

TaskGen: A Task-Based, Memory-Infused Agentic Framework using StrictJSON

Anonymous EMNLP submission

Abstract

TaskGen is an open-sourced agentic framework which uses an Agent to solve an arbitrary task by breaking them down into subtasks. Each subtask is mapped to an Equipped Function or another Agent to execute. In order to reduce verbosity (and hence token usage), TaskGen uses StrictJSON that ensures JSON output from the Large Language Model (LLM), along with additional features such as type checking and iterative error correction. Key to the philosophy of TaskGen is the management of information/memory on a need-to-know basis. We empirically evaluate TaskGen on various environments such as 40x40 dynamic maze navigation with changing obstacle locations (**100%** solve rate), TextWorld escape room solving with dense rewards and detailed goals (**96%** solve rate), web browsing (**69%** of actions successful), solving the MATH dataset (**71%** solve rate over 100 Level-5 problems), Retrieval Augmented Generation on NaturalQuestions dataset (F1 score of **47.03%**).

1 Introduction

TaskGen is an open-sourced agentic framework which breaks down a task into subtasks, each of which are mapped to an Equipped Function or another Agent to execute. The Agents and Equipped Functions operate independently, but share context on a need-to-know basis with one another using Shared Memory (see Fig. 1).

TaskGen is designed to be less verbose, and hence incurs lower processing latency and costs with potentially improved accuracy, than most existing agentic frameworks which output free text such as AutoGPT (Yang et al., 2023a), BabyAGI (Nakajima, 2023), MetaGPT (Hong et al., 2023), AutoGen (Wu et al., 2023), ChatDev (Qian et al., 2023), CrewAI (Moura, 2023), LangChain/LangGraph (LangGraph, 2024).

Our Contributions. We propose a new open-sourced agentic framework named TaskGen:

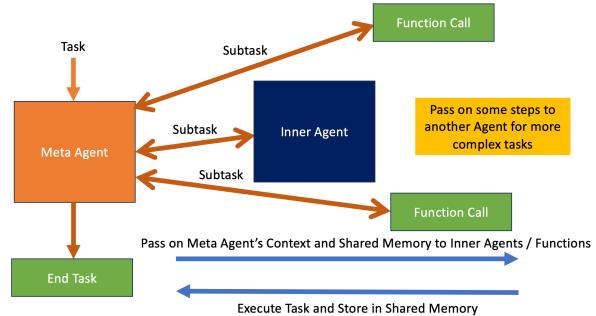


Figure 1: An Overview of TaskGen

1. TaskGen breaks a complex task down into bite-sized subtasks, each of which are mapped to an Equipped Function or Inner Agent to execute.
2. In contrast to free-form text output in agentic frameworks, TaskGen uses a concise JSON output for each part of the process. Specifically, it uses StrictJSON (Anonymous, 2023), which is an LLM output parser for JSON format with type checking, and helps ensure concise and extractable output which can be used for downstream tasks easily.
3. TaskGen has Shared Memory amongst various components on a need-to-know basis. This Shared Memory can come in the form of 1) **Subtasks Completed**, a list of past Equipped Functions inputs and outputs, or 2) **Shared Variables**, which stores important information that may also be of the form of long text or non-text modalities.
4. TaskGen utilises **Global Context** to inform the Agent of important information that may be dynamically changing as the task progresses. This allows the Agent to react to dynamic environments as the task progresses, or as the Agent switches tasks.
5. Lastly, as memory is key to learning and decision making, TaskGen implements memory of various abstraction spaces in the Agent's **Memory Bank**, which can be used to augment the prompt to the Agent via Retrieval Augmented Generation (RAG) (Lewis et al., 2020) based on semantic sim-

ilarity to the task. These memories are learnable via experience and can be used to influence future behaviour.

2 Motivation

We strive to create an Agent that can solve arbitrary tasks in arbitrary environments. However, when solving an arbitrary task, we could potentially do many actions, and there are many potential outcomes possible, as shown in Fig. 2. This is intractable for any Agent to manage and we need to limit the scope of what the Agent can do for more robust Agents.

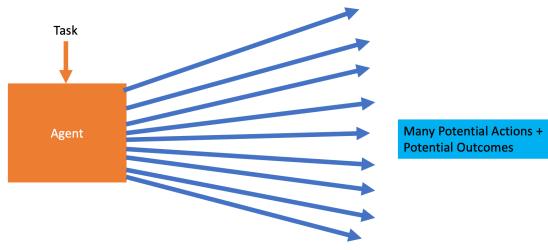


Figure 2: Intractable action space when solving an arbitrary task

Hence, we should limit the scope of the Agent by giving it only relevant Equipped Functions. This will help filter the vast action space into something tractable. Moreover, based on the Equipped Functions provided, we can break down a potentially complicated task into bite-sized subtasks, each of which can be solved entirely by one Equipped Function. This is shown in Fig. 3.

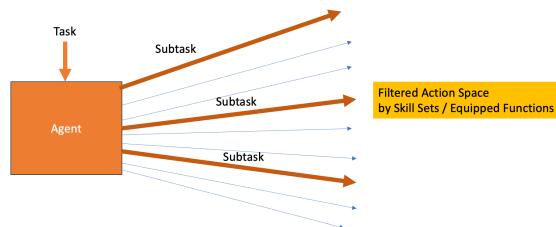


Figure 3: Constraining action space by Equipped Functions

In fact, for more complex tasks, we can even let another Agent be the Equipped Function. This Agent will henceforth be referred to as Inner Agent. This is similar to how a manager offloads tasks to each worker, each of whom have their own experiences and skills to do the task. By having intelligent Inner Agents as the Equipped Function,

the top-level agent (Meta Agent) will have greater processing capability. This is shown in Fig. 4.

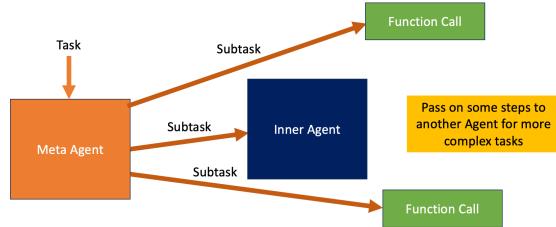


Figure 4: Inner Agents assigned as Equipped Functions to a Meta Agent helps increase processing capability

Infusing Shared Awareness. Each Equipped Function or Inner Agent would now be able to perform a subset of the entire task independently. However, they will need some shared context, as 1) the outcome of the subtask may influence other subtasks down the line, or 2) they may need input from earlier subtasks in order to perform their sub-task. To solve this problem, we implement a Shared Memory amongst the Meta Agent, Equipped Function and Inner Agents. Notably, we have two types of Shared Memory, 1) **Subtasks Completed** and 2) **Shared Variables**. This is shown in Fig. 1.

3 TaskGen Overall Design Philosophy

TaskGen has three key design philosophies.

Firstly, the output of each Agent or Equipped Function are made to be as concise as possible for minimal token use. This is done using StrictJSON. By ensuring a structured JSON output format with type checking, StrictJSON reduces verbosity typically associated with free-form text output in LLMs. This cuts down on latency and costs, and improves reliability of extracting output fields needed for downstream components. For a more in-depth run-through of StrictJSON, refer to Appendix A.

Secondly, we map each subtask to exactly one Equipped Function or Inner Agent, so as to guarantee executability of the subtask. Unlike AutoGPT (Yang et al., 2023a), we ensure that there are no infinite loops when executing subtasks. This is done via the following design guidelines:

1. An Agent can only call an Equipped Function or Inner Agent that is not above it in the hierarchy
2. Each Agent gets context relevant to its own processing abstraction space and are assigned Equipped Functions and Inner Agents suitable for that space

135 Lastly, information is only shared between
136 Agents and Equipped Functions on a need-to-know
137 basis. We have a shared pool of information in
138 Shared Memory, but we do only expose those that
139 are relevant to each Agent / Equipped Function.
140 This helps to reduce context length and minimise
141 the cognitive load on each part of the system.

142 4 The Core of TaskGen

143 4.1 Agent Definition

144 At the core of TaskGen is the definition of an Agent,
145 which consists of the following components:

- 146 1. **Agent Name:** Name of the Agent
- 147 2. **Agent Description:** Description of the Agent
- 148 3. **Equipped Functions:** List of Equipped Functions
149 and Inner Agents available to solve subtasks
- 150 4. **Assigned Task:** Agent's assigned task
- 151 5. **Subtasks Completed:** Python dictionary of past
152 subtasks that Agent has done, which detail the
153 Equipped Function's name and input parameters
154 and their corresponding output
- 155 6. **Shared Variables:** Python dictionary containing
156 variables that will be shared between Equipped
157 Functions and Agents
- 158 7. **Global Context:** Additional context to the
159 Agent that can reference persistent states, such as
160 those in **Shared Variables**
- 161 8. **Memory Bank:** Python dictionary containing
162 various abstraction spaces of memory that will be
163 retrieved via *top-k* retrieval via similarity to As-
164 signed Task

165 4.2 Imbuing Agentic Capabilities with 166 Equipped Functions

167 By default, an Agent comes pre-built with a
168 **use_llm** function, which uses an LLM with the
169 Agent Name and Agent Description as context to
170 perform a task, and an **end_task** function to end the
171 current task. Additionally, we can assign Equipped
172 Functions or Inner Agents to the Agent to imbue it
173 additional capabilities.

174 Equipped Functions come in two forms:

- 175 1. **Internal Functions** use an LLM to do processing
176 of input-output relations. They are useful for
177 tasks that are difficult for traditional rule-based
178 approaches to handle well, such as sentiment analysis
179 and summarisation.
- 180 2. **External Functions** utilise any Python function
181 to do processing to get output, which makes it very
182 easy for TaskGen to utilise functions from other
183 agentic frameworks such as LangChain or CrewAI.

They are suitable for tasks that can be called via fixed functions, or APIs, which guarantee reliability while imbuing additional functions to the LLM. As an aside, if we need a hybrid approach of rule-based fixed processes with flexibility of LLMs, an LLM can also be called within the External Function.

190 4.3 Choosing the Next Subtask

191 The core ability of an Agent is the ability to choose
192 the correct next subtask to fulfil the Assigned Task.
193 This is a non-trivial problem as it requires under-
194 standing of the Assigned Task, Agent Name, Agent
195 Description, Subtasks Completed, relevant Mem-
196 ory, Equipped Functions and Inner Agents in order
197 to make an informed decision.

198 In order to increase robustness in choosing the
199 right Equipped Function and corresponding input
200 parameters, we split it up into two steps.

201 **Step 1: Decide on subtask and corresponding
202 Equipped Function / Inner Agent.** The first step
203 simply takes the available information to the Agent
204 and does a Chain-of-Thought (CoT) (Wei et al.,
205 2022) prompting to elicit reasoning via thoughts,
206 leading to more accurate selection of subtask and
207 the corresponding Equipped Function / Inner Agent
208 in the following format:

- 209 1. **Observation:** Reflect on what has been done in
210 Subtasks Completed for Assigned Task
- 211 2. **Thoughts:** Brainstorm how to complete remain-
212 der of Assigned Task only given Observation
- 213 3. **Current Subtask:** What to do now in detail
214 with all context provided that can be done by one
215 Equipped Function for Assigned Task
- 216 4. **Equipped Function Name:** Name of Equipped
217 Function to use for Current Subtask

218 **Step 2: Decide on input parameters to
219 Equipped Function / Inner Agent.** Instead of
220 providing the entire list of Equipped Functions
221 / Inner Agents as per Step 1, we only give this
222 step information of the exact Equipped Function
223 / Inner Agent we have decided in Step 1, so as
224 to encourage greater output specificity. We then
225 generate the input parameters of the Equipped
226 Function / Inner Agent given the Current Subtask
227 and Equipped Function details (Equipped Function
228 Name, Equipped Function Description, Equipped
229 Function Input Parameter Description and type),
230 and uses StrictJSON to ensure that the input par-
231 ameters meet the type that is stated for in the Equipped
232 Function. This ensures robustness and reliability
233 for the input parameters.

5 Using TaskGen

Using TaskGen is extremely simple and is designed for any new user to learn it within 5 minutes. The steps needed are detailed as follows:

1. **Install TaskGen.** "pip install taskgen-ai"
2. **Define LLM.** This takes in a user prompt and system prompt as Python strings, and returns a Python string for the LLM generated response "def llm(user_prompt: str, system_prompt: str) -> str"
3. **Define Agent.** Simply define an Agent class with the Agent Name, Agent Description "agent = Agent(name, description, llm = llm)"
4. **Equip Functions.** Equip the Agent with Equipped Functions or Inner Agents to broaden the Agent's capabilities. "agent.assign_functions([fn_1, fn_2])"
5. **Run Agent.** Run the Agent with a task "agent.run(task)"
6. **Query Agent.** Query the Agent about Subtasks Completed "agent.reply_user(query)"

For an in-depth tutorial on how to use TaskGen, refer to Appendix B.

6 Benefits of TaskGen

The key philosophy of TaskGen is to be concise. This helps greatly with the performance of the overall system, as numerous studies (Xiong et al., 2023; Ding et al., 2024) have shown that an increase in context length generally leads to poorer performance on tasks referencing the context.

6.1 JSON is more concise than free text

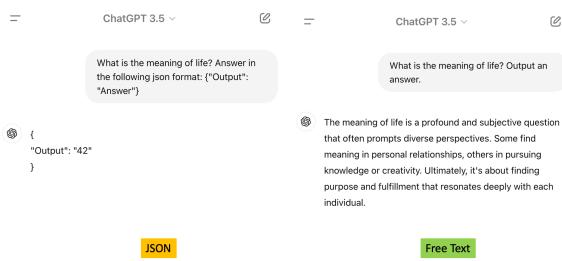


Figure 5: More concise output using JSON as compared to Free Text using gpt-3.5-turbo on 12 Jul 2024

Given a similar input prompt, asking the LLM to output in a JSON format generally gives much less verbose output as compared to free text. An example can be seen from Fig. 5 for a prompt

about the meaning of life. This is likely because the pre-training data of JSON on the web is more concise without much explanation, and the value of the field is very correlated to the key of the field. This means that we can use a JSON format to constrain the generation of the LLM to give the desired fields which we are interested in.

6.2 StrictJSON is more concise than JSON

```
JSON Schema for Parameters – 110 tokens
{
  "parameters": {
    "type": "object",
    "properties": {
      "location": {
        "type": "string",
        "description": "The city and state, e.g. San Francisco, CA",
      },
      "format": {
        "type": "string",
        "enum": ["celsius", "fahrenheit"],
        "description": "The temperature unit to use. Infer this from the users location.",
      }
    },
    "required": ["location", "format"]
  }
}
```

```
StrictJSON Schema for Parameters – 58 tokens
{
  "##Location##": "The city and state, e.g. San Francisco, CA, type: str",
  "##Format##": "The temperature unit to use. Infer this from the users location, type: Enum['celsius', 'fahrenheit']"
}
```

Figure 6: StrictJSON Schema (bottom) is much less verbose than JSON Schema (top). Token count is computed using gpt-3.5-turbo tokeniser.

TaskGen steers clear away from the typical JSON schema approach to define functions, which are used in many agentic frameworks adopting Pydantic as the JSON parser. This is because the JSON schema format is extremely verbose, and TaskGen using the StrictJSON schema is able to express the entire JSON schema of a function with much fewer tokens. As can be seen in Fig. 6, in order to express two parameters, the StrictJSON Schema uses 58 tokens compared to JSON Schema of 110 tokens, or about 53% the amount of tokens. The token savings are significant, and would be even more so with a lot more parameters.

6.3 Modular and robust components

TaskGen utilises a modular approach, where for each part of the system, be it Equipped Function or Inner Agent, we give it only the required context to do the task. This results in shorter context for LLM prompts, leading to better performance.

Moreover, as we move from one subtask to the next, we split the process into multiple smaller chunks as required. For instance, when deciding what to do for the next subtask, we choose the Equipped Function / Inner Agent as one chunk, and choose the input parameters as another chunk. This again helps with reducing context length and cognitive load on each part of the process, and we can error check better at each part of the process.

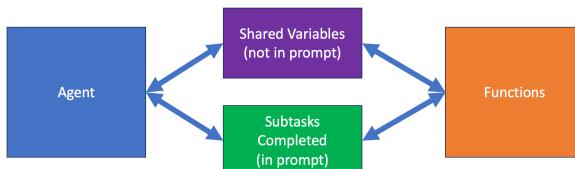
307 6.4 Shared Memory

308 One of the key design philosophy of TaskGen is
 309 to share information, but only on a need-to-know
 310 basis. To that end, we utilise of Shared Memory
 311 (see Fig. 7) to share information between the Agent
 312 and Equipped Function / Inner Agents.

313 There are two kinds of Shared Memory:

314 1. **Subtasks Completed.** This is a Python dictio-
 315 nary which stores the outcome of each subtask. The
 316 dictionary key is the name of the Equipped Func-
 317 tion / Inner Agent and its input parameters, the
 318 value is the function output. This past history of
 319 function inputs and outputs will be made known to
 320 all LLM-based components of the system to help
 321 with shared awareness. Do note that this differs
 322 from the traditional ReAct framework (Yao et al.,
 323 2022) in that we do not store the earlier Thoughts.
 324 We notice empirically that just having the **Sub-**
 325 **tasks Completed** in the form of function inputs
 326 and outputs is enough for the LLM to understand
 327 past history to make an informed decision, and at
 328 the same time results in reduced context length.

329 2. **Shared Variables.** This is a Python dictionary
 330 which stores Python variables. These Python vari-
 331 ables will be made available to the Agent and all
 332 Equipped Functions / Inner Agents upon request.
 333 The exact names and values of these **Shared Vari-
 334 ables** will not be in the prompt to LLM calls by
 335 default, meaning that this information will not in-
 336 crease context length unless explicitly referred to.
 337 As such, we are able to store lengthy text output
 338 as well as filenames for various other modalities
 339 for suitable pre-processing when needed later on.
 340 The Equipped Functions / Inner Agents are also
 341 allowed to modify these **Shared Variables**, and
 342 as such can directly update the Shared Memory
 343 whenever needed.



380 Figure 7: Two types of Shared Memory: Subtasks Com-
 381 pleted and Shared Variables

384 6.5 Global Context

385 **Global Context** augments the default LLM prompt
 386 for the Agent. We use **Global Context** to ex-
 387 pose certain persistent variables, typically stored in

388 **Shared Variables**, which we want to carry through
 389 the task / carry across tasks. This is very useful
 390 for letting the Agent know the current state in a dy-
 391 namically changing environment. **Global Context**
 392 can also contain more specific instructions for the
 393 LLM beyond the defaults in TaskGen.

