift, a fine-tuned LM to predict actions, and Sage, a module that queries ChatGPT for planning when the performance of Swift is inadequate (as determined by some handcrafted rules) (Lin et al., 2023).

Swift-only: this is the Swift part of the SwiftSage method which only has the fine-tuned LM to
predict the actions. We consider this method as
a strong baseline and the most comparable to our
approach as it relies on imitation learning without
the need for querying ChatGPT during inference.

412 4.4 Results and Analysis

413 Main Results: Table 1 compares the perfor-414 mance of the baselines with our approach in the ScienceWorld. The score for each task type is the average score (in percent) obtained for 10 test variations. Our approach demonstrates an overall performance of **63.43%**, surpassing Swift-only by 16.71% (33.9% relative increase), and showing a slight improvement over SwiftSage of 3.3% (5.3% relative). Interestingly, our method is able to solve all test variations (i.e., gets an average score of 100%) for 11 out of the 30 task types. In contrast, SwiftSage solves them only for 2 task types, and Swift-only, only for 4 task types. 415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

Additionally, we measured the performance of the agents with respect to the length of the tasks (a proxy for task complexity). The length of a task is determined by how many actions was needed by the expert to solve it.¹ Following Lin et al. (2023), we group the tasks into three categories: *Short* when the length is less than 20 actions, *Medium* when it falls between 20 and 50 (inclusively), and *Long* if above 50. As shown in Table 1, our approach outperforms other methods on short and medium tasks. On long tasks, we outperform all methods except SwiftSage, which has a substantial advantage here: The longer the task, the higher the chance it triggers one of the rules for Sage to take over.

A key motivation for our approach is costeffectiveness in terms of LLM queries. During training, we make one query to ChatGPT per task to identify the sub-goals within an expert trajectory. Sometimes mismatches occur between the expert trajectory and the actions assigned to each sub-goal by ChatGPT. When that is the case, we employ dynamic programming, with a maximum of 10 attempts per task. This contrasts with other baseline methods, where LLM is queried for each action, incurring considerably higher costs.

Why is it failing on some task types? The performance of our algorithm in some tasks are low, (see Table 4). In Table 1, the scores of two tasks are presented. One contributing factor is the variations in the test are very different from those in the training. For instance, the objects might be very different or the path to complete the task is very different and longer. The main culprit is the subgoal generator which is not able to generate good sub-goals.

As a concrete example (Table 2), in the test variations for task 3-3, the agent needs to go to kitchen and then fill a jug with water. When look-

¹Expert trajectories for test tasks were not seen during training.

Methods	SayCan*	ReAct*	Reflexion*	Swift-only	SwiftSage*	Ours
Overall Average	25.22	19.76	23.40	46.25	62.22	65.43
Solved Task Types	0/30	0/30	4/30	4/30	2/30	11/30
Short [†]	37.24	28.95	39.19	79.68	72.81	91.61
Medium	20.06	21.09	14.73	35.80	55.34	62.83
Long	18.66	11.23	16.27	25.36	57.99	45.35
Task 1-1	33.06	3.52	4.22	15.0	58.0	16.22
Task 3-3	99.56	76.19	72.54	59.5	66.9	5.6

Table 1: The table illustrates the overall average score (%) across all test tasks on the ScienceWorld benchmark for SayCan, ReAct, Reflexion, Swift-only, SwiftSage, and our algorithm (last column). The *Solved Task Types* row represents the number of task types for which an agent manages to solve all the test variations. The table also shows the average scores for tasks with a short, medium, and long length of expert trajectory. The rows *Task 1-1* and *Task 3-3* display the scores for each of them in which our approach does not work well in comparison with the other methods. The * denotes scores reported from (Lin et al., 2023).

ing at the transcript, we see the agent is able to go to kitchen but then when it arrives, the sub-goal generator issues a sub-goal which is not relevant, FocusOn(fountain). The agent attempts to focus on the fountain which is a wrong action and the game terminates with a score of 0.

Another example is task 1-1 (Table 2) in which the agent should boil a substance. It should first find the substance but since the substance is in a totally different location than those seen during training, the sub-goal generator is not able to generate a good sub-goal for this step. Consequently the agent will do other actions and exhaust all the allocated time steps.

Does scale matter? We conduct a comparison across various sizes of language models such as FLAN-T5-XL, FLAN-T5-BASE, and FLAN-T5-SMALL. Additionally, we evaluate T5-3B and T5-LARGE to determine the effectiveness of FLAN-T5 versus T5. The results are illustrated in Fig. 4. In our initial findings, we observed that FLAN-T5 outperforms T5 significantly. Moreover, our results reveal a positive correlation between the LM size and its performance - larger models generally yield better results. Intriguingly, we observe that for smaller models (FLAN-T5-SMALL and FLAN-T5-BASE), not conditioning on sub-goals works slightly better than including them. This might be indicative that the sub-goal generator is not expressive enough to generate meaningful and effective sub-goals which in turn impacts the action generator policy and leads to lower scores.

