
AVIS: Autonomous Visual Information Seeking with Large Language Model Agent

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper, we propose an autonomous information seeking visual question
2 answering framework, AVIS. Our method leverages a Large Language Model
3 (LLM) to dynamically strategize the utilization of external tools and to investigate
4 their outputs, thereby acquiring the indispensable knowledge needed to provide
5 answers to the posed questions. Responding to visual questions that necessitate
6 external knowledge, such as "What event is commemorated by the building depicted
7 in this image?", is a complex task. This task presents a combinatorial search space
8 that demands a sequence of actions, including invoking APIs, analyzing their
9 responses, and making informed decisions. We conduct a user study to collect a
10 variety of instances of human decision-making when faced with this task. This data
11 is then used to design a system comprised of three components: an LLM-powered
12 planner that dynamically determines which tool to use next, an LLM-powered
13 reasoner that analyzes and extracts key information from the tool outputs, and
14 a working memory component that retains the acquired information throughout
15 the process. The collected user behavior serves as a guide for our system in
16 two key ways. First, we create a transition graph by analyzing the sequence of
17 decisions made by users. This graph delineates distinct states and confines the
18 set of actions available at each state. Second, we use examples of user decision-
19 making to provide our LLM-powered planner and reasoner with relevant contextual
20 instances, enhancing their capacity to make informed decisions. We show that AVIS
21 achieves state-of-the-art results on knowledge-intensive visual question answering
22 benchmarks such as Infoseek [7] and OK-VQA [26].

23 1 Introduction

24 Large language models (LLMs), such as GPT3 [5], LaMDA [16], PALM [9], BLOOM [34] and
25 LLaMA [37], have showcased the capacity to memorize and utilize a significant amount of world
26 knowledge. They demonstrate emerging abilities [38] like in-context learning [5], code genera-
27 tion [19], and common sense reasoning [24]. Recently, there is a growing focus towards adapting
28 LLMs to handle multi-modal inputs and outputs involving both vision and language. Noteworthy
29 examples of such visual language models (VLMs) include GPT4 [29], Flamingo [4] and PALI [6].
30 They set the state of the art for several tasks, including image captioning, visual question answering,
31 and open vocabulary recognition.

32 While LLMs excel beyond human capabilities in tasks involving textual information retrieval, the
33 current state of the art VLMs perform inadequately on datasets designed for visual information
34 seeking such as Infoseek [7], Oven [14] and OK-VQA [26]. Many of the visual questions in these
35 datasets are designed in such a way that they pose a challenge even for humans, often requiring the