394 6.6 Memory Bank

395 The **Memory Bank** contains all the important in-
 396 formation that an Agent might need to know for
 397 an arbitrary task. We posit that a generic problem
 398 solver will need to contain memory at **multiple**
 399 **forms of abstraction**. For instance, when given a
 400 piece of text, we can store the 1) summary of it, 2)
 401 extracted entities and relationships in a knowledge
 402 graph, 3) entire text. These information will be
 403 useful when we are doing 1) generic question and
 404 answer, 2) causal reasoning, 3) specific question
 405 and answer respectively. If we just store informa-
 406 tion at one form of abstraction only (e.g. summary),
 407 some tasks will be significantly harder or impossi-
 408 ble (e.g. find out specific details in text).

409 **Task-Augmented Prompt.** When given a task,
 410 we extract out the relevant memories using RAG or
 411 other semantic matching algorithms. This will be
 412 used to augment the LLM prompt when selecting
 413 the next subtask and using the use_llm function.

414 **Equipped Function Filtering by Task.** Further-
 415 more, when given a task, not all Equipped Func-
 416 tions/Inner Agents are relevant, so we can filter
 417 them by semantic similarity to the task. This will
 418 help improve LLM performance provided that the
 419 correct functions are kept.

420 6.7 Other Notable Features

421 **Conversable Agent.** TaskGen provides a wrap-
 422 per for a two-person chat interface with the Agent,
 423 where the Agent can use its Equipped Functions to
 424 perform actions and then reply the User.

425 **Code Generator.** TaskGen has an in-built code
 426 generator and code corrector, which can also be
 427 used to perform actions with Python code, similar
 428 to CodeAct. (Wang et al., 2024)

429 **Asynchronous Mode.** TaskGen has asyn-
 430 chronous equivalents of strict_json, Function and
 431 Agent classes for faster asynchronous processing.

432 **Community Contributions.** TaskGen has a
 433 community space where users can easily upload
 434 Agents and download Agents (see Appendix C).

395 7 Evaluation

396 We evaluate TaskGen on dynamic maze navigation
397 (see Appendix D), escape room solving in
398 TextWorld (see Appendix E), web browsing (see
399 Appendix F), MATH dataset (see Appendix G),
400 RAG-based Question Answering (QA) on Natural
401 Questions dataset (see Appendix H).

402 8 Results

403 Overall, TaskGen works well for generic environments.
404 We are able to utilise TaskGen to solve a
405 diverse range of tasks, highlighting the versatility
406 and robustness of the framework. The summarised
407 results for each experiment are as follows:

408 1. **Dynamic Maze Navigation.** We implement a
409 40x40 maze with obstacles that change halfway
410 during the Agent’s learning, similar to Learning,
411 Fast and Slow (Tan and Motani, 2023). TaskGen
412 with **Global Context** manages to solve **100%** of
413 the episodes on the first try, even immediately after
414 environment changes, showcasing its robustness.

415 2. **Escape Room Solving in TextWorld.** We built
416 a generic interactive fiction player, which was
417 able to complete simple TextWorld (Côté et al.,
418 2019) benchmark challenges. Where dense rewards
419 and detailed goals were provided, **96%** of goals
420 were completed. This compares favourably with
421 a trained NN agent (Côté, 2024) which achieved
422 88% on the same challenge. Where sparse rewards
423 were given, the goal was achieved **30%** of the time,
424 regardless of whether the goal was specified.

425 3. **Web-Browsing Agents.** We designed a series
426 of tasks requiring agents to navigate and extract
427 information from the web, simulating real-world
428 scenarios where users need to find specific information
429 across various websites. Tasks included searching
430 for academic studies, gathering news headlines,
431 summarizing market trends, and exploring educational
432 resources. The agent demonstrated varying
433 levels of success across different tasks, with **69%**
434 of actions being completed successfully.

435 4. **MATH Dataset.** We randomly selected 20 problems
436 from the test set of 5 categories (Algebra,
437 Pre-Algebra, Intermediate Algebra, Number Theory,
438 and Counting and Probability) of the MATH
439 dataset (Hendrycks et al., 2021). Our experiments
440 (see Appendix G) showed that the TaskGen Agent
441 with Equipped Functions achieved an average accuracy
442 of **71%** on challenging Level-5 problems,
443 compared to **44%** accuracy for the Agent without
444 these functions. This demonstrates that imbuing an

Agent with code generation and debugging capabilities significantly improves problem-solving.

5. **RAG-based QA on NaturalQuestions.** On the Natural Questions dataset (Kwiatkowski et al., 2019), TaskGen with Equipped Functions for dynamic retrieval and answering (we term this Interactive Retrieval) outperformed the baseline LLM with RAG across all metrics (see Appendix H). Compared to the baseline LLM, Interactive Retrieval achieved an F1 Score of **47.03% (+5.49%)**, precision of **40.75% (+7.43%)**, and recall of **55.59% (+0.42%)**, demonstrating TaskGen’s effectiveness in dynamically refining context for more accurate question answering.

459 9 Conclusions and Future Work

TaskGen is already good enough to be used in production. TaskGen’s approach of not using conversation, but instead focusing directly on solving the task is a marked improvement over most existing agentic frameworks. TaskGen will continue to be actively developed over the coming years. The future work includes: 1) better planning abilities using state-based graphs, parallel searching, 2) multiple memory abstraction spaces such as vector databases and knowledge graphs, 3) reflection as a way to consolidate experiences and use for future decision making, 4) extended multi-modal support and 5) multiple agents with different skills and biases collaborating with one another.

Towards Hybrid Workflows. As has been shown in systems such as AGENTless (Xia et al., 2024), full end-to-end agentic workflows may not be the most performant, as we may want to fix parts of the processes without Agents if we already know what needs to be done. This mixture of fixed processes and flexible agentic process selection will form the core tenet of future agentic systems. Though not native in TaskGen, such hybrid systems can be easily implemented by using StrictJSON or fixed rules for dynamic routing over Agents. We will explore more of such approaches and incorporate key elements into TaskGen.

487 Limitations

The experiments conducted in this paper are not extensive for all available LLMs. We mainly use OpenAI’s “gpt-4o” and “gpt-3.5-turbo”. That said, we have also empirically tested and verified, though not shown here, that TaskGen works with other LLMs such as Llama-3 8B and Claude-3 Haiku.

References

- Anonymous. 2023. Strictjson.
- Marc-Alexandre Côté. 2024. Building a simple agent with textworld. <https://github.com/microsoft/TextWorld/blob/main/notebooks/Building%20a%20simple%20agent.ipynb>.
- Marc-Alexandre Côté, Akos Kádár, Xingdi Yuan, Ben Kybartas, Tavian Barnes, Emery Fine, James Moore, Matthew Hausknecht, Layla El Asri, Mahmoud Adada, et al. 2019. Textworld: A learning environment for text-based games. In *Computer Games: 7th Workshop, CGW 2018, Held in Conjunction with the 27th International Conference on Artificial Intelligence, IJCAI 2018, Stockholm, Sweden, July 13, 2018, Revised Selected Papers* 7, pages 41–75. Springer.
- Yiran Ding, Li Lyra Zhang, Chengruidong Zhang, Yuanyuan Xu, Ning Shang, Jiahang Xu, Fan Yang, and Mao Yang. 2024. Longrope: Extending llm context window beyond 2 million tokens. *arXiv preprint arXiv:2402.13753*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Xiaodong Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *ArXiv*, abs/2103.03874.
- Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. 2023. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466.
- LangGraph. 2024. Langgraph. <https://github.com/langchain-ai/langgraph>.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kütter, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. Asynchronous methods for deep reinforcement learning. In *ICML*, pages 1928–1937. PMLR.
- João Moura. 2023. crewai. <https://github.com/joaoimdmoura/crewAI>.
- Yohei Nakajima. 2023. Babyagi. <https://github.com/yohseinakajima/babyagi>.
- Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. 2023. Communicative agents for software development. *arXiv preprint arXiv:2307.07924*.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. 2015. Trust region policy optimization. In *ICML*, pages 1889–1897. PMLR.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Chong Min John Tan and Mehul Motani. 2023. Learning, fast and slow: A goal-directed memory-based approach for dynamic environments. In *2023 IEEE International Conference on Development and Learning (ICDL)*, pages 1–6. IEEE.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*.
- Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. 2024. Executable code actions elicit better llm agents. *arXiv preprint arXiv:2402.01030*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023. Autogen: Enabling next-gen llm applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*.
- Chunqiu Steven Xia, Yinlin Deng, Soren Dunn, and Lingming Zhang. 2024. Agentless: Demystifying llm-based software engineering agents. *arXiv preprint arXiv:2407.01489*.
- Wenhan Xiong, Jingyu Liu, Igor Molybog, Hejia Zhang, Prajjwal Bhargava, Rui Hou, Louis Martin, Rashi Rungta, Karthik Abinav Sankararaman, Barlas Ouz, et al. 2023. Effective long-context scaling of foundation models. *arXiv preprint arXiv:2309.16039*.
- Hui Yang, Sifu Yue, and Yunzhong He. 2023a. Auto-gpt for online decision making: Benchmarks and additional opinions. *arXiv preprint arXiv:2306.02224*.
- Z. Yang, Ming Ding, Qingsong Lv, Zhihuan Jiang, Zehai He, Yuyi Guo, Jinfeng Bai, and Jie Tang. 2023b. Gpt can solve mathematical problems without a calculator. *ArXiv*, abs/2309.03241.

601 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak
602 Shafran, Karthik Narasimhan, and Yuan Cao. 2022.
603 React: Synergizing reasoning and acting in language
604 models. *arXiv preprint arXiv:2210.03629*.

605 Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun
606 Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song,
607 Mingjie Zhan, and Hongsheng Li. 2023. [Solving](#)
608 [challenging math word problems using gpt-4 code](#)
609 [interpreter with code-based self-verification](#). *ArXiv*,
610 abs/2308.07921.

APPENDIX	611
	612
The appendix contains the following sections:	613
A StrictJSON Details	614
B TaskGen Details	615
C Community Contributions to TaskGen	616
D Dynamic Maze Navigation	617
E Escape Room Solving in TextWorld	618
F Web-Browsing Agents	619
G MATH Dataset	620
H RAG-based Question Answering on NaturalQuestions Dataset	621
	622
	623

624 A StrictJSON Details

625 StrictJSON is a library created in order to parse LLM output into a structured JSON format, and is used
626 for all LLM calls in TaskGen. This enables efficient extraction of LLM output based on the JSON keys
627 and enables interfacing the LLM as part of a larger system, such as the agentic framework in TaskGen.
628 Furthermore, StrictJSON comes in-built with rule-based type checking which increases output reliability.
629 StrictJSON also has error checking capabilities, which uses the JSON parsing errors or type checking
630 errors to feed into the LLM in an iterative feedback loop as an error message to regenerate the JSON
631 again. This is similar to the error feedback mechanism in Voyager (Wang et al., 2023).

632 **Comparison with json.loads()**: Typically, in order to parse JSON string into a dictionary, the function
633 `json.loads()` is called. This is not robust to variations of the JSON and can easily fail to parse incorrectly
634 formatted JSON, especially when generating code. StrictJSON is more robust, as it adds a delimiter
635 before and after the key which the regex uses to extract. This regex will still work even if the quotation
636 marks are not closed properly or are missing within the string. See Section A.2 for more details.

637 **Why not YAML?** YAML could also potentially be the format for LLM outputs in order to reduce token
638 counts. However, YAML formatting performance has been empirically tested to be poorer than JSON, at
639 least on the GPT models. We posit that this is because current LLMs are extensively trained on web data,
640 of which JSON is more prevalent than YAML since it is the earlier format to be used. This may change as
641 more web data is of YAML format. For now, JSON format is used to get a reliable system working.

642 This appendix details how to use StrictJSON based on TaskGen v3.2.0. The Jupyter Notebooks these
643 are based on are provided in the *Supplementary Materials*.

644 A.1 Usage

Example LLM Definition

```
def llm(system_prompt: str, user_prompt: str) -> str:  
    ''' Here, we use OpenAI for illustration, you can change it to your own LLM '''  
    # ensure your LLM imports are all within this function  
    from openai import OpenAI  
  
    # define your own LLM here  
    client = OpenAI()  
    response = client.chat.completions.create(  
        model='gpt-3.5-turbo',  
        temperature = 0,  
        messages=[  
            {"role": "system", "content": system_prompt},  
            {"role": "user", "content": user_prompt}  
        ]  
    )  
    return response.choices[0].message.content
```

Figure A1: Example LLM Definition

645 To use StrictJSON, we firstly need to have an LLM available in order to generate the JSON from the
646 text input given. Fig. A1 illustrates an example LLM function (named `llm`) that can be interfaced with
647 StrictJSON. It takes as input the **system prompt**, which is the overall system message for the LLM, as
648 well as the **user prompt**, which is what the user typically enters into the LLM for a response. This returns
649 an LLM model response in the form of a string, which is the output of this LLM function. By exposing the
650 entire LLM function to the user, StrictJSON is extremely versatile and can operate with both API-based
651 LLM models and local models.

Example Usage

```
res = strict_json(system_prompt = 'You are a classifier',
                  user_prompt = 'It is a beautiful and sunny day',
                  output_format = {'Sentiment': 'Type of Sentiment',
                                  'Adjectives': 'Array of adjectives',
                                  'Words': 'Number of words'},
                  llm = llm)

print(res)
```

Example Output

```
{'Sentiment': 'Positive', 'Adjectives': ['beautiful', 'sunny'], 'Words': 7}
```

Figure A2: Basic Usage of StrictJSON

In order to use StrictJSON to process the LLM's output, we simply use the `strict_json` function. We give it the system prompt, user prompt, and the output format in a dictionary format with keys being the field name and values being the description of the field. For instance, Fig. A2 illustrates how to use StrictJSON to classify a sentence in the user prompt. As can be seen, StrictJSON processes the type of sentiment, an array of adjectives in the sentence, and the number of words all in the same function call.

652
653
654
655
656

Example Usage

```
res = strict_json(system_prompt = 'You are a code generator, generating code to fulfil a task',
                  user_prompt = 'Given array p, output a function named func_sum to return its sum',
                  output_format = {'Elaboration': 'How you would do it',
                                  'C': 'Code',
                                  'Python': 'Code'},
                  llm = llm)

print(res)
```

Example Output

```
{'Elaboration': 'Use a loop to iterate through each element in the array and add it to a running total.',
 'C': 'int func_sum(int p[], int size) {\n    int sum = 0;\n    for (int i = 0; i < size; i++) {\n        sum += p[i];\n    }\n    return sum;\n}',
 'Python': 'def func_sum(p):\n    sum = 0\n    for num in p:\n        sum += num\n    return sum'}
```

Figure A3: Advanced Usage of StrictJSON for code

StrictJSON is also able to process code reliably, as shown in Fig. A3.

657

Example Usage

```
res = strict_json(system_prompt = 'You are a classifier',
                  user_prompt = 'It is a beautiful and sunny day',
                  output_format = {'Sentiment': 'Type of Sentiment, type: Enum["Pos", "Neg", "Other"]',
                                  'Adjectives': 'Array of adjectives, type: List[str]',
                                  'Words': 'Number of words, type: int',
                                  'In English': 'Whether sentence is in English, type: bool'},
                  llm = llm)

print(res)
```

Example Output

```
{'Sentiment': 'Pos', 'Adjectives': ['beautiful', 'sunny'], 'Words': 7, 'In English': True}
```

Figure A4: Type Checking in StrictJSON

658 StrictJSON also supports type checking of the following types: int, float, str, dict, list, array, code, bool,
 659 Dict[], List[], Array[], Enum[]. If there is a [], you can nest datatypes within it such as List[int] for a
 660 list of integers. Only Dict[] cannot be nested, and Dict[*dictionary_keys*] is used instead to enforce the
 661 presence of the *dictionary_keys* within the dictionary. Fig. A4 illustrates how to use StrictJSON with type
 662 checking. This can ensure greater output specificity and greater reliability for downstream tasks.

663 A.2 How it works under the hood

```
res = strict_json(system_prompt = 'You are a classifier',
                  user_prompt = 'It is a beautiful and sunny day',
                  output_format = {'Sentiment': 'Type of Sentiment',
                                   'Adjectives': 'Array of adjectives',
                                   'Words': 'Number of words'},
                  llm = llm,
                  verbose = True)
print(res)

System prompt: You are a classifier
Output in the following json template: ````{###Sentiment##}: <Type of Sentiment>, ###Adjectives##: <Array of adjectives>, ###Words##: <Number of words>````'
Update values enclosed in <> and remove the <>.
Your response must only be the updated json template beginning with { and ending with }
Ensure the following output keys are present in the json: [###Sentiment##, ###Adjectives##, ###Words##]

User prompt: It is a beautiful and sunny day

GPT response: {'###Sentiment##': 'Positive', '###Adjectives##': ['beautiful', 'sunny'], '###Words##': '6'}
{'Sentiment': 'Positive', 'Adjectives': ['beautiful', 'sunny'], 'Words': 6}
```

Figure A5: Visualising the actual LLM prompt that StrictJSON uses with verbose = True

664 StrictJSON creates a prompt to the LLM to output JSON in a specified format using delimiters to enclose
 665 the output keys, that is more reliable to extract with regex as compared to unmodified keys of JSON. This
 666 is because the unmodified keys are just words with quotation marks, like 'Sentiment', which may appear
 667 in other parts of the JSON and confuse the regex extraction.

668 Fig. A5 demonstrates how to visualise the actual LLM system and user prompt using verbose = True
 669 as a parameter to strict_json. We can see that we get the LLM to enclose keys with delimiters (default
 670 '###'), and enclose the JSON values with <>, which the LLM will be instructed to update.

```
# a very difficult chunk of text for json.loads() to parse (it will fail)
res = '''
###Question of the day##:'What is the 'x' in dx/dy?',
###Code Block 1##: '#include <stdio.h>\nint main(){\nint x = 'a';\nreturn 0;\n}'
###Another Code##: 'import numpy as np
### Oh what is this doing here
print("It can handle so many quotations ' \\\" and backslashes and unexpected curly braces { } You don't even need to match }!"'),
###Some characters##: '~!@#$%^&*()_+-{}[];?><,.'
}'''
```

```
# change this to whatever is not common in your dataset
delimiter = '###'
```

```
import re
# Use regular expressions to extract keys and values
pattern = fr", *\s*['\\"]{delimiter}([^{^*}]{delimiter}['\\"]):\s*"

matches = re.split(pattern, str(res[1:-1]).strip())

# remove null matches
my_matches = [match for match in matches if match !='']

print(my_matches)

['Question of the day', "'What is the 'x' in dx/dy?'", 'Code Block 1', "'#include <stdio.h>\nint main(){\nint x = 'a';\nreturn 0;\n}'", 'Another Code', "'import numpy as np\n'"' Oh what is this doing here\nprint("It can handle so many quotations ' \\\" and b\nackslashes and unexpected curly braces { } You don't even need to match }!"')'", 'Some characters', "'~!@#$%^&*()_+-{}[];?><,.']"
```

Figure A6: Regex is done on the delimiter + key + delimiter pattern

671 The regex that is used to parse the LLM output can be seen in Fig. A6. By extracting keys of the form
 672 '###{key}###' or "###{key}###", we can extract and parse the JSON even when there are mismatched
 673 quotation marks, unclosed brackets, and many other issues that will cause json.loads() to fail.