To further demonstrate the impact of the subgoal generator's size on the overall performance, we try pairing different sizes of sub-goal genera-

Example	e (task 3-3)
With Sub-goal	Expert Trajectory
NavigateTo(kitchen)	-
- go to art studio	- go to art studio
- go to outside	- go to outside
- go to kitchen	- go to kitchen
FocusOn(fountain)	-
-focus on fountain	- move jug to sink
	- activate sink
	- deactivate sink
	- pick up jug
Example	e (task 1-1)
With Sub-goal	Expert Trajectory
NavigateTo(kitchen)	-
- go to art studio	- go to art studio
- go to outside	- go to outside
- go to hallway	- go to hallway
NavigateTo(bedroom)	-
-go to bedroom	- go to workshop
	- pick up metal pot
	containing gallium

Table 2: Two instances where the performance of our algorithm is low. The first column displays the trajectory generated with sub-goals, while the second column presents the expert trajectory. Sub-goals are highlighted in dark red, accompanied by their corresponding actions, and incorrect actions are marked in red.

tor while limiting the action generator to be small. In Fig. 5, the average scores exhibit an upward trajectory. This can be attributed to the larger subgoal generators producing more accurate and relevant sub-goals, subsequently empowering the action generator to generate more correct actions. See Table 5 for a complete breakdown of the score per task type and per model size. 499

500

501

502

503

504

505

507

508

509

5 Discussion

In contrast to SwiftSage, which relies on interactive usage of the ChatGPT API to handle planning, our

493

494

495

496

497

498

464 465

466



Figure 4: Average scores across different model sizes for FLAN-T5 and T5. For T5 model, X-Large refers to T5-3B. The larger models work better and FLAN-T5 performs also better than T5. Dashed lines represent models that are not conditioning on any sub-goals ("no sg"). Those are equivalent to Swift-only.

522

524

531

532

533

535

539

approach makes use of a trained sub-goal generator to guide the action generator. Moreover, our 511 framework empowers the agent to retrieve a nearly optimal trajectory by supplying the appropriate subgoal. Nevertheless our framework has significantly reduced the frequency of API calls, which are both expensive and not universally accessible. ReAct, Reflexion, and SwiftSage require human annotations to correct sub-goals and predict a reasonable action. However in our approach, we do not need human help to predict sub-goals or provide precise prompts.

6 Conclusion

We introduce a straightforward yet highly effective approach for tackling complex text-based environments. Our framework leverages the knowledge of an LLM to extract sub-goals. A hierarchical policy of two LMs proposed: a high-level policy predicts a sub-goal, and a low-level policy, by using the predicted sub-goal, generates elementary actions. Through extensive experiments across 30 task types in ScienceWorld, our approach demonstrates increase performance compared to state-of-the-art baselines, including standard imitation learning and SwiftSage.

As future directions for this work, we aim to delve into further exploration of goal modification strategies when the agent encounters challenges in solving the current sub-goal. This could involve breaking down or transforming a sub-goal



Figure 5: Average scores across different sizes of subgoal generator while the action generator is kept to be base (blue) or small (green). Having larger sub-goal generators can significantly boost performance of small action generators.

into a more achievable form. Another venue for future research involves extending this approach to a multi-module environment. In such scenarios, the sub-goal generator could leverage each module as an independent source to generate diverse and context-specific sub-goals.

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

7 Limitations

We highlight the limitations of our approach and experimental setup.

First, we acknowledge the proposed technique has only been tested on a single type of environments - ScienceWorld - even though it encompasses 30 different types of task each with 10 variations. Also, ScienceWorld naturally provides language observations and actions which are readily compatible with LLMs. Tackling more visual environments would require multi-modal foundation models for the KD phase.

Second, while we did run the Swift-only model ourselves to check for reproducibility, we have not run SwiftSage for budget constraints. We acknowledge that ChatGPT might have improved (or worsen) since when Lin et al. (2023) ran their experiments.

References

Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportuni-

681

ties and risks of foundation models. *arXiv preprint arXiv:2108.07258*.

569

570

571

573

577

579 580

581

582

583

584

586

587

588

592

593

596

608

610

611

612

613

614

615

616

617

618

619

- Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, et al. 2023. Do as i can, not as i say: Grounding language in robotic affordances. In *Conference on Robot Learning*, pages 287–318. PMLR.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901.
- Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. 2006. Model compression. In Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 535–541.
 - Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. 2021. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See https://vicuna. Imsys. org (accessed 14 April 2023).
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. arXiv preprint arXiv:2210.11416.
- Marc-Alexandre Côté, Ákos Kádár, Xingdi Yuan, Ben Kybartas, Tavian Barnes, Emery Fine, James Moore, Ruo Yu Tao, Matthew Hausknecht, Layla El Asri, Mahmoud Adada, Wendy Tay, and Adam Trischler. 2018. Textworld: A learning environment for textbased games. *CoRR*, abs/1806.11532.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. 2021. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129:1789–1819.