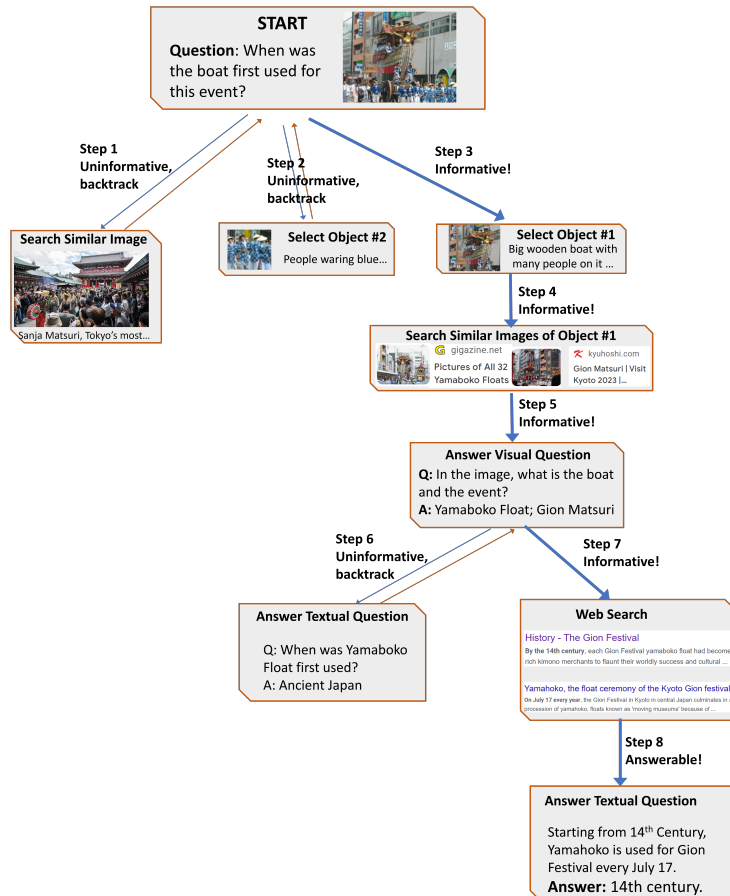


Figure 1: An example of AVIS’s generated workflow for answering a challenging visual question using LLM with tree search to use tools. The input image is taken from the Infoseek dataset.

36 assistance of various APIs and web search to obtain the answer. Examples of such questions include
 37 "where is this church located?", "what species of butterfly is this?", or "what is the brand of this
 38 dress?".

39 Current state-of-the-art vision-language models (VLMs) find it challenging to answer such questions
 40 for several reasons. Firstly, they are not trained with objectives that encourage them to discern
 41 fine-grained categories and details within images. Secondly, they utilize a relatively smaller language
 42 model compared to state-of-the-art Large Language Models (LLMs), which constrains their reasoning
 43 capabilities. Lastly, they do not compare the query image against a substantial corpus of images
 44 associated with varying metadata, unlike systems that employ image search techniques.

45 To overcome these challenges, we introduce a novel method in this paper that achieves state-of-the-
 46 art results on visual information seeking tasks by integrating LLMs with three types of tools: (i)
 47 computer vision tools such as object detection, OCR, image captioning models, and VQA models,
 48 which aid in extracting visual information from the image, (ii) a web search tool that assists in
 49 retrieving open world knowledge and facts, and (iii) an image search tool that enables us to glean
 50 relevant information from metadata associated with visually similar images. Our approach utilizes an
 51 LLM-powered planner to dynamically determine which tool to use at each step and what query to
 52 send to it. Furthermore, we employ an LLM-powered reasoner that scrutinizes the output returned
 53 by the tools and extracts the crucial information from them. To retain the information throughout
 54 the process, we use a working memory component. Figure 1 shows an example information seeking
 55 process performed by our method.

56 Several recent studies [13, 23, 36, 40, 42] have enhanced LLMs with APIs to handle multi-modal
 57 vision-language inputs. These systems generally employ a two-stage strategy, namely *plan* and
 58 *execute*. Initially, the LLM breaks down a question into a plan, typically represented as a structured
 59 program or a sequence of instructions. Following this, the necessary APIs are activated to collect the
 60 required information. While this method has shown potential in elementary visual-language tasks, it
 61 frequently fails in more complex real-world situations. In such cases, a comprehensive plan cannot

62 be inferred merely from the initial question. Instead, it necessitates dynamic modifications based on
63 real-time feedback.

64 The primary innovation in our proposed method lies in its dynamic decision-making capability.
65 Answering visual information seeking questions is a highly complex task, requiring the planner
66 to take multiple steps. At each of these steps, the planner must determine which API to call and
67 what query to send. It is unable to predict the output of complex APIs, such as image search, or to
68 anticipate the usefulness of their responses prior to calling them. Therefore, unlike previous methods
69 that pre-plan the steps and API calls at the beginning of the process, we opt for a dynamic approach.
70 We make decisions at each step based on the information acquired from previous API calls, enhancing
71 the adaptability and effectiveness of our method.

72 We conduct a user study to gather a wide range of instances of human decision-making when using
73 APIs to answer questions related to visual information seeking. From this data, we formulate a
74 structured framework that directs the Large Language Model (LLM) to use these examples for making
75 informed decisions regarding API selection and query formulation. The collected user behavior
76 informs our system in two significant ways. First, by analyzing the sequence of user decisions, we
77 construct a transition graph. This graph delineates distinct states and constrains the set of actions
78 available at each state. Second, we use the examples of user decision-making to guide our planner
79 and reasoner with pertinent contextual instances. These contextual examples contribute to improving
80 the performance and effectiveness of our system.