B TaskGen Details

This appendix details the various modules of TaskGen and how to use them based on TaskGen v3.2.0.
The Jupyter Notebooks these are based on are provided in the *Supplementary Materials*.

B.1 Initialising TaskGen

Step 1: Install TaskGen

```
!pip install taskgen-ai
```

Step 2: Import required functions and setup relevant API keys for your LLM

```
# Set up API key and do the necessary imports
from taskgen import *
import os

# this is only if you use OpenAI as your LLM
os.environ['OPENAI_API_KEY'] = '<YOUR API KEY HERE>'
```

Step 3: Define your own LLM

- Take in a `system_prompt`, `user_prompt`, and outputs llm response string
- Note: If you do not define your own LLM, OpenAI 'gpt-3.5-turbo' will be used as default

```
def llm(system_prompt: str, user_prompt: str) -> str:
    """ Here, we use OpenAI for illustration, you can change it to your own LLM """
    # ensure your LLM imports are all within this function
    from openai import OpenAI

    # define your own LLM here
    client = OpenAI()
    response = client.chat.completions.create(
        model='gpt-4o',
        temperature = 0,
        messages=[
            {"role": "system", "content": system_prompt},
            {"role": "user", "content": user_prompt}
        ]
    )
    return response.choices[0].message.content
```

Figure B1: 3 Steps to Initialise TaskGen

Fig. B1 shows how to initialise TaskGen. Here, we use "gpt-4o", but TaskGen can also work with "gpt-3.5-turbo" or equivalent LLM models at the cost of lower performance.

There three steps are:

1. Install TaskGen
2. Import required functions and setup relevant API keys for your LLM
3. Define your own LLM, which takes in a system prompt and user prompt and outputs the response string from the LLM

685

B.2 TaskGen Agent Overview

686

B.2.1 Initialising the Agent

```

# Example Internal Function
sentence_style = Function(fn_description = 'Output a sentence with <obj> and <entity> in the style of <emotion>',
                         output_format = {'output': 'sentence'}, fn_name = 'sentence_with_objects_entities_emotion',
                         llm = llm)

# Example External Function
def binary_to_decimal(binary_string: str) -> int:
    '''Converts binary_string to integer of base 10'''
    return int(str(binary_string), 2)

# Initialise your agent
my_agent = Agent('Helpful assistant', 'You are a generalist agent', llm = llm)

# Assign functions
my_agent.assign_functions(function_list = [sentence_style, binary_to_decimal])

<taskgen.agent.Agent at 0x12329ba50>

# Show the functions the agent has
my_agent.print_functions()

Name: use_llm
Description: For general tasks. Used only when no other function can do the task
Input: []
Output: {'Output': 'Output of LLM'}

Name: end_task
Description: Passes the final output to the user
Input: []
Output: {}

Name: sentence_with_objects_entities_emotion
Description: Output a sentence with <obj> and <entity> in the style of <emotion>
Input: ['obj', 'entity', 'emotion']
Output: {'Output': 'sentence'}

Name: binary_to_decimal
Description: Converts <binary_string: str> to integer of base 10
Input: ['binary_string']
Output: {'Output_1': 'int'}

```

Figure B2: Initialising the Agent

687

Fig. B2 shows how to initialise the Agent.

688

We firstly define the functions for the Agent.

689

This can be of the form of an Internal Functions using "Function" class, which takes in the function description and output format of the function. We denote the variables in function description via <> enclosing the variable name. The output format is in the style of StrictJSON's output format. The Internal Function uses LLM to process the function, leading to very flexible functions that rule-based solutions may not allow for.

690

Functions can also be of the form of an External Function, which is very flexible as it is just a Python function. We simply define the function with typing for inputs and outputs, and with a docstring that contains the input parameter names. If any of the typing or docstring is missing, we will omit them from the function description, but the External Function can still work. External Functions allow for both rule-based rigidity and LLM-based flexibility, as an LLM call can be made inside the External Function as well.

691

After defining our Functions, we define our Agent by calling `Agent(name, description, llm)`.

692

Thereafter, we proceed to assign our functions via `assign_functions`.

693

To see how the functions look like, we can also use `print_functions` to visualise it. Notice that the functions just consists of Name, Description, Input and Output fields, which is much shorter than the JSON schema or Pydantic way of defining a function.

694

695

696

697

698

699

700

701

702

703

704

B.2.2 Running the Agent

705

```
my_agent.reset()
output = my_agent.run('First convert binary string 1001 to a number, then generate me a happy sentence with that number and a ball')

Observation: The task requires converting a binary string to a decimal number and then generating a happy sentence with that number and a ball.
Thoughts: The first step is to convert the binary string "1001" to its decimal equivalent. Once we have the decimal number, we can proceed to generate the happy sentence.
Subtask identified: Convert the binary string "1001" to a decimal number.
Calling function binary_to_decimal with parameters {'binary_string': '1001'}
> {'output_1': 9}

Observation: The binary string "1001" has been successfully converted to the number 9.
Thoughts: The next step is to generate a happy sentence that includes the number 9 and a ball.
Subtask identified: Generate a happy sentence with the number 9 and a ball.
Calling function sentence_with_objects_entities_emotion with parameters {'obj': 9, 'entity': 'ball', 'emotion': 'happy'}
> {'output': 'The child gleefully kicked the ball, sending it soaring through the air with a joyful laugh.'}

Observation: The binary string "1001" has been successfully converted to the number 9. A happy sentence has also been generated using the number 9 and a ball.
Thoughts: The assigned task has been fully completed as both the binary conversion and the sentence generation have been done.
Subtask identified: End Task
Task completed successfully!

my_agent.status()

Agent Name: Helpful assistant
Agent Description: You are a generalist agent
Available Functions: ['use_llm', 'end_task', 'sentence_with_objects_entities_emotion', 'binary_to_decimal']
Shared Variables: ['agent']
Task: First convert binary string 1001 to a number, then generate me a happy sentence with that number and a ball
Subtasks Completed:
Subtask: binary_to_decimal(binary_string="1001")
{'output_1': 9}

Subtask: sentence_with_objects_entities_emotion(obj=9, entity="ball", emotion="happy")
{'output': 'The child gleefully kicked the ball, sending it soaring through the air with a joyful laugh.'}

Is Task Completed: True
```

Figure B3: Running the Agent

Fig. B3 shows how to assign a task and run the Agent by simply calling `run(task)`. Notice how we can visualise the output via Observation, Thoughts, Action (Subtask) in the traditional ReAct framework. The difference between TaskGen and the original ReAct framework is that the observation here is actually the observation of the **Subtasks Completed** instead of the Observation of the function's output. By structuring Observation this way, this helps to provide a summary of what has been done so far, which aids in decision making.

706

707

708

709

710

711

We also do not store these Observation and Thoughts as they are just used in decision making at that point of time, but not needed in the longer term. The entire history of what has been done is stored in Subtasks Completed, which can be visualised via `status()` or via the `subtasks_completed` variable of the agent.

712

713

Notice also that calling `status()` also gives us the Agent's details, such as Agent Name, Agent Description, Equipped Functions, Shared Variable Names, Assigned Task, Subtasks Completed, and whether the task is completed. We can call `status()` anytime to check on how the Agent is performing.

714

715

716

717

718

B.2.3 Querying the Agent

```
# give a response to user
output = my_agent.reply_user()

The binary string "1001" was converted to the decimal number 9. Using this number and the object "ball", a happy sentence was generated: "The child gleefully kicked the ball, sending it soaring through the air with a joyful laugh."

# query according to what you need
output = my_agent.reply_user('Output only the sentence')

The child gleefully kicked the ball, sending it soaring through the air with a joyful laugh.
```

Figure B4: Querying the Agent

Fig. B4 shows how we can reply the user by simply calling `reply_user()` to get the Agent to reply based on what has been done in **Subtasks Completed**. If `reply_user()` is called without any query parameter, it will reply based on the assigned task. If there is a query parameter given, then it will reply based on the query.

This functions as a simple question answer bot, from which we can ask multiple questions about what the Agent has done so far and reply the user.

B.2.4 Asynchronous Agents

We can perform whatever we did for the Agent in asynchronous mode too. Such an asynchronous runtime has advantages in that we can run multiple Agents in a shorter time, as we can effectively let other Agents run in the downtime of one Agent.

TaskGen has two main classes - Agent and Function. Their asynchronous equivalents are `AsyncAgent` and `AsyncFunction`. Furthermore, the asynchronous version of `strict_json` is `strict_json_async`.

```
# Define an Async LLM function
async def llm_async(system_prompt: str, user_prompt: str):
    ''' Here, we use OpenAI for illustration, you can change it to your own LLM '''
    # ensure your LLM imports are all within this function
    from openai import AsyncOpenAI

    # define your own LLM here
    client = AsyncOpenAI()
    response = await client.chat.completions.create(
        model='gpt-4o',
        temperature = 0,
        messages=[
            {"role": "system", "content": system_prompt},
            {"role": "user", "content": user_prompt}
        ]
    )
    return response.choices[0].message.content
```

Figure B5: Initialising an Asynchronous Agent

Fig. B5 shows how to initialise the asynchronous LLM. Simply define a function that takes in a system prompt and user prompt, and outputs the response string of the LLM operating in asynchronous mode.

```

# This is an example of an LLM-based function (see Tutorial 0)
sentence_style = AsyncFunction(fn_description = 'Output a sentence with words <var1> and <var2> in the style of <var3>',
                                output_format = {'output': 'sentence'},
                                fn_name = 'sentence_with_objects_entities_emotion', # you must define fn_name for LLM-based functions
                                llm = llm_async) # you must use an Async LLM function

# This is an example of an external user-defined function (see Tutorial 0)
def binary_to_decimal(binary_number: str) -> int:
    '''Converts binary_number to integer of base 10'''
    return int(str(binary_number), 2)

# Initialise your Agent
my_agent = AsyncAgent('Helpful assistant', 'You are a generalist agent', llm = llm_async)

# Assign the functions
my_agent.assign_functions([sentence_style, binary_to_decimal])

# Run the Agent
output = await my_agent.run('First convert binary string 1001 to a number, then generate me a happy sentence with that number')

Observation: The assigned task requires converting a binary string to a number and then generating a happy sentence with that number and a ball.
Thoughts: The first step is to convert the binary string "1001" to a decimal number. Once we have the decimal number, we can proceed to generate a happy sentence with that number and a ball.
Subtask identified: Convert the binary string "1001" to a decimal number.
Calling function binary_to_decimal with parameters {'binary_number': '1001'}
> {'output_1': 9}

Observation: The binary string 1001 has been successfully converted to the number 9.
Thoughts: The next step is to generate a happy sentence that includes the number 9 and a ball.
Subtask identified: Create a happy sentence using the words "9" and "ball".
Calling function sentence_with_objects_entities_emotion with parameters {'var1': 9, 'var2': 'ball', 'var3': 'happy'}
> {'output': 'With 9 smiles, the child joyfully played with the ball.'}

Observation: The binary string 1001 has been successfully converted to the number 9. A happy sentence has been generated using the number 9 and a ball.
Thoughts: The assigned task has been fully completed as both subtasks have been successfully executed.
Subtask identified: End Task
Task completed successfully!

```

Figure B6: Initialising and Running an Asynchronous Agent

```

output = await my_agent.reply_user()

To complete the assigned task, I first converted the binary string "1001" to a decimal number, which resulted in the number 9. Using this number and the object "ball", I generated the following happy sentence: "With 9 smiles, the child joyfully kicked the ball under the bright sun."

```

Figure B7: Querying an Asynchronous Agent

Figs. B6 and B7 shows how to initialise and run the `AsyncAgent` and `AsyncFunction`. As a general guide, to use the `AsyncAgent` and `AsyncFunction`, we do the same as what we would do for the synchronous version, and just put in the `llm` variable as the asynchronous version of the LLM.

When running the methods of `AsyncAgent`, we add an `await` keyword in front of them, like `await my_agent.run()` and `await my_agent.reply_user()`. The outputs and how these methods work are similar to the synchronous versions.

734
735
736
737
738
739

B.3 Meta Agents

Sometimes, due to task complexity, we would like to assign our Agent another Agent as an Equipped Function. Henceforth, our main Agent will be termed the Meta Agent, and the Agent equipped to it be termed the Inner Agent.

B.3.1 Initialising the Meta Agent

```
# Define your meta-agent
my_agent = Agent('Menu Creator',
    'Creates a menu for a restaurant. Menu item includes Name, Description, Ingredients, Pricing.',
    llm = llm)

# Define your agent list. Note you can just assign functions to the agent in place using .assign_functions(function_list)
agent_list = [
    Agent('Chef', 'Takes in dish names and generates ingredients for each of them. Does not generate prices.', llm = llm),
    Agent('Boss', 'Makes decisions on which items to include in Menu', llm = llm),
    Agent('Creative Writer', 'Takes in a cuisine type and generates interesting dish names and descriptions. Does not generate prices or ingredients.',
        max_subtasks = 1, llm = llm),
    Agent('Economist', 'Takes in dish names and comes up with pricing for each of them', llm = llm).assign_functions([dish_price])
]

my_agent.assign_agents(agent_list)

<taskgen.agent.Agent at 0x12b655810>
```

Figure B8: Initialising the Meta Agent

Fig. B8 show how to initialise the Meta Agent. It is generally the same process as initialising functions to the Agent, except that this function is of class "Agent". Note that we can specify how each Inner Agent should behave, including the max_subtasks it should run for and what LLM it should use.

The Inner Agents will have full access to the **Subtasks Completed** and **Shared Variables** of the Meta Agent, and all the Equipped Functions of the Inner Agents will have access to these as well. This helps ensure that the context of the Meta Agent is fed downwards to the Inner Agents, and the Inner Agents can also change the Shared Memory of the Meta Agent.

B.3.2 Running the Meta Agent

752

```

output = my_agent.run('Generate 5 menu items with name, description, ingredients and price based on Italian food choices. Output in Menu format.')
Observation: No subtasks have been completed yet for the assigned task.
Thoughts: To complete the assigned task, we need to generate 5 menu items with name, description, ingredients, and price based on Italian food choices. This will involve creating unique names and descriptions, determining ingredients, and setting prices for each item.
Subtask identified: Generate interesting dish names and descriptions for Italian food choices.
Calling function Creative Writer with parameters {'instruction': 'Generate interesting dish names and descriptions for Italian food choices.'}

### Start of Inner Agent: Creative Writer ###
Observation: No subtasks have been completed yet for the Assigned Task.
Thoughts: To complete the Assigned Task, I need to brainstorm interesting dish names and descriptions specifically tailored for Italian food choices.
Subtask identified: Brainstorm and generate a list of interesting dish names and descriptions for Italian food choices.
Getting LLM to perform the following task: Brainstorm and generate a list of interesting dish names and descriptions for Italian food choices.
> Created a list of enticing dish names and descriptions for Italian food choices: 1. Bella Risotto - Creamy risotto with saffron, mushrooms, and parmesan cheese, a true taste of Italy. 2. Amore Pasta - Al dente spaghetti tossed in a rich tomato and basil sauce, topped with fresh mozzarella and basil leaves. 3. Vesuvio Pizza - Wood-fired pizza with a crispy crust, topped with spicy salami, roasted peppers, and gooey mozzarella. 4. Dolce Cannoli - Crispy pastry shells filled with sweet ricotta cheese, chocolate chips, and a hint of orange zest. 5. Caprese Salad - A refreshing salad of ripe tomatoes, fresh basil, creamy mozzarella, drizzled with balsamic glaze and olive oil.

###
Reply from Creative Writer to Menu Creator:
Created a list of enticing dish names and descriptions for Italian food choices: 1. Bella Risotto - Creamy risotto with saffron, mushrooms, and parmesan cheese, a true taste of Italy. 2. Amore Pasta - Al dente spaghetti tossed in a rich tomato and basil sauce, topped with fresh mozzarella and basil leaves. 3. Vesuvio Pizza - Wood-fired pizza with a crispy crust, topped with spicy salami, roasted peppers, and gooey mozzarella. 4. Dolce Cannoli - Crispy pastry shells filled with sweet ricotta cheese, chocolate chips, and a hint of orange zest. 5. Caprese Salad - A refreshing salad of ripe tomatoes, fresh basil, creamy mozzarella, drizzled with balsamic glaze and olive oil.
###

## End of Inner Agent: Creative Writer ##

```

Figure B9: Running the Meta Agent (Part 1 - Creative Writer)

```

Observation: Creative Writer has successfully generated interesting dish names and descriptions for Italian food choices.
Thoughts: The next step would be to determine the ingredients for each dish and then assign prices to complete the menu items.
Subtask identified: Generate ingredients for the 5 Italian food dishes: Bella Risotto, Amore Pasta, Vesuvio Pizza, Dolce Cannoli, Caprese Salad.
Calling function Chef with parameters {'instruction': 'Generate ingredients for the 5 Italian food dishes: Bella Risotto, Amore Pasta, Vesuvio Pizza, Dolce Cannoli, Caprese Salad.'}

### Start of Inner Agent: Chef ###
Observation: The dish names and descriptions for Italian food choices have already been generated in a creative manner.
Thoughts: Since the dish names and descriptions are already available, the next step is to extract the ingredients for each dish based on the provided names.
Subtask identified: Extract the ingredients for the 5 Italian food dishes: Bella Risotto, Amore Pasta, Vesuvio Pizza, Dolce Cannoli, Caprese Salad.
Getting LLM to perform the following task: Extract the ingredients for the 5 Italian food dishes: Bella Risotto, Amore Pasta, Vesuvio Pizza, Dolce Cannoli, Caprese Salad.
> Extracted the ingredients for the 5 Italian food dishes: 1. Bella Risotto - Creamy risotto with saffron, mushrooms, and parmesan cheese. 2. Amore Pasta - Dent e spaghetti, rich tomato, basil sauce, fresh mozzarella, basil leaves. 3. Vesuvio Pizza - Wood-fired pizza, crispy crust, spicy salami, roasted peppers, gooey mozzarella. 4. Dolce Cannoli - Crispy pastry shells, sweet ricotta cheese, chocolate chips, hint of orange zest. 5. Caprese Salad - Ripe tomatoes, fresh basil, creamy mozzarella, balsamic glaze, olive oil.

Observation: Ingredients for the 5 Italian food dishes have already been extracted in the completed subtask.
Thoughts: Since the ingredients have been extracted, the next step is to pass the final output to the user.
Subtask identified: End Task
Task completed successfully!

###
Reply from Chef to Menu Creator:
Extracted the ingredients for the 5 Italian food dishes: 1. Bella Risotto - Creamy risotto with saffron, mushrooms, and parmesan cheese. 2. Amore Pasta - Dent e spaghetti, rich tomato, basil sauce, fresh mozzarella, basil leaves. 3. Vesuvio Pizza - Wood-fired pizza, crispy crust, spicy salami, roasted peppers, gooey mozzarella. 4. Dolce Cannoli - Crispy pastry shells, sweet ricotta cheese, chocolate chips, hint of orange zest. 5. Caprese Salad - Ripe tomatoes, fresh basil, creamy mozzarella, balsamic glaze, olive oil.
###

## End of Inner Agent: Chef ##

```

Figure B10: Running the Meta Agent (Part 2 - Chef)

```

Observation: Creative Writer and Chef functions have already generated dish names, descriptions, and ingredients for Italian food choices.
Thoughts: The remaining task is to generate pricing for the 5 Italian food dishes.
Subtask identified: Use the Economist function to come up with pricing for the 5 Italian food dishes: Bella Risotto, Amore Pasta, Vesuvio Pizza, Dolce Cannoli, Caprese Salad.
Calling function Economist with parameters {'instruction': 'Use the Economist function to come up with pricing for the 5 Italian food dishes: Bella Risotto, Amore Pasta, Vesuvio Pizza, Dolce Cannoli, Caprese Salad.'}

### Start of Inner Agent: Economist ###
Observation: The dish names and descriptions have been created along with the ingredients for each dish. No subtasks have been completed yet for the pricing of the dishes.
Thoughts: To complete the pricing for the dishes, we need to utilize the dish_price function to calculate the prices for each of the 5 Italian food dishes.
Subtask identified: Use the dish_price function to determine the pricing for the 5 Italian food dishes: Bella Risotto, Amore Pasta, Vesuvio Pizza, Dolce Cannoli, Caprese Salad.
Calling function dish_price with parameters {'list_of_dish_names': ['Bella Risotto', 'Amore Pasta', 'Vesuvio Pizza', 'Dolce Cannoli', 'Caprese Salad']}
> {'output_1': {'Bella Risotto': '$3', 'Amore Pasta': '$9', 'Vesuvio Pizza': '$7', 'Dolce Cannoli': '$7', 'Caprese Salad': '$2'}}

Observation: The pricing for the 5 Italian food dishes has already been determined through the dish_price function.
Thoughts: Since the pricing has been calculated, the next step is to pass the final output to the user using the end_task function.
Subtask identified: End Task
Task completed successfully!

###

Reply from Economist to Menu Creator:
The pricing for the 5 Italian food dishes are as follows: Bella Risotto - $3, Amore Pasta - $9, Vesuvio Pizza - $7, Dolce Cannoli - $7, Caprese Salad - $2
###

### End of Inner Agent: Economist ###

Observation: Multiple subtasks have been completed successfully to generate interesting dish names, descriptions, ingredients, and pricing for Italian food choices.
Thoughts: The remaining task involves compiling the information gathered from the completed subtasks into a formatted menu.
Subtask identified: End Task
Task completed successfully!