- Xu Han, Zhengyan Zhang, Ning Ding, Yuxian Gu, Xiao Liu, Yuqi Huo, Jiezhong Qiu, Yuan Yao, Ao Zhang, Liang Zhang, et al. 2021. Pre-trained models: Past, present and future. *AI Open*, 2:225–250.
- Matthew Hausknecht, Prithviraj Ammanabrolu, Côté Marc-Alexandre, and Yuan Xingdi. 2019. Interactive fiction games: A colossal adventure. *CoRR*, abs/1909.05398.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv* preprint arXiv:1503.02531.
- Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. 2022a. Language models as zeroshot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*, pages 9118–9147. PMLR.
- Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. 2022b. Inner monologue: Embodied reasoning through planning with language models. arXiv preprint arXiv:2207.05608.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2022. Decomposed prompting: A modular approach for solving complex tasks. *arXiv preprint arXiv:2210.02406*.
- Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. 2023. Code as policies: Language model programs for embodied control. In 2023 IEEE International Conference on Robotics and Automation (ICRA), pages 9493–9500. IEEE.
- Bill Yuchen Lin, Yicheng Fu, Karina Yang, Prithviraj Ammanabrolu, Faeze Brahman, Shiyu Huang, Chandra Bhagavatula, Yejin Choi, and Xiang Ren. 2023. Swiftsage: A generative agent with fast and slow thinking for complex interactive tasks. *arXiv preprint arXiv:2305.17390*.
- Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. 2023. Dissociating language and thought in large language models: a cognitive perspective. *arXiv preprint arXiv:2301.06627*.
- Keerthiram Murugesan, Mattia Atzeni, Pavan Kapanipathi, Pushkar Shukla, Sadhana Kumaravel, Gerald Tesauro, Kartik Talamadupula, Mrinmaya Sachan, and Murray Campbell. 2021. Text-based RL Agents with Commonsense Knowledge: New Challenges, Environments and Baselines. In *Thirty Fifth AAAI Conference on Artificial Intelligence*.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted questionanswering with human feedback. *arXiv preprint arXiv:2112.09332*.

769

770

771

772

773

774

776

735

736

- OpenAI. 2022. Openai: Introducing chatgpt.
- OpenAI. 2023. Gpt-4 technical report.

682

690

703

705

707

710

711

713

714

716

717

718

719

720

721

722

724

725

726

727

728

729 730

731

732

733

734

- Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–22.
- Roma Patel and Ellie Pavlick. 2021. Mapping language models to grounded conceptual spaces. In *International Conference on Learning Representations*.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*.
 - Archiki Prasad, Alexander Koller, Mareike Hartmann, Peter Clark, Ashish Sabharwal, Mohit Bansal, and Tushar Khot. 2023. Adapt: As-needed decomposition and planning with language models. *arXiv preprint arXiv:2311.05772*.
- Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. 2021. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, et al. 2022. Bloom: A 176b-parameter openaccess multilingual language model. *arXiv preprint arXiv:2211.05100*.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface. *arXiv preprint arXiv:2303.17580*.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew Hausknecht. 2021. ALFWorld: Aligning Text and Embodied Environments for Interactive Learning. In Proceedings of the International Conference on Learning Representations (ICLR).
- Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox,

Jesse Thomason, and Animesh Garg. 2023. Progprompt: Generating situated robot task plans using large language models. In 2023 IEEE International Conference on Robotics and Automation (ICRA), pages 11523–11530. IEEE.

- Simeng Sun, Yang Liu, Shuohang Wang, Chenguang Zhu, and Mohit Iyyer. 2023. Pearl: Prompting large language models to plan and execute actions over long documents. *arXiv preprint arXiv:2305.14564*.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models. https://crfm. stanford. edu/2023/03/13/alpaca. html*, 3(6):7.
- Faraz Torabi, Garrett Warnell, and Peter Stone. 2018. Behavioral cloning from observation. *arXiv preprint arXiv:1805.01954*.
- Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. 2022a. Scienceworld: Is your agent smarter than a 5th grader? *arXiv preprint arXiv:2203.07540*.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022b. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*.
- Peter West, Chandra Bhagavatula, Jack Hessel, Jena D Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. 2021. Symbolic knowledge distillation: from general language models to commonsense models. *arXiv preprint arXiv:2110.07178*.
- Shunyu Yao, Rohan Rao, Matthew Hausknecht, and Karthik Narasimhan. 2020. Keep calm and explore: Language models for action generation in text-based games. *arXiv preprint arXiv:2010.02903*.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.