81 The primary contributions of this paper can be summarized as follows:

- 82 • We propose a novel visual question answering framework that leverages a large language
83 model (LLM) to dynamically strategize the utilization of external tools and to investigate
84 their outputs, thereby acquiring the necessary knowledge needed to provide answers to the
85 posed questions.
- 86 • We leverage the human decision-making data collected from a user study to develop a
87 structured framework. This framework guides the Large Language Model (LLM) to utilize
88 examples of human decision-making in making informed choices concerning API selection
89 and query construction.
- 90 • Our method achieves state-of-the-art results on knowledge-based visual question answering
91 benchmarks such as Infoseek [7] and OK-VQA [26]. Notably, We achieve an accuracy of
92 50.7% on the Infoseek (unseen entity split) dataset which is significantly higher than the
93 results achieved by PALI [6] with accuracy of 16.0%.

94 2 Related Work

95 **Augmenting LLMs with Tools.** Large Language Models (LLMs) have shown impressive language
96 understanding [33], and even reasoning capabilities [39]. Nevertheless, certain limitations of LLMs
97 are evident, due to their intrinsic characteristics. Such limitations include providing up-to-date
98 answers based on external knowledge or performing mathematical reasoning. Consequently, a recent
99 surge of techniques have integrated LLMs with various external tools [27]. For example, TALM [31]
100 and ToolFormer [35] use in-context learning to teach the language model how to better leverage
101 various tools on benchmarks such as question answering and mathematical reasoning.

102 In the computer vision domain, LLMs also show significant improvements when combined with
103 external visual tools. For example, Visual ChatGPT [40] and MM-ReAct [42] enable LLMs to call
104 various vision foundation models as tools to understand visual inputs, and even better control the
105 image generation. VisProg [13] and ViperGPT [36] explore the decomposition of visual language
106 tasks into programs, where each line corresponds to general code or a visual API. Chameleon [23]
107 uses an LLM as a natural language planner to infer the appropriate sequence of tools to utilize, and
108 then executes these tools to generate the final response.

109 Most of these previous works follow a plan-then-execute paradigm, i.e., i) they pre-plan the sequence
110 of actions (API calls) that they will take (either hard coded or using code generation); and ii) they
111 execute the generated plan. One drawback of such an approach is that it cannot update and improve
112 its plan based on the output of the tools it calls. This is not a trivial problem, as it requires to predict
113 the output quality of each tools beforehand. In contrast, our proposed method allows the system to
114 dynamically decide its next steps based on the output it receives from the tools at each step.

115 **Decision Making with LLM as an Agent.** There has also been a surge of interest in applying
 116 Large Language Models (LLMs) as autonomous agents. These agents are capable of interacting with
 117 external environments, making dynamic decisions based on real-time feedback, and consequently
 118 achieving specific goals. For example, WebGPT [28] enables an LLM to access real-time information
 119 from the web search engines. ReAct [43] further improves external search engine usage via the self-
 120 reasoning of LLM in an interleaved manner. Similar ideas have also been adopted for robotic action
 121 planning. SayCan [3], for instance, uses LLMs to directly predict robot actions, and PALM-E [10]
 122 further fine-tunes LLMs to make better decisions based on instructions and open web media.

123 When compared to works that follow a plan-then-execute paradigm, these AI agents exhibit increased
 124 flexibility, adjusting their actions based on the feedback that they receive. However, many of these
 125 methods do not restrict the potential tools that can be invoked at each stage, leading to an immense
 126 search space. This becomes particularly critical for web search APIs [1, 2] that return extensive result
 127 lists and span a combinatorial search space of multiple tools. Consequently, even the most advanced
 128 LLMs today can fall into infinite loops or propagate errors. To alleviate this issue, we propose
 129 restricting and guiding LLMs to mimic human behavior when solving complex visual questions with
 130 APIs. This idea is similar to the AI alignment research [21, 30] that teaches LLMs to follow human
 131 instructions. The difference is that our model only uses the human prior at the decision-making stage
 132 via prompt guidance, instead of re-training the model.