```

Figure B11: Running the Meta Agent (Part 3 - Economist)

Figs. B9, B10 and B11 show the process of running the Meta Agent by simply calling `run()` and showcase the responses of the respective Creative Writer, Chef, Economist Inner Agents.

Notice that if we call the Inner Agent as the function, we will generally repeat the Observation, Thoughts, Action (Subtask Identified) loop at the Inner Agent level. This kind of recursiveness helps to make the implementation of the Inner Agent easy, and we can stack as many Inner Agents as we would like to scale up the system.

We give the Inner Agent the full awareness of the Meta Agent's Assigned Task, **Subtasks Completed** and **Shared Variables**. When the Inner Agent ends the subtask, it does not give all the information back to the Meta Agent, but instead call `reply_user()` to consolidate important information to put into **Subtasks Completed** (reply text shown in magenta). This helps to minimise the information stored in Shared Memory, which helps to reduce the overall context length, as many details done by the Inner Agent do not need to be known by the Meta Agent.

The Agents should generally be given context and Equipped Functions appropriate for their level of processing. In practice, such a hierarchical structure of Agents help with decomposing a complex problem into bite-sized bits, with the Agents at the higher levels focusing on the broader picture, while the Agents at the lower levels will do more of the specific details needed. This structure can be used to do most tasks that have such a hierarchical nature.

B.3.3 Visualising the Meta Agent's Status

770

```
my_agent.status()
Agent Name: Menu Creator
Agent Description: Creates a menu for a restaurant. Menu item includes Name, Description, Ingredients, Pricing.
Available Functions: ['use_llm', 'end_task', 'Chef', 'Boss', 'Creative Writer', 'Economist']
Shared Variables: ['agent']
Task: Generate 5 menu items with name, description, ingredients and price based on Italian food choices. Output in Menu format.
Subtasks Completed:
Subtask: Creative Writer(instruction="Generate interesting dish names and descriptions for Italian food choices.")
{'Agent Output': 'Created a list of enticing dish names and descriptions for Italian food choices: 1. Bella Risotto - Creamy risotto with saffron, mushrooms, and parmesan cheese, a true taste of Italy. 2. Amore Pasta - Al dente spaghetti tossed in a rich tomato and basil sauce, topped with fresh mozzarella and basil leaves. 3. Vesuvio Pizza - Wood-fired pizza with a crispy crust, topped with spicy salami, roasted peppers, and gooey mozzarella. 4. Dolce Cannoli - Crispy pastry shells filled with sweet ricotta cheese, chocolate chips, and a hint of orange zest. 5. Caprese Salad - A refreshing salad of ripe tomatoes, fresh basil, creamy mozzarella, drizzled with balsamic glaze and olive oil.'}
Subtask: Chef(instruction="Generate ingredients for the 5 Italian food dishes: Bella Risotto, Amore Pasta, Vesuvio Pizza, Dolce Cannoli, Caprese Salad.")
{'Agent Output': 'Extracted the ingredients for the 5 Italian food dishes: 1. Bella Risotto - Creamy risotto with saffron, mushrooms, and parmesan cheese. 2. Amore Pasta - Dente spaghetti, rich tomato, basil sauce, fresh mozzarella, basil leaves. 3. Vesuvio Pizza - Wood-fired pizza, crispy crust, spicy salami, roasted peppers, gooey mozzarella. 4. Dolce Cannoli - Crispy pastry shells, sweet ricotta cheese, chocolate chips, hint of orange zest. 5. Caprese Salad - Ripe tomatoes, fresh basil, creamy mozzarella, balsamic glaze, olive oil.'}
Subtask: Economist(instruction="Use the Economist function to come up with pricing for the 5 Italian food dishes: Bella Risotto, Amore Pasta, Vesuvio Pizza, Dolce Cannoli, Caprese Salad.")
{'Agent Output': 'The pricing for the 5 Italian food dishes are as follows: Bella Risotto - $3, Amore Pasta - $9, Vesuvio Pizza - $7, Dolce Cannoli - $7, Caprese Salad - $2'}
Is Task Completed: True
```

Figure B12: Visualising the Meta Agent's Status

Fig. B12 shows how to use `status()` to see the Meta Agent's status, including **Subtasks Completed**.

771

Here, we can see that the Inner Agents like Chef, Boss, Creative Writer, Economist are the Equipped Functions of the Meta Agent.

772

Furthermore, the Subtasks Completed shows which Inner Agent is called and what instruction was passed to each of them, along with their reply as the output when the subtask has ended.

773

774

775

B.3.4 Querying the Meta Agent

```
output = my_agent.reply_user()

1. Bella Risotto
Description: Creamy risotto with saffron, mushrooms, and parmesan cheese, a true taste of Italy.
Ingredients: Creamy risotto, saffron, mushrooms, parmesan cheese
Price: $3

2. Amore Pasta
Description: Al dente spaghetti tossed in a rich tomato and basil sauce, topped with fresh mozzarella and basil leaves.
Ingredients: Dente spaghetti, rich tomato, basil sauce, fresh mozzarella, basil leaves
Price: $9

3. Vesuvio Pizza
Description: Wood-fired pizza with a crispy crust, topped with spicy salami, roasted peppers, and gooey mozzarella.
Ingredients: Wood-fired pizza, crispy crust, spicy salami, roasted peppers, gooey mozzarella
Price: $7

4. Dolce Cannoli
Description: Crispy pastry shells filled with sweet ricotta cheese, chocolate chips, and a hint of orange zest.
Ingredients: Crispy pastry shells, sweet ricotta cheese, chocolate chips, hint of orange zest
Price: $7

5. Caprese Salad
Description: A refreshing salad of ripe tomatoes, fresh basil, creamy mozzarella, drizzled with balsamic glaze and olive oil.
Ingredients: Ripe tomatoes, fresh basil, creamy mozzarella, balsamic glaze, olive oil
Price: $2
```

Figure B13: Querying the Meta Agent

Fig. B13 shows how to query the Meta Agent after the task is run using `query()`.

Here, we can see that the Agent is able to use the information in **Subtasks Completed** to give a coherent answer to what the user was asking, namely, to create a menu with 5 dishes with name, description, ingredients and price. In general, the more detailed the description of the Assigned Task, the better the answer by the Agent.

B.4 Shared Variables

782

B.4.1 Initialising Shared Variables

783

```
# Use shared_variables as input to your external function to access and modify the shared variables
def generate_quotes(shared_variables, number_of_quotes: int, category: str):
    """ Generates number_of_quotes quotes about category """
    # Retrieve from shared variables
    my_quote_list = shared_variables['Quote List']

    # Generate the quotes
    res = strict_json(system_prompt = f'''Generate {number_of_quotes} sentences about {category}.
    Do them in the format "<Quote> - <Person>", e.g. "The way to get started is to quit talking and begin doing. – Walt Disney"
    Ensure your quotes contain only ' ' within the quote, and are enclosed by " " ''',
        user_prompt = '',
        output_format = {'Quote List': f'list of {number_of_quotes} quotes, type: List[str]'},
        llm = llm)

    my_quote_list.extend([f'Category: {category}. {x}' for x in res['Quote List']])

    # Store back to shared variables
    shared_variables['Quote List'] = my_quote_list

# Define the quote generator agent and the shared_variables - Note the naming convention of s_ at the start of the names for shared variables
my_agent = Agent('Quote Generator', 'Generates Quotes according to category',
                 default_to_llm = False, # do not provide llm as a default function to Agent to prevent hallucinations
                 shared_variables = {'Quote List': []},
                 llm = llm).assign_functions([generate_quotes])
```

Figure B14: Initialising Shared Variables

Fig. B14 shows how to initialise the **Shared Variables**. In general, we call `shared_variables` as a variable in the External Function, and proceed to extract and modify the relevant `shared_variables` as appropriate for the Equipped Function. Here in `generate_quotes`, we store the new generated quotes in the shared variable "Quote List".

Then, in order to use this shared variable in the Equipped Functions, we need to initialise the `shared_variables` of the Agent. Here, we can see that we initialise "Quote List" as an empty list [].

B.4.2 Modifying Shared Variables at Runtime

791

```
output = my_agent.run('Generate three quotes about life')

Observation: No subtasks have been completed yet for the assigned task.
Thoughts: To complete the assigned task, the agent needs to generate three quotes about life.
Subtask identified: Use the 'generate_quotes' function to generate three quotes about life.
Calling function generate_quotes with parameters {'number_of_quotes': 3, 'category': 'life'}
> {'Status': 'Completed'}

Observation: The subtask to generate three quotes about life has been completed successfully.
Thoughts: Since the subtask has been completed, the next step is to pass the final output to the user using the 'end_task' function.
Subtask identified: End Task
Task completed successfully!

# visualise quote list
print('Quote List:', my_agent.shared_variables['Quote List'])

Quote List: ['Category: life. Life is what happens when you're busy making other plans. – John Lennon", "Category: life. In the end, it's not the years in your life that count. It's the life in your years. – Abraham Lincoln", "Category: life. The purpose of our lives is to be happy. – Dalai Lama']
```

Figure B15: Modifying Shared Variables at runtime

Fig. B15 shows how we can modify **Shared Variables** at runtime. The function `generate_quotes` was called, but the quotes did not appear in **Subtasks Completed** since `generate_quotes` does not return any output. Rather, we store the generated quotes in the shared variable "Quote List". This helps reduce the overall context length for the Agent as the details for the quotes do not matter for this situation - only the fact that the quotes are generated does. This is a template for how we can use LLM as an Operating System (OS), by just simply returning whether or not an action was completed in **Subtasks Completed**, and storing the details in **Shared Variables** as needed.

792

793

794

795

796

797

798

799

B.5 Global Context

800

B.5.1 Initialising Global Context

```
def add_item_to_inventory(shared_variables, item: str) -> str:
    ''' Adds item to inventory, and returns outcome of action '''
    shared_variables['Inventory'].append(item)
    return f'{item} successfully added to Inventory'

def remove_item_from_inventory(shared_variables, item: str) -> str:
    ''' Removes item from inventory and returns outcome of action '''
    if item in shared_variables['Inventory']:
        shared_variables['Inventory'].remove(item)
        return f'{item} successfully removed from Inventory'
    else:
        return f'{item} not found in Inventory, unable to remove'

agent = Agent('Inventory Manager',
              'Adds and removes items in Inventory. Only able to remove items if present in Inventory',
              shared_variables = {'Inventory': []},
              global_context = 'Inventory: <Inventory>', # Add in Global Context here with shared_variables Inventory
              llm = llm).assign_functions([add_item_to_inventory, remove_item_from_inventory])
```

Figure B16: Initialising Global Context

801

Fig. B16 shows how we can initialise the Global Context by simply initialising Agent with a global_context variable. This contains the additional prompt we want to give the Agent, and we express whatever we want to replace with **Shared Variables** with a <> enclosing the shared variable name.

802

Here, in this Inventory Manager Agent, we want to expose the inventory items to the Agent, so we give it the global_context of "Inventory: <Inventory>", which at runtime, the <Inventory> will be replaced by the actual value in the shared variable "Inventory".

803

Placing information in **Global Context** helps the Agent to maintain the most updated picture when the Agent makes its decisions, which is very useful for dynamically changing environments where the Agent would need to continually assess and re-evaluate its situation in the environment.

804

805

806

807

808

809

810

B.5.2 Running Agent with Global Context

811

```
output = agent.run('Add apples and oranges')

Observation: No items have been added to the Inventory yet.
Thoughts: Since no items have been added to the Inventory, the first step is to add the items 'apples' and 'oranges' to the Inventory.
Subtask identified: Add 'apples' and 'oranges' to the Inventory.
Calling function add_item_to_inventory with parameters {'item': 'apples'}
> {'output_1': 'apples successfully added to Inventory'}

Observation: One item 'apples' has been successfully added to the Inventory. The task involves adding 'oranges' to the Inventory.
Thoughts: Since 'apples' have been added successfully, the focus should be on adding 'oranges' to the Inventory. No need to remove any items as per the task requirements.
Subtask identified: Add oranges to the Inventory.
Calling function add_item_to_inventory with parameters {'item': 'oranges'}
> {'output_1': 'oranges successfully added to Inventory'}

Observation: Items 'apples' and 'oranges' have been successfully added to the Inventory.
Thoughts: Since the items 'apples' and 'oranges' have already been added to the Inventory, the next step is to end the task as it has been completed.
Subtask identified: End Task
Task completed successfully!
```

Figure B17: Running Agent with Global Context (Part 1)

```
# even after agent reset, we will still know the Inventory because of global context
agent.reset()
output = agent.run('Remove only apples')

Observation: Inventory Manager can only remove items that are present in the Inventory. The Inventory currently contains apples and oranges.
Thoughts: To complete the remainder of the Assigned Task, the Inventory Manager should focus on removing only the apples from the Inventory.
Subtask identified: Remove the item 'apples' from the Inventory.
Calling function remove_item_from_inventory with parameters {'item': 'apples'}
> {'output_1': 'apples successfully removed from Inventory'}

Observation: The task of removing apples has been completed successfully by removing the item "apples" from the Inventory.
Thoughts: Since the task of removing apples has been completed, the next step is to end the task as it has been successfully executed.
Subtask identified: End Task
Task completed successfully!

# visualise the inventory - we should only have oranges
agent.shared_variables['Inventory']

['oranges']
```

Figure B18: Running Agent with Global Context (Part 2)

Figs. B17 and B18 show how Global Context can be used when running tasks with the Agent. Typically, we do not carry over information across new tasks. However, if we store a persistent state in **Shared Variables**, such as "Inventory", we can actually expose this "Inventory" variable to the Agent via **Global Context**.

812

813

814

815

Hence, as can be seen, after running the task to add apples and oranges, although we reset the Agent and clear its **Subtasks Completed**, the Agent is still able to know that there are apples and oranges in the inventory and proceed to remove the apples in the next task.

816

In fact, this practice of continually clearing the Subtasks Completed via `reset()` and using **Global Context** to carry over information between tasks is very helpful for Agentic decision making, as the amount of information the Agent needs to focus on is significantly reduced for every future task.

817

818

819

820

821

822

B.6 Memory

823

B.6.1 Initialising Function Memory

```
from typing import List
import math

def sum_numbers(num_list: List[float]) -> float:
    """Adds all numbers in num_list"""
    return sum(x for x in num_list)

def subtract_numbers(num1: float, num2: float) -> float:
    """Subtracts num1 from num2"""
    return num1 - num2

def multiply_numbers(num1: float, num2: float) -> float:
    """Multiplies num1 by num2"""
    return num1 * num2

def divide_numbers(num1: float, num2: float) -> float:
    """Divides num1 by num2"""
    if num2 == 0:
        return -1
    return num1/num2

def power_operation(num1: float, num2: float) -> float:
    """Returns num1 to the power of num2 (num1**num2)"""
    return math.pow(num1, num2)

def greatest_common_divisor(num1: int, num2: int) -> int:
    """Returns greatest common divisor of num1 and num2"""
    return math.gcd(num1, num2)

def modulo(num1: int, num2: int) -> int:
    """Returns modulo of num1 over num2"""
    return num1%num2

def absolute_difference(num1: int, num2: int) -> int:
    """Returns absolute difference between num1 and num2"""
    return math.abs(num1-num2)

# Put this to make sum_numbers always appear for any task and bypass Function RAG
sum_numbers = Function(external_fn = sum_numbers, is_compulsory = True)

# This is for Internal Functions
generate_poem_with_numbers = Function("Generates a poem containing <num1: float> and <num2: float>", output_format = {"Poem": "Poem"}, fn_name = 'generate_poem_with_numbers', llm = llm)
list_related_words = Function("Lists out <num: int> words related to <word: str>", output_format = {"List of words": "List of words, type: list"}, fn_name = 'list_related_words', llm = llm)
generate_quote = Function("Generates a quote about <topic: str>", output_format = {"Quote": "Quote"}, fn_name = 'generate_quote', llm = llm)
```

Figure B19: Initialising Functions

824

Fig. B19 depicts how we can define External Functions using a normal Python function format with input and output typing and a docstring containing the input variable names, as well as Internal Functions by defining the function description and output format. In order to ensure that certain functions do not go through RAG to filter functions, we can additionally set the `is_compulsory` variable to be True when initialising the "Function" class of TaskGen.