A Appendix

A.1 Few-shot Prompt for the Large Language Model

We employed the ChatGPT API as the large language model in our study. The structure of the ChatGPT prompt is comprised of three main components. Firstly, there is a general description of the environment. The second part includes two examples, each containing the task description, an expert trajectory, and a set of sub-goals with their corresponding action sequences. Lastly, the prompt presents a new task description, along with an expert trajectory, then we ask the LLM to generate sub-goals for this new task.

Here is the first part:

Description of the Environment

You are a helpful assistant. You are in a simulated environment as an agent. A task and its description will be given to you. Suggest the best actions the agent can take based on the things you see and the items in your inventory to complete the task. Only use valid actions and objects. If you know what are around, then suggest the following actions. You are allowed to do the following actions with the objects. Open or close OBJ meaning open or close a container, Deactivate or activate OBJ meaning activate or deactivate a device, connect OBJ to OBJ meaning connect electrical components, disconnect OBJ meaning disconnect electrical components, use OBJ [on OBJ] meaning use a device/item, look around meaning describe the current room, look at OBJ meaning describe an object in detail, look in OBJ meaning describe a container's contents, read OBJ meaning read a note or book, move OBJ to OBJ meaning move an object to a container, pick up OBJ meaning move an object to the inventory, put down OBJ meaning drop an inventory item, pour OBJ into OBJ meaning pour a liquid into a container, dunk OBJ into OBJ meaning dunk a container into a liquid, mix OBJ meaning chemically mix a container, go to LOC meaning move to a new location, teleport to LOC meaning teleport to a specific room, eat OBJ meaning eat a food, flush OBJ meaning flush a toilet, focus on OBJ meaning signal intent on a task object, wait [DURATION] meaning take no action for some duration, task meaning describe current task, inventory meaning list agent's inventory, OBJ means objects. LOC means location. There are 10 locations centered around a house theme. These are: kitchen, bathroom, workshop, art studio, greenhouse, outside, bedroom, living room, foundry.

785

787

779

778

780 781

782

783

Here are the two examples for task 4 - 1, which is *find-living-thing*. This is the second part of the prompt:

Two-shot examples of Task 4-1 (find-living-thing)

Example 1 [Task Description] Your task is to find a(n) living thing. First, focus on the thing. Then, move it to the red box in the kitchen.

[Expert trajectory] Here is the goal path to achieve to the goal: open door to greenhouse, go to greenhouse, open door to outside, go to outside, focus on dove, pick up dove, open door to kitchen, go to kitchen, move egg dove egg in inventory to red box

Based on the given goal path, provide me with the functional format of high-level sub-tasks to complete this task and their corresponding actions.

[sub-goals] 1- navigate_to(greenhouse): {'open door to greenhouse', 'go to greenhouse'} 2- navigate_to(outside): {'open door to outside', 'go to outside'} 3- Focus_on(dove): {'focus on dove'} 4pick_up(dove): {'pick up dove'} 5- navigate_to(kitchen): {'open door to kitchen', 'go to kitchen'} 6move(dove egg, red box): {'move dove egg in inventory to red box'}

Example 2 [Task Description] Your task is to find a(n) living thing. First, focus on the thing. Then, move it to the green box in the kitchen.

[Expert trajectory] Here is the goal path to achieve to the goal: open door to kitchen, go to kitchen, open door to outside, go to outside, focus on egg turtle, pick up egg turtle, open door to kitchen, go to kitchen, move egg turtle egg in inventory to green box

Based on the given goal path, provide me with the functional format of high-level sub-tasks to complete this task and their corresponding actions.

[sub-goals] 1- navigate_to(kitchen): {'open door to kitchen', 'go to kitchen'} 2- navigate_to(outside): {'open door to outside', 'go to outside'} 3- Focus_on(egg turtle): {'focus on egg turtle'} 4- pick_up(egg turtle): {'pick up egg turtle'} 5- navigate_to(kitchen): {'open door to kitchen', 'go to kitchen'} 6move(egg turtle, green box): {'move egg turtle in inventory to green box'}

The third part is just a new task description with an expert trajectory:

Request for sub-goal generating

[Task Description] Your task is to find a(n) living thing. First, focus on the thing. Then, move it to the green box in the living room.

[Expert trajectory] Here is the goal path to achieve to the goal: open door to hallway, go to hallway, open door to greenhouse, go to greenhouse, open door to outside, go to outside, focus on baby baby beaver, pick up baby beaver, open door to greenhouse, go to greenhouse, open door to hallway, go to hallway, open door to living room, go to living room, move baby baby beaver in inventory to green box

Based on the given goal path, provide me with the functional format of high-level sub-tasks to complete this task and their corresponding actions.

790

791 A.2 Input Prompt for the Policies

Here, we show the inputs and outputs for both the action generator and sub-goal generator. We followed
the format used by SwiftSage (Lin et al., 2023), incorporating sub-goals and making minor textual
adjustments accordingly.