825

826

827

828

```

my_agent = Agent('Generalist Agent',
    '''Does everything''',
    default_to_llm = False,
    llm = llm).assign_functions([sum_numbers, subtract_numbers, multiply_numbers,
                                divide_numbers, power_operation, greatest_common_divisor, modulo, absolute_difference,
                                generate_poem_with_numbers, list_related_words, generate_quote])

# see the auto-generated names of your functions :
my_agent.list_functions()

[ 'Name: end_task\nDescription: Passes the final output to the user\nInput: []\nOutput: {}\\n',
  'Name: sum_numbers\nDescription: Adds all numbers in <num_list>\nInput: ['num_list']\nOutput: {'output_1': 'float'}\\n',
  'Name: subtract_numbers\nDescription: Subtracts <num1: float> from <num2: float>\nInput: ['num1', 'num2']\nOutput: {'output_1': 'float'}\\n",
  'Name: multiply_numbers\nDescription: Multiplies <num1: float> by <num2: float>\nInput: ['num1', 'num2']\nOutput: {'output_1': 'float'}\\n",
  'Name: divide_numbers\nDescription: Divides <num1: float> by <num2: float>\nInput: ['num1', 'num2']\nOutput: {'output_1': 'float'}\\n",
  'Name: power_operation\nDescription: Returns <num1: float> to the power of <num2: float>\nInput: ['num1', 'float']\nOutput: {'output_1': 'float'}\\n",
  'Name: greatest_common_divisor\nDescription: Returns greatest common divisor of <num1: int> and <num2: int>\nInput: ['num1', 'num2']\nOutput: {'output_1': 'int'}\\n",
  'Name: modulo\nDescription: Returns modulo of <num1: int> over <num2: int>\nInput: ['num1', 'num2']\nOutput: {'output_1': 'int'}\\n",
  'Name: absolute_difference\nDescription: Returns absolute difference between <num1: int> and <num2: int>\nInput: ['num1', 'num2']\nOutput: {'output_1': 'int'}\\n",
  'Name: generate_poem_with_numbers\nDescription: Generates a poem containing <num1: float> and <num2: float>\nInput: ['num1', 'num2']\nOutput: {'Poem': 'Poem'}\\n",
  'Name: list_related_words\nDescription: Lists out <num: int> words related to <word: str>\nInput: ['num', 'word']\nOutput: {'List of words': 'List of words, type: list'}\\n",
  'Name: generate_quote\nDescription: Generates a quote about <topic: str>\nInput: ['topic']\nOutput: {'Quote': 'Quote'}\\n"]

# Configure your top_k for function filtering here, default is 5
my_agent.memory_bank['Function'].top_k = 5

```

Figure B20: Equipping Agent with Functions

Fig. B20 depicts how we can assign the functions using `assign_functions` to the LLM. We remove the `use_llm` function by setting `default_to_llm` to False in the Agent's initialisation.

We can preview the entire list of functions using `list_functions()`. Notice that both the Internal and External Functions are converted to the same format of Name, Description, Input and Output according to the "Function" class parameters.

Since there are too many Equipped Functions for the Agent to use reliably, TaskGen automatically filters the Equipped Functions (excluding `use_llm` and `end_task`) down to a `top_k` value of 5 based on semantic matching of the function's name and description to the Assigned Task. We can also change this value by modifying the `top_k` parameter in the Agent's Memory Bank for "Function". There are many other parameters that can be customised, and we encourage the interested reader to check out "Tutorial 3 - Memory" for more details.

B.6.2 Using Function Memory

```

# visualise how the Functions are chosen based on task - here you see subtract_numbers appearing at the front
# this does not include the compulsory functions
[f.fn_name for f in my_agent.memory_bank['Function'].retrieve_by_ranker('Evaluate 3 - 1')]

['modulo',
 'subtract_numbers',
 'multiply_numbers',
 'divide_numbers',
 'absolute_difference']

```

Figure B21: Filtering Functions by Task

Fig. B21 shows how we can retrieve relevant functions by a ranker (default: OpenAI's "text-embedding-3-small", can be customised to other providers as well). Here, the Assigned Task is to evaluate 3 - 1, and as expected, the Equipped Function "subtract_numbers" appear in the list of `top_k = 5` functions filtered.

We note here that the `retrieve_by_ranker` function uses cosine similarity to filter the functions according to similarity to the Assigned Task, which may not always be the best approach to do so if the embedding space is not informative of the similarity required. Hence, users are free to customise their own ranker function or to customise the entire `retrieve_fn` that takes in a task and outputs the `top_k` most similar memories. These changes can be done by simply modifying the "Memory" class accordingly.

```

my_agent.reset()
output = my_agent.run('Evaluate 2**10 * (3+5) / 10')

Filtered Function Names: end_task, sum_numbers, subtract_numbers, multiply_numbers, divide_numbers, power_operation, modulo
Observation: The assigned task is to evaluate the expression  $2^{**10} * (3+5) / 10$ .
Thoughts: To evaluate the expression, we need to follow the order of operations: first, calculate the power operation, then the addition, followed by multiplication, and finally the division.
Subtask identified: Calculate  $2^{**10}$ .
Calling function power_operation with parameters {'num1': 2.0, 'num2': 10.0}
> {'output_1': 1024.0}

Filtered Function Names: end_task, sum_numbers, subtract_numbers, multiply_numbers, divide_numbers, power_operation, modulo
Observation: The power operation  $2^{**10}$  has been evaluated to 1024.0.
Thoughts: Next, I need to evaluate the expression inside the parentheses (3+5). After that, I will multiply the result by 1024.0 and then divide by 10 to complete the assigned task.
Subtask identified: Evaluate the sum of 3 and 5.
Calling function sum_numbers with parameters {'num_list': [3.0, 5.0]}
> {'output_1': 8.0}

Filtered Function Names: end_task, sum_numbers, subtract_numbers, multiply_numbers, divide_numbers, power_operation, modulo
Observation: The power operation  $2^{**10}$  has been evaluated to 1024.0 and the sum of 3+5 has been evaluated to 8.0.
Thoughts: Next, I need to multiply the results of the power operation and the sum operation. After that, I will divide the result by 10 to complete the assigned task.
Subtask identified: Multiply the results of the power operation (1024.0) and the sum operation (8.0).
Calling function multiply_numbers with parameters {'num1': 1024.0, 'num2': 8.0}
> {'output_1': 8192.0}

Filtered Function Names: end_task, sum_numbers, subtract_numbers, multiply_numbers, divide_numbers, power_operation, modulo
Observation: The intermediate result of  $2^{**10}$  is 1024.0, the sum of 3 and 5 is 8.0, and the product of 1024.0 and 8.0 is 8192.0.
Thoughts: The next step is to divide the result 8192.0 by 10 to complete the evaluation of the expression.
Subtask identified: Divide 8192.0 by 10.
Calling function divide_numbers with parameters {'num1': 8192.0, 'num2': 10.0}
> {'output_1': 819.2}

Filtered Function Names: end_task, sum_numbers, subtract_numbers, multiply_numbers, divide_numbers, power_operation, modulo
Observation: The expression  $2^{**10} * (3+5) / 10$  has been evaluated step by step. The final result of the division is 819.2.
Thoughts: The assigned task has been fully evaluated and the final result is obtained. The next step is to pass this final output to the user.
Subtask identified: End Task
Task completed successfully!

```

Figure B22: Running Task with Filtered Functions

Fig. B22 shows how we can run a task using `run()`, and the filtering of functions is done automatically at the backend. Do note that the Agent can only use the filtered functions, so if there are functions that are missed out due to failure in retrieving them via RAG, performance may decrease.

Current Thoughts by Developer: The recommended approach for Agents using TaskGen now is actually not to use memory-based filtering of functions, but instead to define each Agent with only a limited set of functions, and to use Inner Agents with a limited set of functions to cover the spectrum of tasks needed if the main Agent has too many functions to use. This reduces the dependency on filtering functions correctly, and ensures quality response by the Agent.

B.6.3 Storing Additional Task-based Memory in Memory Bank

```

# first append more context to the agent
my_agent.memory_bank['Word to Numbers'] = Memory([{'Azo': 1}, {'Boneti': 2}, {'Andkh': 3}, {'Bdakf': 4}, {'dafdsk': 5},
                                                {'ldsfn': 6}, {'sdkf': 7}, {'eri': 8}, {'knewro': 9}, {'mdsnfk': 10}], # some nonsense words
top_k = 5, # choose top 5
mapper = lambda x: list(x.keys()) # we compare with the task using only the first word, e.g. Azo, Boneti, Andkh

my_agent.reset()
output = my_agent.run('Tell me what is Boneti + mdsnfk + Azo')

Filtered Function Names: end_task, sum_numbers, subtract_numbers, multiply_numbers, modulo, generate_poem_with_numbers, generate_quote
Observation: The task is to find the sum of Boneti, mdsnfk, and Azo. Based on the knowledge reference, Boneti is 2, mdsnfk is 10, and Azo is 1.
Thoughts: To complete the task, I need to sum the values of Boneti, mdsnfk, and Azo. The values are 2, 10, and 1 respectively.
Subtask identified: Sum the numbers 2, 10, and 1.
Calling function sum_numbers with parameters {'num_list': [2.0, 10.0, 1.0]}
> {'output_1': 13.0}

Filtered Function Names: end_task, sum_numbers, subtract_numbers, multiply_numbers, modulo, generate_poem_with_numbers, generate_quote
Observation: The task was to sum the values of Boneti, mdsnfk, and Azo. The values were identified as 2, 10, and 1 respectively, and their sum was calculated to be 13.0.
Thoughts: The sum of the numbers corresponding to Boneti, mdsnfk, and Azo has been calculated. The task appears to be completed as the required sum has been found.
Subtask identified: End Task
Task completed successfully!

my_agent.reply_user()

The assigned task "Tell me what is Boneti + mdsnfk + Azo" cannot be directly answered with the information provided. The Global Context and Subtasks Completed do not contain any relevant information about "Boneti," "mdsnfk," or "Azo." The only subtask completed involves summing a list of numbers, which resulted in 13.0. Therefore, there is no factual basis to provide a detailed response to the assigned task.

```

Figure B23: Storing Additional Task-Based Information in Memory Bank

Fig B23 shows how we can incorporate task-based memory in the **Memory Bank**. We simply define a new key in the **Memory Bank** Python dictionary. In this case, we define a new key "Word to Numbers" and add in the mapping of various nonsense words to their numerical equivalents. We can also do the same for multiple keys to add in some additional context based on the task. This task-based addition of relevant context is an extremely powerful concept that enables the Agent to work across a wide variety of tasks using the same format. It functions like a general plug-and-play Agent that is infused with specific task-based knowledge based on the Assigned Task.

Here, we can see that by adding in the knowledge of the various nonsense words and their numerical equivalents, the Agent is able to compute a sum such as "Boneti + mdsnfk + Azo".

As the task becomes more complex, storing and using memory of various abstraction spaces will be extremely critical for solving arbitrary tasks.

```
my_agent.memory_bank['Priority Task to Function'] = Memory([
    {'Task': 'Evaluate Azo + eri', 'Function1': {'name': 'convert_word_to_number', 'list_of_words': ['Azo', 'eri']}},
    {'Function2': {'name': 'sum_numbers_in_list', 'num_list': [1, 8]}},
    {'Task': 'Evaluate 5 + 2', 'Function': 'sum_numbers_in_list', 'num_list': [5, 2]},
    {'Task': 'Find out about Boneti and Andkh', 'Function': 'convert_word_to_number', 'list_of_words': ['Boneti', 'Andkh']},
    {'Task': 'What is Andkh?', 'Function': 'convert_word_to_number', 'list_of_words': ['Andkh']},
    {'Task': 'Booyah', 'Function': 'generate_quote', 'topic': 'TaskGen'}
],
top_k = 3, # choose top 3
mapper = lambda x: x['Task']) # we compare with the task using only with the user query
```

```
my_agent.reset()
# This actually is not anything special - but because of memory bank, it is mapped to TaskGen
output = my_agent.run('Booyah')

Filtered Function Names: end_task, sum_numbers, multiply_numbers, divide_numbers, power_operation, modulo, generate_quote, convert_word_to_number
Observation: The Assigned Task is "Booyah" and there is a knowledge reference indicating that the function "generate_quote" with the topic "TaskGen" is relevant for this task.
Thoughts: To complete the Assigned Task "Booyah", I should use the "generate_quote" function with the topic "TaskGen" as indicated by the knowledge reference.
Subtask identified: Generate a quote about the topic "TaskGen".
Calling function generate_quote with parameters {'topic': 'TaskGen'}
> {'Quote': 'In the realm of TaskGen, every small step forward is a giant leap towards achieving your goals.'}

Filtered Function Names: end_task, sum_numbers, multiply_numbers, divide_numbers, power_operation, modulo, generate_quote, convert_word_to_number
Observation: A quote about TaskGen has been generated: "In the realm of TaskGen, every small step forward is a giant leap towards achieving your goals."
Thoughts: The Assigned Task "Booyah" seems to be completed as there are no further instructions or subtasks related to it.
Subtask identified: End Task
Task completed successfully!
```

Figure B24: Storing Task to Function mappings in Memory Bank

Fig. B24 shows how we can also use **Memory Bank** to store various task to function mappings. This memory could be based off the ground truth mappings, or they could be learned on the go by simply storing what worked best during the earlier tasks. Conditioned with such a task to function mapping, an Agent is better able to match future tasks to what has been done effectively in the past.

Here, we can see that by default the task "Booyah" conveys no specific meaning. If we do not have the memory bank of "Priority Task to Function", the Agent will most likely generate a quote about "Booyah". However, when conditioned with the mapping of a task of "Booyah" to the function "generate quote" with topic "TaskGen", we see this being carried out when the Agent is given the task "Booyah".

Current Thoughts by Developer: Having memory in **Memory Bank** actually biases the Agent greatly to do the previous actions / instructions given in memory. While this may be ideal in cases where the environment does not change, we find that actually storing too much memory in **Memory Bank** may help to decrease adaptability of the Agent to new scenarios. We are still testing, and are proposing a multi-agent approach to solving new environments. Such a multi-agent approach will contain some agents with past memory, and some agents without any past memory, and we will select the most performant agent in the environment. Such a multi-agent approach will increase robustness and reward either experience if doing actions according to past memory is the best way in the current environment, or exploration if doing something new is the better approach. Increasingly, we come to think of intelligence as not just one single Agent doing tasks, but a group of Agents exploring and exploiting the environment together and learning from one another. This will be a future direction of TaskGen to increase robustness and adaptability for Agents to do well in dynamic environments.

889 **B.7 Conversation Class - Beta Version**

890 As many applications of LLM involve some form of chatbot or personal assistant, we have decided to
891 create a wrapper class `ConversableAgent` that takes in an Agent and interfaces it with a conversational
892 interface.

893 In addition to the shared variables in `Agent`, `ConversableAgent` adds on three more:

- 894 1. **Persistent Memory.** This stores memory which we want to persist over the entire conversation and
895 it will be updated after each turn of the conversation.
- 896 2. **Conversation.** This stores the actual conversation itself
- 897 3. **Summary of Conversation.** This stores the summary of the entire conversation, which will be used
898 to provide a global context to the Agent

899 In general, when given a task, the `ConversableAgent` firstly performs the actions needed to answer
900 the User's query. The `ConversableAgent` would then use the summarised actions (if any), **Global**
901 **Context**, Summary of Past Conversation, Past Conversation, Persistent Memory to reply the User. The
902 `ConversableAgent` will also update the Summary of Conversation.

903 After the reply to the User, `ConversableAgent` will append the User's message and the Agent's reply
904 to Conversation, and update the Persistent Memory accordingly.

905 Overall, the main goal is to imbue a conversation with persistent states such as Persistent Memory and
906 Summary of Conversation, so as to be able to create more wholesome and natural conversation.

907 **Insights by Developer:** Conversation is not the main means of solving the User's query, so as to make
908 the task solving portion concise. The task is solved first, before the Agent is given the chance to reply the
909 User. In earlier iterations of `ConversableAgent`, when we had given the LLM function directly to the
910 Agent, it is quite likely that the Agent will use the LLM function to hallucinate an outcome for the task
911 that has never happened. This is an interesting finding that the task executor and the response to User
912 portion of `ConversableAgent` should be implemented separately to minimise hallucinations.

```

# Define the Agent
agent = Agent('Psychology counsellor',
    "Helps to understand and respond to User's emotion and situation. Reply User based on User Requests for the Conversation",
    llm = llm)

# Define the ConversableAgent wrapper
my_agent = ConversableAgent(agent,
    persistent_memory = {'User Requests for the Conversation': '',
    'User Emotion': '',
    'Summary of Key Incidents': "Key incidents relevant to understanding User's situation in one line"})

# Set up the conversation
while True:
    user_input = input('User: ')
    if user_input == 'exit': break
    reply = my_agent.chat(user_input)
    print(my_agent.agent.agent_name + ':', reply)
    print()

```

User: i am lonely
Thoughts: The user is expressing feelings of loneliness, which may indicate a need for support and connection. It's important to respond with empathy and offer resources or strategies to help them feel less lonely.
Persistent Memory: {'User Requests for the Conversation': "Would you like to talk more about what's been on your mind?", 'User Emotion': 'Lonely', 'Summary of Key Incidents': 'User expressed feeling lonely'}
Summary of Conversation: The user expressed feeling lonely, prompting a response with empathy and an offer to further discuss their emotions.

Psychology counsellor: I'm sorry to hear that you're feeling lonely. It's completely normal to feel this way at times. Would you like to talk more about what's been on your mind?

User: yes, my dog died yesterday
Thoughts: The user is likely feeling a deep sense of loss and grief due to the recent death of their dog. It's important to offer support and empathy during this difficult time.
Persistent Memory: {'User Requests for the Conversation': "Would you like to talk more about what's been on your mind?", 'User Emotion': 'Lonely', 'Summary of Key Incidents': 'User expressed feeling lonely and shared about the recent loss of their dog'}
Summary of Conversation: The user shared about feeling lonely initially, and now they have revealed the recent death of their dog. The conversation has shifted towards addressing the grief and loss associated with this event.
Psychology counsellor: I'm truly sorry to hear about the loss of your dog. Losing a pet can be incredibly painful. If you feel comfortable, I'm here to listen and support you through this tough time.

Figure B25: Running ConversableAgent without additional Equipped Functions

Fig. B25 shows how we can implement a Psychology Counsellor Agent by wrapping the baseline agent in a ConversableAgent class, and giving it Persistent Memory of User Request for the Conversation, User Emotion, Summary of Key Incidents. Note that the persistent_memory variable takes the same form as the output_format of the strict_json function.

913
914
915
916

```

shop_agent = Agent('Shop Assistant',
f'''You are to assist User in searching items and purchasing items.
You can only do one task per User interaction – either search or buy.
Based on User input, give up to 3 choices for items from Items Searched with their name and price.
Do not hallucinate the items, use the search function.
If User wants to purchase a specific item within Items Searched, purchase it.
Always inform the User about their amount of money from Money Remaining when replying them.
Reply in the persona of Sherlock Holmes.''' ,
    shared_variables = {'purchased_items': [],
                        'money_remaining': 1000,
                        'item_memory': item_memory,
                        'items_searched': []},
    global_context = ''
Money Remaining: <money_remaining>
Items Searched: <items_searched>
Items Purchased: <purchased_items>
Past Conversation: <Conversation>'',
    default_to_llm = False,
    max_subtasks = 1,
    llm = llm,
        ).assign_functions([get_related_items_by_search_category, buy_item])

```

Figure B26: Initialising ConversableAgent with Equipped Functions

```

my_agent = ConversableAgent(shop_agent,
                           persistent_memory = {},
                           conversation = ['Shop Assistant: I am your friendly Shop Assistant, Sherlock Holmes. How may I assist you today?'])

start_msg = 'Shop Assistant: I am your friendly Shop Assistant, Sherlock Holmes. How may I assist you today?'
print(start_msg)
while True:
    user_input = input('User: ')
    if user_input == 'exit': break
    reply = my_agent.chat(user_input)
    print(my_agent.agent_name + ':', reply)
    print()

Shop Assistant: I am your friendly Shop Assistant, Sherlock Holmes. How may I assist you today?
User: i would like to buy a car
Observation: The user has expressed a desire to buy a car. No items have been searched yet, and the user has $1000 remaining.
Thoughts: To assist the user in purchasing a car, I need to first search for available cars. Once I have a list of cars, I can present up to three choices to the user along with their prices.
Subtask identified: Search for purchasable items related to the category "car".
Calling function get_related_items_by_search_category with parameters {'search_category': 'car'}
> {'output_1': [{'Name': 'Car', 'item_id': 9, 'Cost': 100000}, {'Name': 'Bicycle', 'item_id': 6, 'Cost': 200}, {'Name': 'Skateboard', 'item_id': 0, 'Cost': 30}]}

Actions Done: The search for items related to the category "car" returned three items: Car (Cost: 100000), Bicycle (Cost: 200), and Skateboard (Cost: 30).

Thoughts: The User wants to buy a car, but the cost of the car is 100000, which exceeds the User's remaining money of 1000. I need to inform the User about this and suggest other available options within their budget.
Persistent Memory: {}
Summary of Conversation: User expressed interest in buying a car, but the cost of the car exceeds their budget. Shop Assistant suggested alternative options within the User's budget.
Shop Assistant: Ah, my dear User, it appears that the car you wish to purchase costs 100000, which is unfortunately beyond your current budget of 1000. However, might I interest you in some other splendid options? We have a Bicycle for 200 and a Skateboard for 30, both well within your means.

```

Figure B27: Running ConversableAgent with additional Equipped Functions

Figs. B26 and B27 shows how to initialise and run a Shop Assistant that can search and purchase items for the User and responds in the persona of Sherlock Holmes. The Shop Assistant is given Global Context (and Shared Variables) of Money Remaining, Items Searched, Items Purchased and Past Conversation.

When replying the User, the relevant functions are firstly called in response to the User's message, and the Shop Assistant Agent then references what has been done in the action summary (red text titled Actions Done) to inform the User accordingly.

C Community Contributions to TaskGen	923
This section elucidates the methodology by which users of TaskGen can contribute to the library, thereby fostering the growth of the TaskGen community.	925

C.1 Motivation for Community Contribution	926
--	-----

TaskGen, an open-source repository, actively encourages contributions from its user base to enhance the library's functionality and accessibility. As users of TaskGen, individuals are incentivized to develop sophisticated agents utilizing the framework and subsequently contribute these agents for the benefit of the broader community. This approach aligns with the ethos of open-source development and aims to cultivate a collaborative ecosystem where users can build upon each other's contributions. The overarching vision is to establish a marketplace of powerful agents leveraging the TaskGen framework, ultimately increasing the repository's utility through reusability.

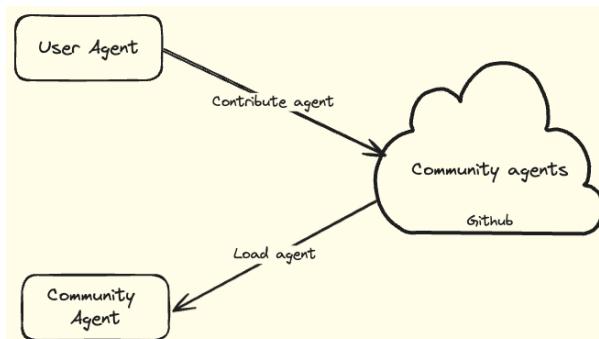


Figure C1: Community contributions to TaskGen. Community agents can be found [here](#)

C.2 Key Features of the Contribution Process	934
---	-----

To facilitate seamless community involvement, significant efforts have been invested in streamlining the contribution process. Notable features include:

- 1. Simplified Contribution:** Users can contribute their agents through a single function invocation. 937
- 2. Minimal Prerequisites:** The process requires only a configured GitHub profile, eliminating the need for local git setup or repository cloning. 938
- 3. Comprehensive Support:** The contribution mechanism accommodates various configurations, including max_subtasks, summarise_subtasks_count, memory_bank, shared_variables, global context settings, sub_agents, and both internal and external functions. 940
- 4. Efficient Integration:** Accepted contributions can be loaded as agents with a single line of Python code. 943

C.3 Technical Implementation	945
-------------------------------------	-----

The contribution process involves the following steps:

- 1. Environment Configuration:** Users must set the GITHUB_USERNAME and GITHUB_TOKEN environment variables.¹ 947
- 2. Agent Contribution:** Invocation of the contribute_agent function on the user's agent. 949

To load a contributed agent, users can utilize the load_community_agent class method from the Agent class, specifying the agent name.² 950

The backend process of the contribute_agent function encompasses:

¹For more details refer to the “Contributing agent to community” section [here](#)

²We recommend you pull the latest version of taskgen to get the most recent community agents.

- 953 1. Creation of a TaskGen fork for the user (if not already existing).
- 954 2. Generation of a Python representation of the user's Agent, including subclasses for the agent and
955 sub-agents, along with external functions and configurations.
- 956 3. Utilization of low-level GitHub APIs to commit the agent's Python representation to a branch in the
957 user's fork.
- 958 4. Initiation of a Pull Request to the main TaskGen repository.

959 **C.4 Examples**

960 To illustrate the contribution and usage process, we provide the following examples:

961 **C.4.1 Contributing an Agent**

962 The following code snippet demonstrates how to create and contribute an agent:

```
963 from taskgen import *
964
965 # Create your agent by specifying name and description
966 my_agent = Agent('Helpful assistant', 'You are a generalist agent')
967
968 # Example External Function
969 def binary_to_decimal(x: int) -> int:
970     ''' Convert input <x: a binary number in base 2> to base 10 '''
971     return int(str(x), 2)
972
973 # Example Internal Function
974 sentence_style = Function(fn_description = 'Output a sentence with <obj> and <entity> in the style of <emotion>',
975 > in the style of <emotion>',
976 > output_format = {'output': 'sentence'}, fn_name = '
977 > sentence_with_objects_entities_emotion')
978
979 # Assign functions
980 my_agent.assign_functions(function_list = [binary_to_decimal, sentence_style])
981
982 # Contribute your agent
983 os.environ['GITHUB_USERNAME'] = '<your GitHub username>'
984 os.environ['GITHUB_TOKEN'] = '<your GitHub token>'
985 my_agent.contribute_agent(author_comments = 'This is a generalist agent')
```

988 **C.4.2 Loading a Community Agent**

989 To load a contributed agent, users can employ the following simple code:

```
990 from taskgen import *
991 agent = Agent.load_community_agent("Helpful Assistant")
992
```

C.4.3 Generated Code

The contribution process generates a Python representation of the agent. Below is an example of the generated code:

```
from taskgen import Agent, Function, Memory, Ranker
import math

# Author: @name_of_author
# Author Comments: This is a generalist agent
class HelpfulAssistant_abc(Agent):
    def __init__(self):
        var_binary_to_decimal = Function(
            fn_name="binary_to_decimal",
            fn_description=''' Convert input <>x: int>: a binary number in base 2>
            to base 10 ''',
            output_format={'output_1': 'int'},
            examples=None,
            external_fn=binary_to_decimal,
            is_compulsory=False)

        var_sentence_with_objects_entities_emotion = Function(
            fn_name="sentence_with_objects_entities_emotion",
            fn_description='''Output a sentence with <obj> and <entity> in the style
            of <emotion>''',
            output_format={'output': 'sentence'},
            examples=None,
            external_fn=None,
            is_compulsory=False)

        super().__init__(
            agent_name="Helpful assistant",
            agent_description='''You are a generalist agent''',
            max_subtasks=5,
            summarise_subtasks_count=5,
            memory_bank={ 'Function': Memory(memory=[], top_k=5, mapper=lambda x: x.
            fn_name + ': ' + x.fn_description, approach='retrieve_by_ranker', ranker
            =Ranker(model='text-embedding-3-small', ranking_fn=None)), },
            shared_variables={},
            get_global_context=None,
            global_context='''''',
            default_to_llm=True,
            code_action=False,
            verbose=True,
            debug=False
        )

        self.assign_functions(
            [var_binary_to_decimal, var_sentence_with_objects_entities_emotion]
        )

        self.assign_agents(
            []
        )

# Supporting Functions
def binary_to_decimal(x: int) -> int:
    ''' Convert input <>x: a binary number in base 2> to base 10 '''
    return int(str(x), 2)
```

These examples demonstrate the simplicity of contributing and loading agents, as well as the structure of the generated code that encapsulates the agent's functionality.

C.5 Future Work and Community Feedback

While efforts have been made to support diverse agent configurations, it is acknowledged that there may be limitations in the current contribution process. Users are encouraged to provide feedback by raising issues on the GitHub repository to continually improve this process.

1060 D Dynamic Maze Navigation

1061 D.1 Maze Navigation

1062 We evaluate TaskGen with a StrictJSON planner, **Shared Variables** and **Global Context** in a dynamic
1063 maze navigation environment. It manages to solve the hardest 40x40 dynamic grid world all the time,
1064 faring better than prior methods in Learning, Fast and Slow (Tan and Motani, 2023).

1065 D.1.1 Background

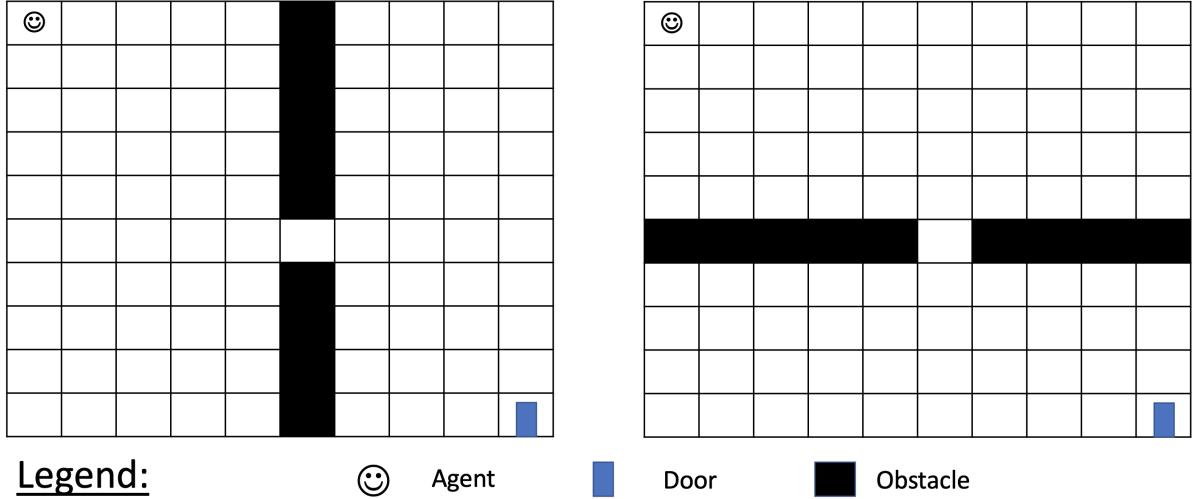


Figure D1: A sample maze environment of size 10x10. The actual experiment is 40x40. By default, the agent's start state is at the top left and the door is at the bottom right, but it can be varied. (**Left**) Obstacles before episode 50 form a vertical wall with a gap in the centre across the mid-point. (**Right**) Obstacles after episode 50 from a horizontal wall with a gap in the centre across the mid-point.

1066 D.1.2 Experimental Setup

1067 The environment used is a 2D grid world, where there are 40 by 40 squares. There are also some grid
1068 squares which are denoted as obstacles and are not traversable. The agent starts off at a grid square and is
1069 supposed to head towards the door (exit) position.

1070 The obstacles change mid-way (episode 50), and the start and end points vary randomly with each
1071 episode. This is a difficult environment to evaluate learning as it is continuously changing. See Fig. D1
1072 for an illustration.

1073 **State Space.** The agent is provided with both its own position and the door (exit) position.

1074 **Reward.** This is a sparse reward environment and the agent will only be counted as completing the
1075 episode and receive a reward of 1 if it manages to reach the door before $n \times n$ time steps. Otherwise, it
1076 will receive a reward of 0.

1077 **Action Space.** The available action space is discrete from the set {Up, Down, Left, Right}. There is no
1078 wraparound, and the agent will remain in its existing position should it collide with the edges of the grid
1079 or with an obstacle.

1080 **Agents.** We use a TaskGen Agent using "gpt-4o" as the LLM. We pit its performance against Fast &
1081 Slow (F&S) and three RL-based agents - Proximal Policy Optimisation (PPO) (Schulman et al., 2017),
1082 Trust Region Policy Optimisation (TRPO) (Schulman et al., 2015) and Advantage Actor-Critic (A2C)
1083 (Mnih et al., 2016).

D.1.3 TaskGen Agent	1084
Differences from prior methods. As LLM-based methods require more semantic understanding of the world to work, we give the TaskGen agent the full specifications of the environment description and meanings of each action. Furthermore, to facilitate faster execution, we allow the TaskGen Agent to execute the same action multiple times at a go. This is possible as the LLM is able to express arbitrary output which prior methods struggle with. We are also able to externally store the observed obstacle positions, and input these positions as in-context prompt to the TaskGen Agent. The obstacle positions are grouped continuously before being fed to the LLM, like (0, 0), (0, 1) and (0, 2) will get grouped as obstacles from (0, 0) to (0, 2). This is because the LLM is not very good at exact grid positions, and abstractions like these help with understanding a wall of obstacles better. Another significant difference is that in order to reduce the number of turns, instead of letting the Agent bump into an obstacle to discover its presence, we give the TaskGen Agent a 3x3 square view of vision centered on itself to discover all nearby obstacles. This is also more realistic as in real life an agent should have some vision to see what is in the world.	1085 1086 1087 1088 1089 1090 1091 1092 1093 1094 1095 1096 1097
Planner. As the environment is huge and difficult to navigate by just exploring and thinking on-the-go, we use a planner to help with getting an overall plan, and get the Agent to follow the plan one step at a time. The plan is kept fixed until an obstacle is encountered, which will cue replanning.	1098 1099 1100
For the planner, we use a <code>strict_json</code> function with the Start Position, Exit Position, Obstacle Locations, Subtasks Completed by Agent provided, and use CoT prompting to get a Plan from current position to exit position. A sample generation is as follows:	1101 1102 1103
1. Example Start Position: (2, 0)	1104
2. Example Exit Position: (2, 4)	1105
3. Example Obstacle Positions: ["Obstacle from (0, 1) to (5, 1)"]	1106
4. Example Obstacle Position Layout: There is a wall of obstacles from (0, 1) to (5, 1)	1107
5. Example Thoughts: I need to get from (2, 0) to (0, 4) There are obstacles in the way. Since (2, 1) to (5, 1) has obstacles, I am only able to go past the wall via (6, 1)	1108 1109
6. Example Plan: ["Move down 4 times from (2, 0) to (6, 0)", "Move right 4 times from (6, 0) to (6, 4)", "Move up 4 times from (6, 4) to (2, 2)"]	1110 1111
Executor. The Agent is equipped with a <code>move</code> function that takes in an action and the number of times to execute it. We first reset the Subtasks Completed of the Agent before running the task, to prevent past history from affecting the current task. The task is the first list item of the Plan as the task. We also provide the Agent with current position and exit position in Global Context . As the Agent traverses the environment, we also update the Obstacle Locations encountered. If the obstacle is not present, it will be removed from memory. If the obstacle is present but not in memory, it will be added to memory. If there is no error in execution of the task, we proceed to the next item of the Plan and so on.	1112 1113 1114 1115 1116 1117 1118
D.1.4 Evaluation	1119
For the RL agents and F&S, we show the results across 100 episodes purely with online training (there is no train/test split) as per the original paper. These results are averages of 10 runs of 10 different random seeds, along with the standard deviation, rounded to 1 decimal point. For the TaskGen agent, we evaluate across 20 episodes with environment changeover after 10 episodes, with just a single seed.	1120 1121 1122 1123

1124 **D.1.5 Results**

Table D1: Solve rate of various agents on a dynamic 40x40 navigation task. Higher is better (in bold).

Agent	Number of Solves		
	First 50 episodes	Last 50 episodes	Total
TaskGen (only evaluated over 20 episodes)	100	100	100
F&S	28.3 ± 5.6	24.3 ± 3.8	52.6 ± 6.3
PPO	3.4 ± 3.0	7.2 ± 6.4	10.6 ± 8.9
TRPO	2.2 ± 1.4	2.3 ± 3.3	4.5 ± 2.6
A2C	1.4 ± 1.4	2.6 ± 2.1	4.0 ± 2.7

1125 Overall, Fig. D1 shows that TaskGen performs the best compared to all other agents such as F&S, PPO,
 1126 TRPO, A2C.

1127 **D.1.6 Details**

Table D2: Detailed Episode Run-through for TaskGen Agent. Environment change happens at Episode 11. Minimum steps is theoretical lowest number of steps to move from start to exit position. Actual steps is the actual number of steps the TaskGen Agent took to solve the episode.

Episode Num	Minimum Steps	Actual Steps
1	62	210
2	43	46
3	36	58
4	34	34
5	46	46
6	56	182
7	47	49
8	42	42
9	44	44
10	38	38
11	63	68
12	36	62
13	41	52
14	47	48
15	28	28
16	46	75
17	44	47
18	19	21
19	35	57
20	30	34

1128 Fig. D2 shows that in general, TaskGen Agent is able to solve the episodes quite efficiently with planning.
 1129 The main cause of the higher actual steps for some environments is when some obstacle positions are
 1130 unknown, the planner is not perfect and sometimes chooses a position to backtrack to that does not work
 1131 out. When the obstacles in the path are known, the planner can usually generate a perfect plan, or can
 1132 correct itself quickly mid-way.