Input for the action generator

[Task Description] Your task is to find a(n) living thing. First, focus on the thing. Then, move it to the red box in the kitchen;

Time: 4; Score: 16; < /s >

Completed subtasks are: navigate_to(greenhouse), navigate_to(outside). The current subtask is Focus_on(dove); </s>

Action history: $< extra_id_4 > \text{look around (+0)} \rightarrow N/A | < extra_id_3 > \text{go to greenhouse (+16)} \rightarrow You move to the greenhouse. | <math>< extra_id_2 > \text{open door to outside (+0)} \rightarrow \text{The door is already open.}$ | $< extra_id_1 > \text{go to outside (+0)} \rightarrow \text{You move to the outside.}$ | < /s >

Current environment: This outside location is called the outside. Here you see: | the agent | a substance called air | an axe | a blue jay egg | a butterfly egg | a dove egg | a fire pit | a fountain (containing a substance called water) | the ground | a substance called wood | You also see: | A door to the foundry | A door to the greenhouse | A door to the kitchen | < /s >;

Current inventory: In your inventory, you see: | an orange | </s>;

Visited rooms: hallway, greenhouse, outside $\langle s \rangle$;

What action should you do next? </s>

Output for the action generator

focus on dove

Input for the sub-goal generator

[Task Description] Your task is to find a(n) living thing. First, focus on the thing. Then, move it to the red box in the kitchen;

Time: 4; Score: 16; < /s >

The previous subtasks are: navigate_to(greenhouse). The current subtask is navigate_to(outside) ;< /s >

Action history: $\langle extra_id_4 \rangle$ look around (+0) \rightarrow N/A | $\langle extra_id_3 \rangle$ go to greenhouse (+16) \rightarrow You move to the greenhouse. | $\langle extra_id_2 \rangle$ open door to outside (+0) \rightarrow The door is already open. | $\langle extra_id_1 \rangle$ go to outside (+0) \rightarrow You move to the outside. | $\langle /s \rangle$

Current environment: This outside location is called the outside. Here you see: | the agent | a substance called air | an axe | a blue jay egg | a butterfly egg | a dove egg | a fire pit | a fountain (containing a substance called water) | the ground | a substance called wood | You also see: | A door to the foundry | A door to the greenhouse | A door to the kitchen | < /s >;

Current inventory: In your inventory, you see: | an orange | < /s >;

Visited rooms: hallway, greenhouse, outside $\langle s \rangle$;

What subtask should you do next? </s>

Output for the sub-goal generator

Focus_on(dove)

A.3 Matching Generated Trajectory

The trajectory produced by the LLM may not always align with the expert trajectory. In such instances, we employ dynamic programming to identify the minimum set of actions that need to be added to or removed from the generated trajectory. For actions absent in the generated trajectory that should be added, we initially determine if they are sequential. Then, for each sequence identified, we prompt the LLM once more to generate a sub-goal corresponding to that specific sequence, see Fig. 6.

A.4 Implementation Details

We set the maximum number of steps per episode to 100. Additionally, we implemented an alternative termination condition: if the scores remain unchanged over the last 50 steps, we stop the episode. This

795

797

800 801

802 803 804

805 806



Figure 6: This figure illustrates instances where the trajectory generated by the LLM deviates from the provided expert trajectory. In this example which is for boiling task, certain actions are omitted in the generated trajectory, indicated in blue in the left box. To address these missing actions, we group them into sequences and prompt the LLM to generate sub-goals for them. If the generated trajectory includes additional actions, such as the green action on the right, we simply remove them to align with the expert trajectory.

prevents the agent from getting stuck in a repetitive loop of actions that do not yield any changes, such as navigating between rooms. We selected this threshold to ensure that it is not too restrictive, considering that certain tasks with lengthy trajectories may require a long sequence of actions to achieve a reward.

The learning rate in all of the experiments is 1e - 4. The values for $max_source_length = 1024$ and for $max_target_length = 30$. The batch size for training is 8. We set the number of epochs to 20 during training. In the evaluation phase, checkpoints were selected based on their loss values, encompassing the checkpoint with the lowest loss and the subsequent three checkpoints. The final choice for test was determined among these checkpoints, with priority given to the one demonstrating the highest score in the test set.

For both the action generator and sub-goal generator, we used greedy decoding. When the generated actions is invalid we attempt to find the closes match from the list of admissible commands provided by ScienceWorld.

A.5 Dataset Statistics

808

810

811

812

814

818

819

820

821

822

823

In Table 3, we provide tasks' names and their variations for all of the tasks in ScienceWorld (Wang et al., 2022a). For each task, variations are partitioned into 50% training, 25% development, and 25% test sets. In the development and test sets, variations include substances, animals, or plants that are not seen in the training.