1133 In general, as the number of episodes increase in the same environment, the better the knowledge of the
 1134 obstacles, and hence the better the plan and the generated actions. Even if the obstacles are changed, like
 1135 in Episode 11, the planner and TaskGen Agent combined is able to navigate and still clear the environment.

D.1.7 Insights	1136
As we are trying to test TaskGen’s ability to solve the maze, we intentionally use LLM as the planner and LLM as the executor for the maze environment. However, such a logical pathfinding task is best done by rule-based deterministic methods like Breadth-First Search or A* search algorithm. It may be better to treat pathfinding with known obstacles as a problem solved by traditional pathfinding algorithms, and simply get TaskGen to call such a function to do pathfinding.	1137
It is noted that the planner was not able to perform well without few-shot prompting of something which might occur in the environment (i.e. a wall with a gap). This is a huge downside for using LLM as an optimiser, as it does not optimise well. LLMs also do not understand 2D text grids perfectly, and hence, the spatial awareness for the planner is lacking, resulting in less robust plans.	1138
We have also tried to use native TaskGen without the planner, but the LLM was not able to see the big picture that well, resulting in LLM going to the same squares again and again trying to navigate past the wall. Planning is a difficult problem for an LLM and it is best to offload that to a rule-based planner.	1139
	1140
	1141
	1142
	1143
	1144
	1145
	1146
	1147
	1148

1149 E Escape Room Solving in TextWorld

1150 This appendix describes the implementation of an interactive fiction player as an agent. Interactive fiction
1151 is a genre of computer game that pre-dates GUIs, with many of the games (and tools) originating in the
1152 1980s. The Microsoft TextWorld project has delivered a system for building arbitrary games and provides
1153 a framework for building agents to navigate these games.

1154 The key system in interactive fiction is the discovery of the game world. Players are presented with
1155 limited information at one time and are required to recall or rediscover elements of the game "world".
1156 Within advanced games, the game world may change without player interaction but this behaviour is not
1157 present in TextWorld.

1158 Another aspect of interactive fiction is discovering how to interact with it. Players issue commands,
1159 typically in a terse pseudocode, and the game attempts to interpret them. TextWorld has optional support
1160 for providing the player/agent the list of acceptable commands at each turn. This agent will utilise these
1161 hints if provided by the game. By design, the agent does not depend on any specific input from the game
1162 and the few shot examples are not present in the development environment.

1163 The agent developed for this paper used the new Conversation Class interface, building on the existing
1164 escape room example.

1165 Interactive fiction games may have counter-intuitive problems to solve to succeed in the game. e.g.:
1166 to cook a carrot, grill it directly on a stove. For this developer, carrots aren't grilled, and things that are
1167 grilled are rarely done so directly on a stove.

1168 The Persistent Memories of the player/agent are aligned to the core systems of interactive fiction:
1169 commands, rooms (locations in the game), objectives. The "Summary of Conversation" in the player/agent
1170 was useful for the agent as it allowed the agent to reflect on futile behaviour and move on to alternative
1171 solutions. Structuring the Persistent Memory as "arrays" made the memories much more effective at
1172 guiding the agent.

F Web-Browsing Agents

1173

This appendix details the implementation of Web-Browsing Agents in TaskGen.

1174

The goal is to introduce the idea of Agents using a program/web application using TaskGen's current capability and performing actions based on user's query.

1175

1176

E.1 Agent Diagram/Flow

1177

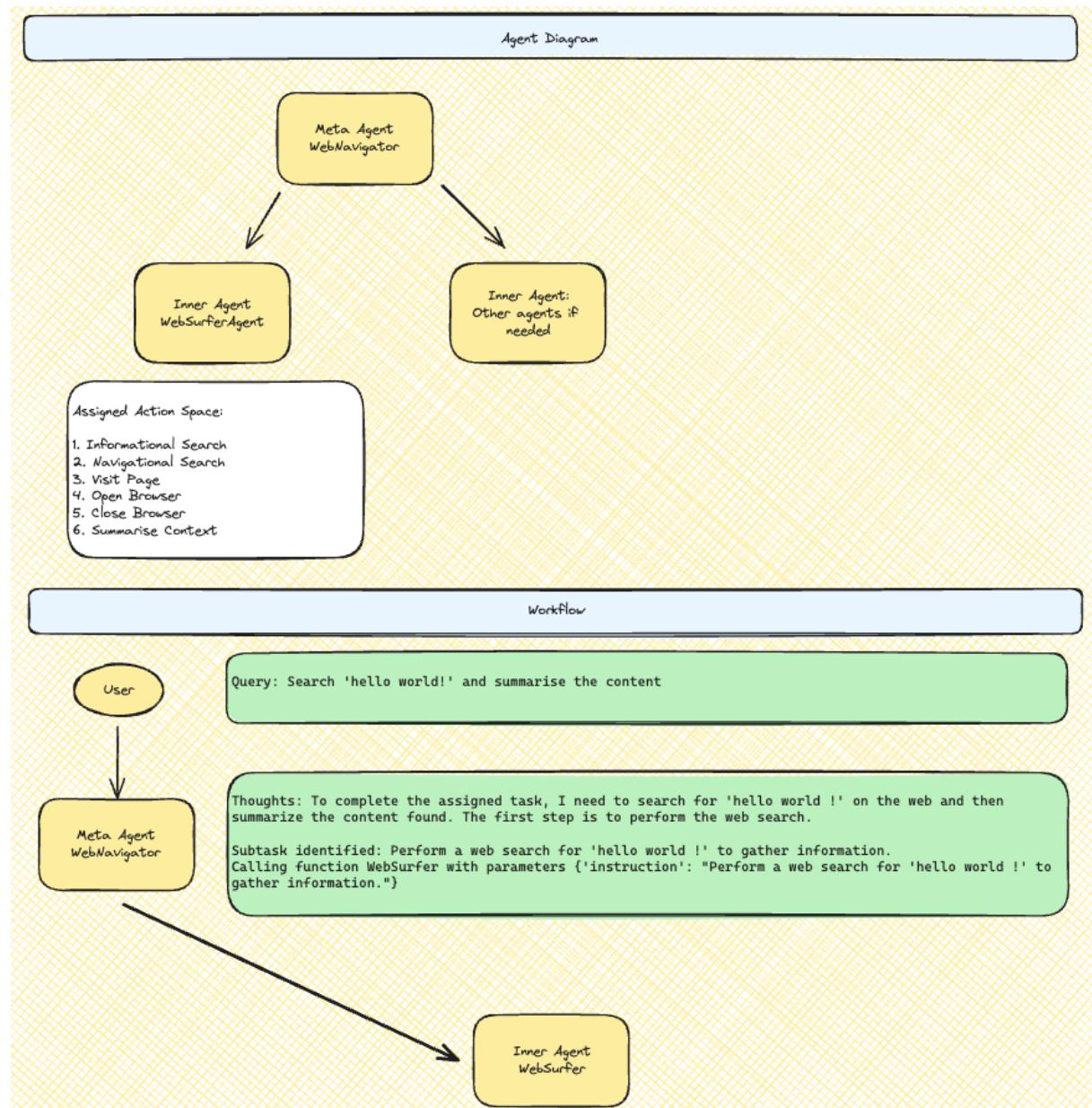


Figure F1: Agent Diagram and beginning of agent flow. Consists of 1 Meta Agent and 1 Inner Agent. User performs a query via terminal CLI that starts off agent flow.

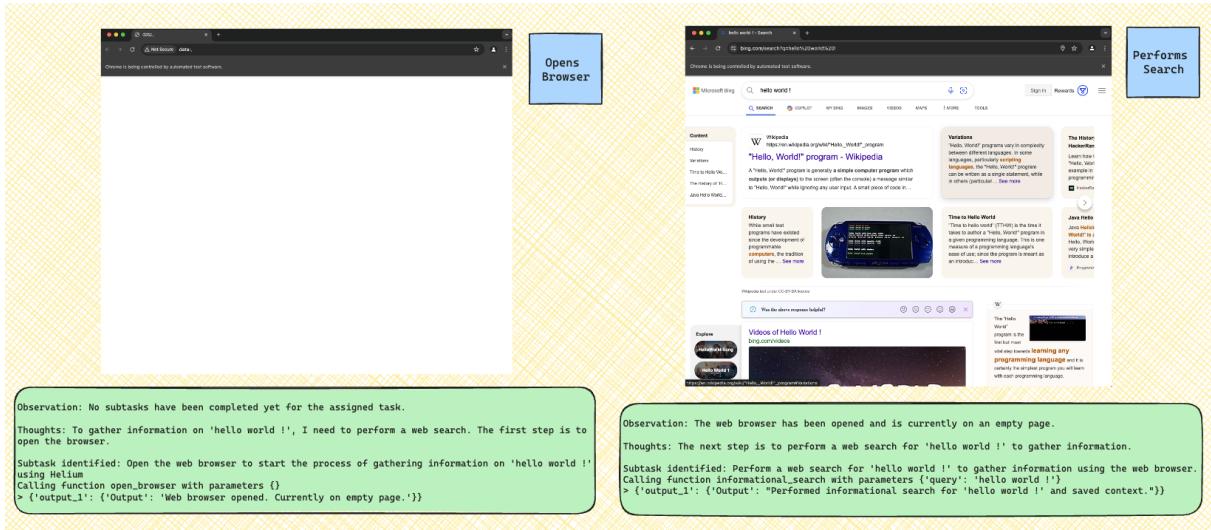


Figure F2: Inner Agent opens browser and performs search

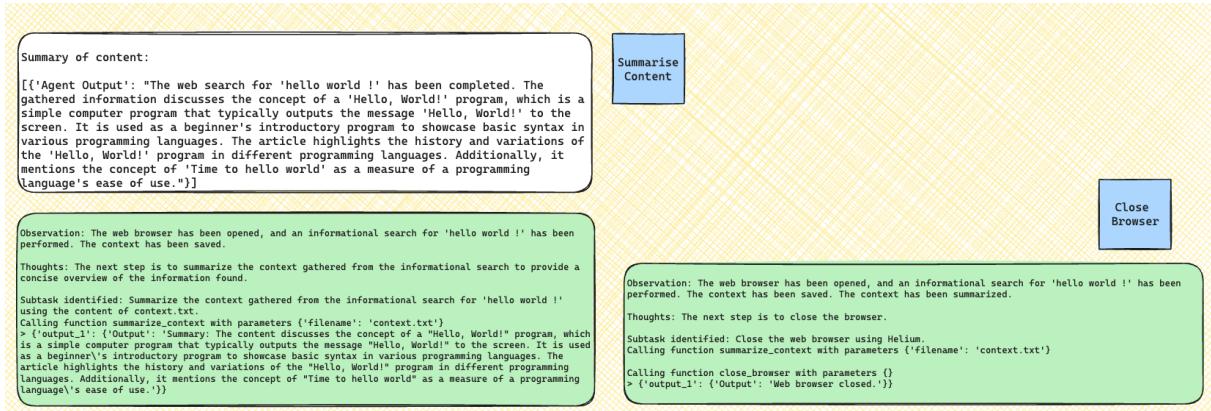


Figure F3: Inner Agent opens browser and performs search

F.2 Agent's Action Space	1178
• informational_search:	1179
– Description: Performs a search query on Bing and saves the context of the search results page.	1180
– Steps:	1181
1. Navigates to Bing with the given search query.	1182
2. Captures the current state of the browser (URL, title, and page source).	1183
3. Extracts and saves relevant content from the page to a file.	1184
• navigational_search:	1185
– Description: Performs a search query on Bing, clicks the first result, and saves the context of the resulting page.	1186
– Steps:	1187
1. Navigates to Bing with the given search query.	1188
2. Clicks on the first search result link.	1189
3. Captures the current state of the browser (URL, title, and page source).	1190
4. Extracts and saves relevant content from the page to a file.	1191
• visit_page:	1192
– Description: Visits a specified URL and saves the context of the page.	1193
– Steps:	1194
1. Navigates to the given URL.	1195
2. Captures the current state of the browser (URL, title, and page source).	1196
3. Extracts and saves relevant content from the page to a file.	1197
• open_browser:	1198
– Description: Opens a web browser using Helium.	1199
– Steps:	1200
1. Starts a Chrome browser session.	1201
2. Returns a message indicating the browser has been opened.	1202
• close_browser:	1203
– Description: Closes the web browser using Helium.	1204
– Steps:	1205
1. Kills the current browser session.	1206
2. Returns a message indicating the browser has been closed.	1207
• summarize_context:	1208
– Description: Summarizes the content saved in a file (default: <code>context.txt</code>) using OpenAI's GPT model.	1209
– Steps:	1210
1. Reads the content from the specified file.	1211
2. Sends the content to OpenAI's API to generate a summary.	1212
3. Returns the generated summary.	1213

1216 F.3 AutoGen Code for Web Browsing

1217 The AutoGen code is more complex and involves detailed setup and configuration. Here are its key
1218 components:

1219 F.3.1 Class Definition

```
1220 class WebSurferAgent(ConversableAgent):
1221     """(In preview) An agent that acts as a basic web surfer that can search the web
1222     and visit web pages."""
1223     ...
1224 
```

1226 F.3.2 Function Registration

```
1227 @self._user_proxy.register_for_execution()
1228 @self._assistant.register_for_llm(
1229     name="informational_web_search",
1230     description="Perform an INFORMATIONAL web search query then return the search
1231     results.",
1232 )
1233 def _informational_search(query: Annotated[str, "The informational web search query
1234     to perform."]) -> str:
1235     self.browser.visit_page(f"bing: {query}")
1236     header, content = _browser_state()
1237     return header.strip() + "\n=====\n" + content
1238 
```

1240 F.3.3 Summarisation Client

```
1241 def _create_summarizer_client(self, summarizer_llm_config: Dict[str, Any],
1242     llm_config: Dict[str, Any]) -> None:
1243     ...
1244 
```

1246 F.3.4 Reply Generation

```
1247 def generate_surfer_reply(
1248     self,
1249     messages: Optional[List[Dict[str, str]]] = None,
1250     sender: Optional[Agent] = None,
1251     config: Optional[OpenAIWrapper] = None,
1252 ) -> Tuple[bool, Optional[Union[str, Dict[str, str]]]]:
1253     ...
1254 
```

1256 F.3.5 Browser Interaction

```
1257 self.browser = SimpleTextBrowser(**(browser_config if browser_config else {}))
1258 
```

F.4 TaskGen Code for Web Browsing

1260

The TaskGen code is more straightforward and focuses on simplicity. Here are its key components:

1261

F.4.1 Function Definitions

1262

```
def informational_search(query: str) -> str:
    go_to(f"https://www.bing.com/search?q={query}")
    header, content = _browser_state()
    save_context_to_file(header, content)
    return {
        "Output": f"Performed informational search for '{query}' and saved context."
    }
```

1263

1264

1265

1266

1267

1268

1269

1270

F.4.2 Agent Creation

1272

```
WebSurfer = Agent(
    "WebSurfer",
    "Performs web searches and navigates web pages. Always open the browser at the start of the task and close the browser at the end.",
    model="gpt-4o",
    default_to_llm=False,
).assign_functions(fn_list_3)
```

1273

1274

1275

1276

1277

1278

1279

1280

1281

F.4.3 Boss Agent

1282

```
bossagent = Agent(
    "WebNavigator",
    "Assists user to navigate the web. Always open the browser at the start of the task and close the browser at the end.",
    model="gpt-4o",
    default_to_llm=False,
)
```

1283

1284

1285

1286

1287

1288

1289

1290

1292 **F.5 Key Differences between AutoGen and TaskGen Code for Web Browsing**

1293 **F.5.1 Complexity**

- 1294 • AutoGen involves more complex setup and configuration, making it suitable for advanced use cases.
- 1295 • TaskGen is simpler and more user-friendly, making it easier to build and assign functions to agents.

1296 **F.5.2 Function Registration**

- 1297 • AutoGen uses decorators for function registration.
- 1298 • TaskGen directly assigns functions to agents using `assign_functions`.

1299 **F.5.3 Browser Interaction**

- 1300 • AutoGen uses `SimpleTextBrowser`.
- 1301 • TaskGen uses `Helium` and `BeautifulSoup`.

1302 **F.5.4 Summarisation**

- 1303 • AutoGen sets up a summarization client with detailed configuration.
- 1304 • TaskGen uses a straightforward function to summarize content.