A.6 Scores of Each Task

The average score of each task of the ScienceWorld is shown in Table 4. The methods are SayCan, ReAct, Reflexion, Swift-only, SwiftSage and our algorithm. In all of them ChatGPT is used as the LLM. The language models for Swift-only, SwiftSage and our algorithm are FLAN-T5-LARGE. For the methods, SayCan, ReAct, Reflexion, and SwiftSage we used the results from Lin et al. (2023). However, we reproduced the results for Swift-only and we found a lower performance than what was reported in the paper. So here we presented our scores.

1-1Matter1-2Matter1-3Matter1-4Matter1-4Matter2-1Measurement2-2Measurement2-3Measurement3-1Electricity3-2Electricity3-3Electricity	boil melt freeze e-the-state-of-matter-of use-thermometer lting-point-known-substance ting-point-unknown-substance power-component renewable-vs-nonrenewable-energy	$ \begin{array}{r} 30 \\ 30 \\ 30 \\ 30 \\ 540 \\ 436 \\ 300 \\ 20 \\ 20 \end{array} $
1-2Matter1-3Matter1-4Matter1-4Matter2-1Measurement2-2Measurement2-3Measurement3-1Electricity3-2Electricity2-3Electricity	melt freeze e-the-state-of-matter-of use-thermometer elting-point-known-substance ting-point-unknown-substance power-component renewable-vs-nonrenewable-energy	30 30 30 540 436 300 20
1-3Matter1-4Matter2-1Measurement2-2Measurement2-3Measurement3-1Electricity3-2Electricity2-3Electricity	freeze e-the-state-of-matter-of use-thermometer elting-point-known-substance ting-point-unknown-substance power-component renewable-vs-nonrenewable-energy	30 30 540 436 300 20
1-4Matterchange2-1Measurement2-2Measurement2-3Measurement3-1Electricity3-2Electricity2-3Floetricity	e-the-state-of-matter-of use-thermometer lting-point-known-substance ting-point-unknown-substance ower-component renewable-vs-nonrenewable-energy	30 540 436 300 20
 2-1 Measurement 2-2 Measurement measure-met 2-3 Measurement measure-met 3-1 Electricity power-component- 3-2 Electricity 	use-thermometer elting-point-known-substance ting-point-unknown-substance power-component renewable-vs-nonrenewable-energy	540 436 300 20
2-2Measurementmeasure-met2-3Measurementmeasure-melt3-1Electricityg3-2Electricitypower-component-3-3Electricitypower-component-	elting-point-known-substance ting-point-unknown-substance power-component renewable-vs-nonrenewable-energy	436 300 20
2-3Measurementmeasure-melt3-1Electricityp3-2Electricitypower-component-3-3Electricityp	ting-point-unknown-substance oower-component renewable-vs-nonrenewable-energy	300 20
3-1Electricityp3-2Electricitypower-component-3-3Electricity	power-component renewable-vs-nonrenewable-energy	20
3-2 Electricity power-component-	renewable-vs-nonrenewable-energy	
2.2 Electricity		20
5-5 Electricity	test-conductivity	900
3-4 Electricity test-conduct	ivity-of-unknown-substances	600
4-1 Classification	find-living-thing	300
4-2 Classification fir	nd-non-living-thing	300
4-3 Classification	find-plant	300
4-4 Classification	find-animal	300
5-1 Biology	grow-plant	126
5-2 Biology	grow-fruit	126
6-1 Chemistry	chemistry-mix	32
6-2 Chemistry chemistry-	-mix-paint-secondary-color	36
6-3 Chemistry chemistry	y-mix-paint-tertiary-color	36
7-1 Biology life	espan-longest-lived	125
7-2 Biology life	espan-shortest-lived	125
7-3 Biology lifespan-long	gest-lived-then-shortest-lived	125
8-1 Biology id	entify-life-stages-1	14
8-2 Biology ide	entify-life-stages-2	10
9-1 Forces inclined	1-plane-determine-angle	168
9-2 Forces inclined-pla	ane-friction-named-surfaces	1386
9-3 Forces inclined-plan	ne-friction-unnamed-surfaces	162
10-1 Biology mendeli	an-genetics-known-plant	120
10-2 Biology mendelia		400

Table 3: ScienceWorld's tasks names and numbers of variations.

The scores for SayCan, ReAct, Reflexion, and SwiftSage when they use GPT4 as the LLM are higher according to the results reported in Lin et al. (2023). However, due to limited access to GPT-4, we utilized ChatGPT, and thus, we present the results obtained using ChatGPT.

Task Type	Length	SayCan	ReAct	Reflexion	Swift-only	SwiftSage	Ours
1-1	107.7	0	3.52	4.22	15.0	58.0	16.22
1-2	78.6	0	13.70	10.61	24.4	58.5	0.2
1-3	88.9	0	7.78	7.78	32.2	38.5	30.33
1-4	75.2	0	9.88	0.92	57.4	62.5	65.0
2-1	21.4	1.7	7.19	5.92	9.4	47.9	98.55
2-2	35.2	14.1	6.10	28.59	6.7	53.3	46.0
2-3	65.0	93.7	22.37	22.37	5.7	48.6	45.0
3-1	13.6	19.3	56.0	100.0	70.0	72.7	100
3-2	20.8	8.7	54.33	17.45	48.3	50.3	78.2
3-3	25.6	22.0	76.19	72.54	59.5	66.9	5.6
3-4	29.0	36.4	36.4	70.22	69.0	78.1	46.0
4-1	14.6	11.7	26.67	64.93	100.0	100.0	100
4-2	8.8	76.0	80.0	87.27	100.0	97.5	100
4-3	12.6	11.4	53.33	16.42	94.4	58.3	100
4-4	14.6	9.5	27.50	100.0	100.0	100.0	100
5-1	69.5	11.3	11.1	5.8	13.4	57.5	100
5-2	79.6	75.0	18.8	47.6	44.6	50.9	65
6-1	33.6	13.5	35.0	22.4	26.2	43.2	66.87
6-2	15.1	25.0	20.0	10.0	53.3	63.3	100
6-3	23.0	58.4	16.7	40.0	11.1	27.4	100
7-1	7.0	75.0	37.5	75.0	83.3	75.0	100
7-2	7.0	100.0	50.0	75.0	100.0	60.0	100
7-3	8.0	31.7	31.7	28.1	77.8	68.3	100
8-1	40.0	5.6	4.2	2.8	33.0	75.6	60.0
8-2	16.3	12.8	7.0	8.2	8.0	33.0	25.0
9-1	97.0	38.0	28.5	100.0	73.3	54.0	52.5
9-2	84.9	4.2	10.0	17.5	73.3	63.3	53.6
9-3	123.1	0.0	0.0	1.7	53.3	77.0	64.4
10-1	130.1	1.3	24.5	1.3	17.0	76.0	30.76
10-2	132.1	0.3	11.7	6.0	17.0	51.1	30.76
Average	49.26	25.22	19.76	23.40	49.22	62.22	65.43

Table 4: The table illustrates the results for each task of the ScienceWorld for SayCan, ReAct, Reflexion, Swift-only, SwiftSage and our algorithm. Each row shows the average score of the test variations for a task type. Column *Length* shows the average lengths of the expert trajectories.

A.7 Task Scores Across Various Model Sizes

In Table 5, the average scores for models of different sizes are presented. In this experiment, the *base* model was employed for the action generator, while the size of the sub-goal generator was varied across small, base, large, and x-large. The first four columns of the table display their respective scores. Additionally, using the *small* model for the action generator, the experiment was repeated with various sizes of the sub-goal generator, and the results are shown in the second set of four columns. The last column illustrates the performance when both the action generator and the sub-goal generator are x-large, achieving the highest score.

Action generator	Base			Small			X-Large		
Sg generator Tasks	Small	Base	Large	X-Large	Small	Base	Large	X-Large	X-Large
1-1	0.2	0.0	0.0	4.2	0.0	0.0	0.0	0.0	16.55
1-2	0.0	0.0	0.0	0.22	0.0	0.0	0.44	0.22	17.11
1-3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
1-4	0.1	0.0	0.0	7.77	0.33	0.0	7.77	19.44	20.44
2-1	10.0	20.0	100	90.0	30.0	10.0	80.0	70.0	100
2-2	13.2	20.0	20.0	77.5	0.0	20.0	20.0	65.7	77.5
2-3	0.0	20.0	0.0	64.0	0.0	20.0	20.0	60.0	45.0
3-1	0.0	78.0	84.0	79.4	0.0	48.0	77.8	71.2	100
3-2	10.6	36.4	47.0	77.2	0.0	27.8	49.0	56.4	77.2
3-3	26.0	20.0	22.0	6.4	7.3	10.5	3.5	1.5	28.0
3-4	40.0	30.0	40.0	30.0	50.31	30.0	60.0	30.0	47.05
4-1	100	100	100	100	100	100	100	90	100
4-2	100	100	100	100	100	100	100	100	100
4-3	100	100	92.5	93	100	100	92.5	94.0	100
4-4	100	100	100	100	90.9	100	100	100	100
5-1	66.2	56.18	100	72.7	82.8	67.1	100	72.7	100
5-2	31.8	7.45	39.6	69.4	35.8	0.0	47.1	45.7	70.0
6-1	3.12	31.25	43.62	36.62	12.5	28.12	38.75	34.62	56.25
6-2	20.0	10.0	66.66	74.44	18.88	32.22	33.33	35.55	91.11
6-3	23.0	13.33	33.66	52.33	1.1	6.33	18.22	17.88	54.88
7-1	80.0	100	100	100	90.9	100	100	100	100
7-2	20.0	78.12	80.0	100	40.0	50.0	50.0	60.0	100
7-3	100	100	100	100	70.0	70.0	70.0	70.0	100
8-1	11.6	18.4	60.0	40.0	0.0	18.4	60.0	40.0	60.0
8-2	0.0	0	0.0	42.5	0.0	0.0	0.0	25.0	67.5
9-1	28.0	35.65	24	70	42.0	44.0	33.0	74.0	84.0
9-2	26.0	32.6	23	65	35.6	44.0	32.0	70.0	82.0
9-3	27.0	29.89	29	66	30.0	43.0	31.0	69.0	83.0
10-1	27.0	26.9	30.1	28.2	24.9	25.2	25.2	25.2	28.2
10-2	30.2	27.75	30.4	27.7	23.0	19.5	29.9	18.9	30.1
Average	33.13	39.73	48.85	59.15	32.87	37.13	45.98	50.56	67.86