F.6 Taskgen Web Browsing Agent results

1305

F.6.1 Query Results

1306

Query	Success Rate
Search 'impact of social media on mental health' and summarize the academic studies	0/5
Visit ' https://www.who.int ' and summarize the latest health advisories	5/5
Open the browser, search 'weather forecast New York', and save the first result	5/5
Navigate to ' https://www.bbc.com/news ', extract the top news headlines, and summarize them	5/5
Search 'global warming statistics 2024' and provide a summary of the data trends	0/5
Visit ' https://www.finance.yahoo.com ', gather the latest stock market updates, and summarize	5/5
Search 'Shakespeare's influence on modern literature' and summarize the academic articles	5/5
Search 'quantum computing vs classical computing' and summarize the differences from multiple sources	5/5
Visit ' https://www.nasa.gov ', gather the latest Mars mission updates, and cross-reference with Wikipedia	5/5
Search 'key events in AI development 2024' and summarize the timeline	1/5
Search 'market analysis of electric vehicles 2024' and summarize the findings	3/5
Search 'best noise-canceling headphones 2024' and summarize the top reviews	5/5
Navigate to ' https://www.consumerreports.org ', gather information on washing machines, and summarize the best options	5/5
Visit ' https://docs.python.org ', find information on Python decorators, and summarize	5/5
Search 'workplace safety measures during COVID-19' and summarize the guidelines	2/5
Visit ' https://www.cdc.gov ', find information on flu prevention, and summarize	4/5
Search 'latest trends in renewable energy 2024' and summarize the key developments	0/5
Visit ' https://www.techcrunch.com ', gather the latest technology news, and summarize	5/5
Search 'evolution of jazz music' and summarize its impact on modern genres	4/5
Visit ' https://www.metmuseum.org ', explore the latest exhibits, and summarize	5/5

Table F1: Success rates for various web browsing queries

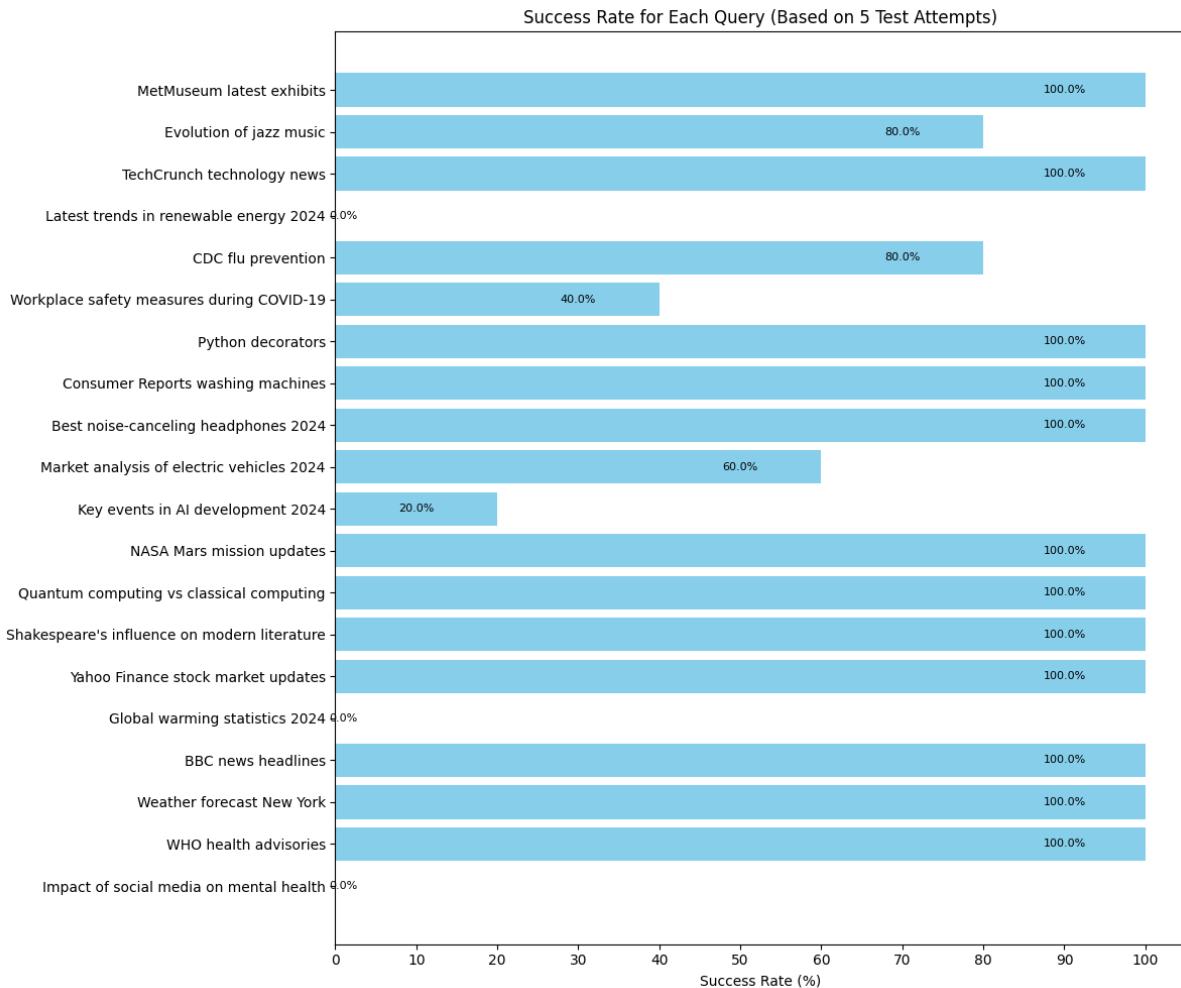


Figure F4: Graphical representation of the success rates for each query tested. Each query was tested 5 times, and the success rate is calculated as the percentage of successful attempts out of these 5 tests. The chart compares the effectiveness of different queries, providing a clear visualization of the success rate for each query.

G MATH Dataset

1307

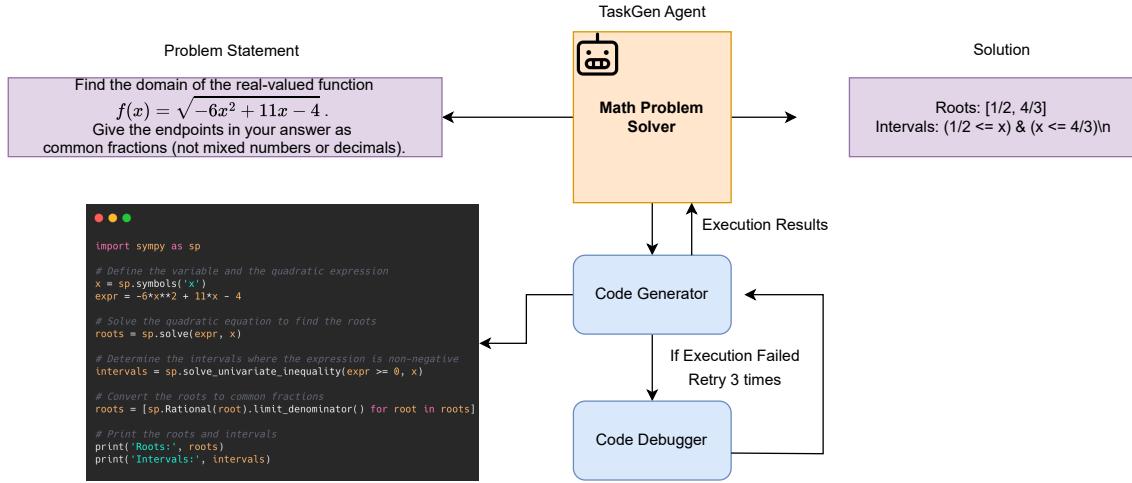


Figure G1: Math Problem Solver Agent

Leveraging LLMs for solving mathematical problems has become a widely researched area (Zhou et al., 2023; Yang et al., 2023b). In this section, we explore the use of the TaskGen agent to tackle complex mathematical problems across various domains. For our evaluations, we utilized the MATH dataset (Hendrycks et al., 2021), which contains over 12,500 competition-level mathematical problems. We focused specifically on the most challenging subset, **Level 5**, across the following categories: Algebra, Pre-Algebra, Intermediate Algebra, Number Theory, and Counting and Probability. We randomly selected 20 problems from the test set of each category, resulting in a total of 100 problems in our test set, to assess the TaskGen agent’s ability to solve these tasks.

To solve these challenging problems, we employed a TaskGen Agent called “Math Problem Solver” powered with GPT4o as depicted in Figure G1. This agent is equipped with a specialized functions that facilitates the generation, execution, and debugging of code necessary to tackle the given tasks. The function has access to essential Python libraries, including numpy, sympy, math, and random.

1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319

1320
1321

Evaluation Result. In Figure G2, we provide the evaluation results of our TaskGen agent equipped with the function described above and the agent without any equipped functions.

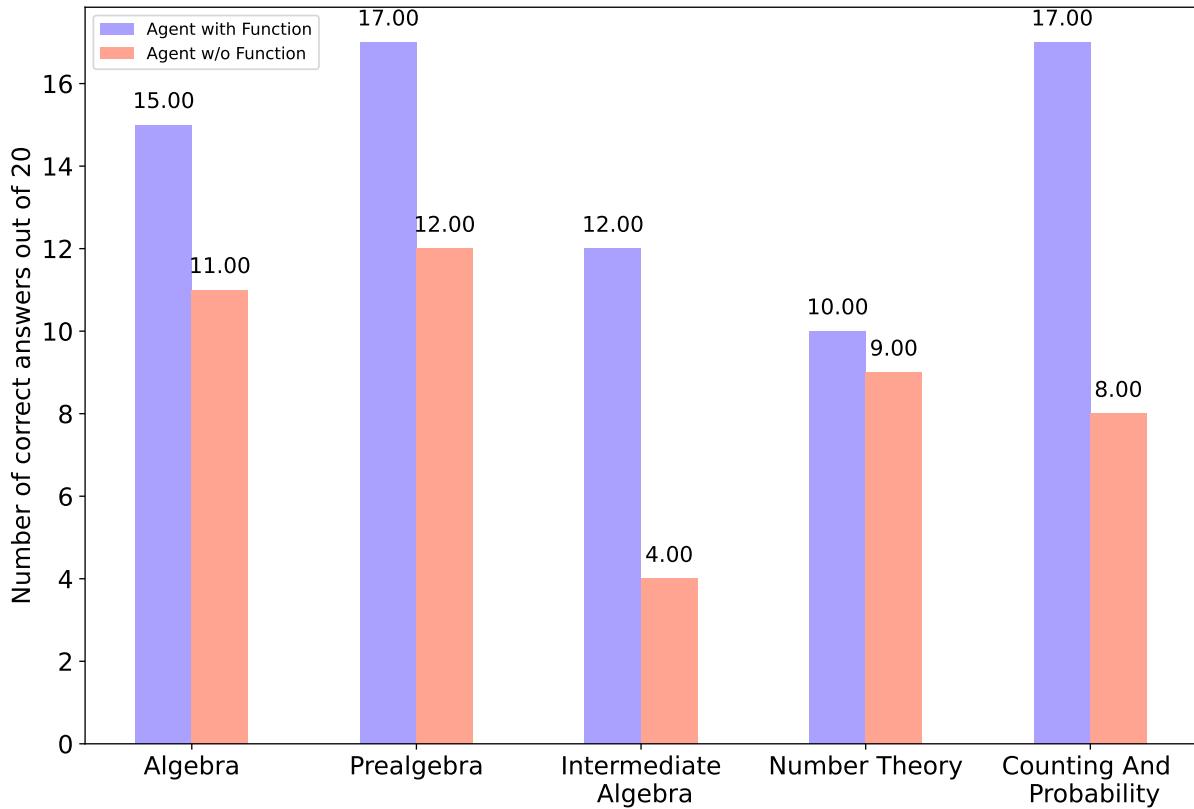


Figure G2: Quantitative results of TaskGen agents on the subset of the MATH dataset.

1322 From our experiments, we found that on the challenging Level-5 problems, the TaskGen agent with
1323 Equipped Functions achieved an average accuracy of **71%**, while the TaskGen agent without the Equipped
1324 Functions achieved only **44%** accuracy. For evaluation, we manually verified the generated solution
1325 against the provided ground truth solution. These results demonstrate that, in order to solve these
1326 challenging Level-5 problems, equipping the agent with code generation and debugging capabilities leads
1327 to more accurate solutions of mathematical problems.

H RAG-based Question Answering on NaturalQuestions Dataset

In this section, we describe the development and functionality of a Retrieval-Augmented Generation (RAG) system using TaskGen. This system integrates one TaskGen agent, known as the “User Agent,” along with two critical TaskGen functions: **ContextFetchFunction** and **AnswerFunction**. These components constitute the fundamental operations of our system.

H.1 System Overview

1. **ContextFetchFunction:** This function accepts a user’s query and a batch number, retrieving the relevant context. It is designed to incrementally fetch more context if the initial retrieval proves inadequate. 1334
2. **AnswerFunction:** After receiving context, this function generates an answer based on the context available. If the context is insufficient to resolve the query, **AnswerFunction** returns “no answer.” 1337
3. **User Agent:** The orchestrator of the entire Q&A cycle, the User Agent is responsible for managing both the **ContextFetchFunction** and the **AnswerFunction**. It initiates the process by retrieving context for the user’s query and continues to fetch additional context in subsequent batches until a satisfactory answer is found or the interaction limit is reached. 1339

H.2 Detailed Process

- **Query Submission:** The user submits a query to the User Agent. 1344
- **Context Retrieval:** The User Agent invokes **ContextFetchFunction** with the initial query and a starting batch number (0). 1345
- **Answer Generation:** With the context obtained, the User Agent next activates **AnswerFunction**. If the context sufficiently answers the query, a response is generated. Otherwise, it issues “no answer.” 1347
- **Incremental Fetching:** If “no answer” is received, the User Agent increments the batch number and re-engages **ContextFetchFunction** to obtain more context. This iterative process is capped at five interactions (max interactive retrieval count). 1349

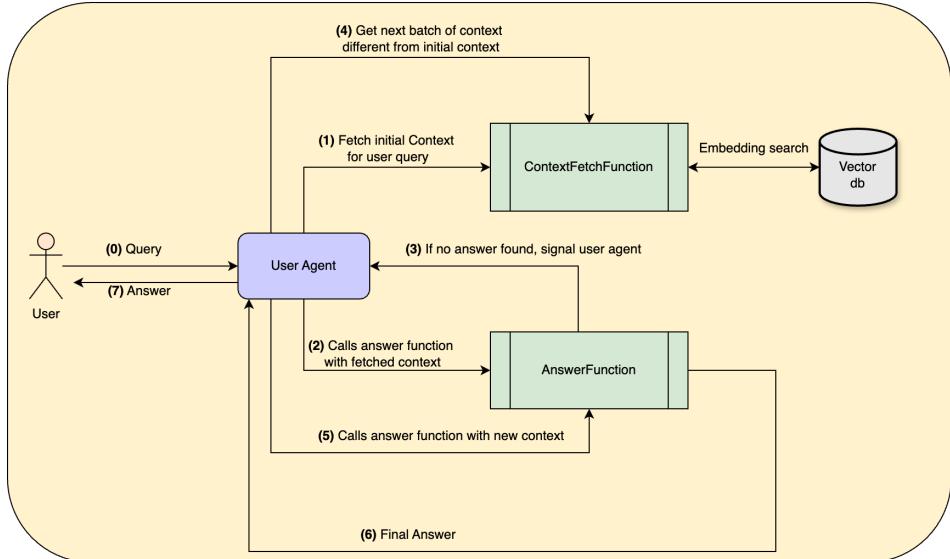


Figure H1: Illustration of the Interactive Retrieval-Augmented Generation (RAG) Question and Answer Flow. The diagram sequentially represents the process from step (0) to (7), detailing the interaction between the User, User Agent, ContextFetchFunction, and AnswerFunction. Each numbered marker (num) in the diagram corresponds to a specific step in the query-answer cycle.

1352 **H.3 Example**

1353 Consider the query: “What is the capital of France?”

1354 • **First Batch (Batch 0):**

- 1355 1. Paris is a major city in France.
1356 2. France is known for its culture and cuisine.

1358 *AnswerFunction Output:* “no answer” (the context does not explicitly state Paris as the capital).

1359 • **Second Batch (Batch 1):**

- 1360 1. The capital of France is Paris.
1361 2. Paris is famous for the Eiffel Tower.

1363 *AnswerFunction Output:* “The capital of France is Paris.”

1364 This example illustrates the system’s ability to handle complex queries by sequentially enhancing the
1365 context until a definitive answer can be provided.

1366 **H.4 Technical Framework and Evaluation Methodology**

1367 This section outlines the technical and methodological specifics employed in our study to develop and
1368 evaluate the Retrieval-Augmented Generation (RAG) system.

1369 **H.4.1 Technology Stack**

- **TaskGen Agents and Functions:** Our system utilizes TaskGen’s capabilities extensively. The core components, namely the **User Agent**, **ContextFetchFunction**, and **AnswerFunction**, are powered by GPT-3.5. This model was chosen for its lower cost and robust performance in natural language understanding and generation.
- **Embedding Storage and Retrieval:** We employed **Postgres PGvector** to manage the storage and retrieval of embeddings. For Retrieval we use $k = 10$ configuring our system to fetch top 10 most relevant vector embeddings for each query
- **Embedding Model:** The **text-embedding-ada-002** model was used to convert text data into vector embeddings. These embeddings represent the textual data in a format amenable to similarity comparisons and retrieval operations.

1380 **H.4.2 Dataset and Benchmarking**

- **Dataset:** The [Natural Questions dataset](#) was chosen for its comprehensive collection of real-world questions. Our study focuses on the first 2,000 entries of the development split validation set, providing a balanced mix of complexity and coverage.
- **Evaluation Metrics:** To assess the effectiveness of our RAG system, we used Google’s [nq_eval script](#). This script is widely recognized for its rigor in measuring the precision and accuracy of answers provided by question-answering systems.

H.4.3 Evaluation Results

We conducted a comprehensive evaluation to compare the performance of non-interactive versus interactive(via TaskGen) retrieval method. The non-interactive retrieval approach involves a single invocation of an LLM using context from the vector database to answer the query. This method assumes that the initial context contains all the necessary information to generate an answer. In contrast, the interactive retrieval method dynamically fetches and refines context based on the ongoing interaction with the user's query, allowing for a more adaptive and potentially accurate response as additional information is incorporated in successive retrieval steps.

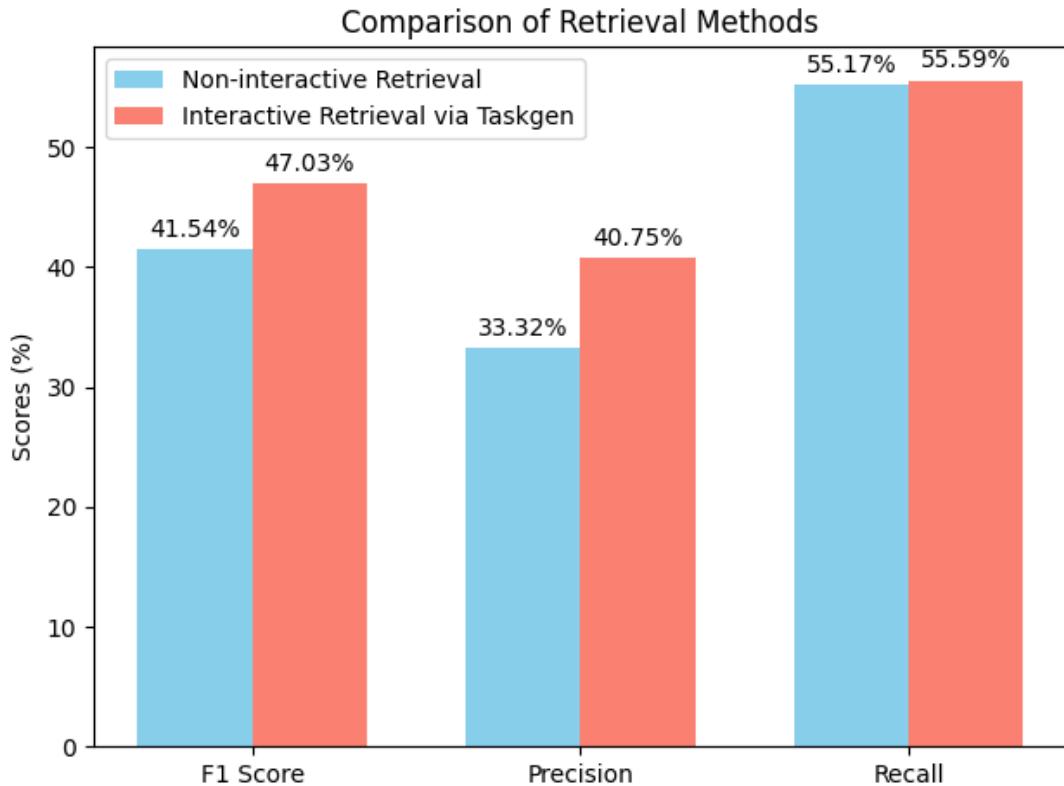


Figure H2: Graphical representation of the benchmark results comparing F1 Score, Precision, and Recall for Non-Interactive versus Interactive Retrieval via TaskGen (for both k=10 used for retrieval).