Table 5: The table displays the average scores for the models of different sizes. The first four columns depict scores when a model with a base size is used as the action generator. The subsequent four columns illustrate scores when the action generator is small size. The last column shows the results when both the action generator and sub-goal generator are x-large.

A.8 Scores for Different Model Sizes Without sub-goals

843 844

845

Table 6 presents task scores without the utilization of sub-goals, akin to the Swift-only method but employing language models of varying sizes. All models are FLAN-T5.

Task Type	Small	Base	Large	X-Large
1-1	0.4	0.0	14.0	14.0
1-2	11.8	0.0	22.88	22.0
1-3	0.0	0.0	24.55	20.0
1-4	7.6	0.0	0.88	0.88
2-1	40.0	80.0	70.32	91.66
2-2	43.2	45.6	6.56	45.82
2-3	40.0	43.0	5.76	41.0
3-1	59.8	85.2	85.2	74
3-2	41.4	51.0	41.1	43.4
3-3	19.2	9.4	64.88	33.29
3-4	50.0	33.0	81.26	58.75
4-1	100	95.0	100	100
4-2	100	97.7	100	100
4-3	80.0	92.0	80.0	81.25
4-4	100	95.0	100	100
5-1	64.2	64.4	24.0	60.42
5-2	34.5	35.0	43.57	42.53
6-1	12.75	29.62	26.0	58.12
6-2	0.0	11.11	36.0	34.11
6-3	3.6	0.0	11.33	7.77
7-1	90.9	100	100	100
7-2	40.0	100	100	100
7-3	74.7	78.34	87.6	100
8-1	8.2	35.2	46.0	34.4
8-2	0.0	0.0	8.0	41.0
9-1	20.0	20.0	40.7	40.0
9-2	19.0	20.0	43.6	43.0
9-3	20.0	20.0	44.52	44.0
10-1	22.8	23.0	19.66	23.0
10-2	22.9	22.9	20.22	23.0
Average	37.56	42.88	46.25	52.88

Table 6: Scores for each task without the utilization of sub-goals across various model sizes. All models are FLAN-T5.

The results for the T5 model, including T5-LARGE and T5-3B, are presented in Table 7. The outcomes

are shown for both scenarios—with and without the integration of sub-goals. In the experiments where

sub-goals were employed, equivalent sizes were utilized for both the action generator and sub-goal

A.9 Results for T5 Language Model With and Without sub-goals

847 848

849

850

generator.

	No sub-goals		With sub-goals		
Task Type	Large	X-large	Large	X-large	
1-1	16.66	0.2	0.55	18.0	
1-2	0.0	56.0	5	19.66	
1-3	0.0	8.55	4.77	0.0	
1-4	0.0	30.77	16.22	38.0	
2-1	9.1	91.66	70.0	70.0	
2-2	37.8	41.0	20.0	27.7	
2-3	10.0	30.0	34.0	44.0	
3-1	85.2	78.6	39.2	48.6	
3-2	46.0	55.2	64.8	30.6	
3-3	50.2	14.61	1.0	31.76	
3-4	30.0	35.71	40.0	50.0	
4-1	100	100	100	100	
4-2	100	100	100	100	
4-3	90.0	93.7	100	100	
4-4	100	100	100	74.33	
5-1	100	64.2	91.5	100	
5-2	8.9	37.8	2.1	100	
6-1	41.62	43.62	21.87	31.25	
6-2	15.55	100	55.55	73.33	
6-3	4.77	80.0	29.2	35.0	
7-1	90.0	100	100	100	
7-2	100	100	30.0	100	
7-3	80.0	60.0	100	100	
8-1	55.4	35.2	60.0	69.0	
8-2	0.0	0.0	0.0	0.0	
9-1	20.0	20.0	36.0	33.33	
9-2	21.0	20.0	34.2	36.1	
9-3	22.0	20.0	37.1	39.0	
10-1	27.9	23.0	29.69	30.4	
10-2	28.0	22.9	28.75	28.0	
Average	43.00	52.09	45.05	54.26	

Table 7: Scores for the T5 model are depicted under two conditions: without sub-goals and with sub-goals.