133 3 Method

134 3.1 General Framework

135 Our approach employs a dynamic decision-making strategy designed to respond to visual information-
 136 seeking queries. Our system is comprised of three primary components. First, we have a planner \mathcal{P} ,
 137 whose responsibility is to determine the subsequent action, including the appropriate API call and
 138 the query it needs to process. Second, we have a working memory \mathcal{M} that retains information about
 139 the results obtained from API executions. Lastly, we have a reasoner \mathcal{R} , whose role is to process the
 140 outputs from the API calls. It determines whether the obtained information is sufficient to produce
 141 the final response, or if additional data retrieval is required.

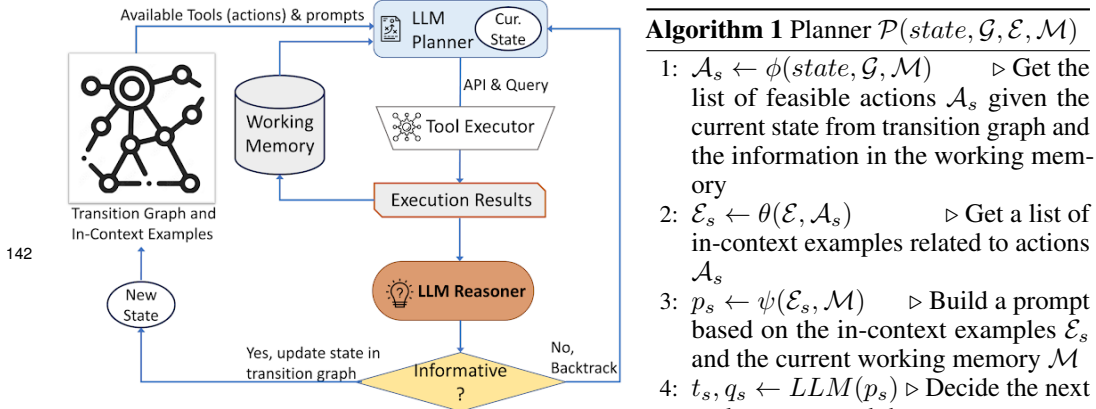


Figure 2: AVIS employs dynamic decision-making to **plan** (find optimal tool and query), execute results, and then **reason** (estimate whether continue or backtrack).

Algorithm 1 Planner $\mathcal{P}(state, \mathcal{G}, \mathcal{E}, \mathcal{M})$

- 1: $\mathcal{A}_s \leftarrow \phi(state, \mathcal{G}, \mathcal{M})$ \triangleright Get the list of feasible actions \mathcal{A}_s given the current state from transition graph and the information in the working memory
 - 2: $\mathcal{E}_s \leftarrow \theta(\mathcal{E}, \mathcal{A}_s)$ \triangleright Get a list of in-context examples related to actions \mathcal{A}_s
 - 3: $p_s \leftarrow \psi(\mathcal{E}_s, \mathcal{M})$ \triangleright Build a prompt based on the in-context examples \mathcal{E}_s and the current working memory \mathcal{M}
 - 4: $t_s, q_s \leftarrow LLM(p_s)$ \triangleright Decide the next tool t_s to use and the query q_s to pass by feeding the prompt p_s to LLM
-

143 Considering the potential intricacy of the task, we conduct a user study to gather a broad range of
 144 examples of human decision-making process, when using tools to respond to visual information-
 145 seeking queries (we introduce the details of data collection in Sec. 3.3). This helps us to establish a
 146 structured framework for decision-making. We utilize the data collected from this study to construct
 147 a transition graph \mathcal{G} shown in Figure 3, which outlines all the possible actions at each given state.
 148 Additionally, we employ real-life decision-making examples \mathcal{E} , i.e., users choose which tool at
 149 different states, to guide the planner in choosing the appropriate action at each stage of the process.

150 The Algorithm 1 presents the operations of the planner \mathcal{P} . The planner undertakes a series of steps
 151 each time a decision is required regarding which tool to employ and what query to send to it. Firstly,

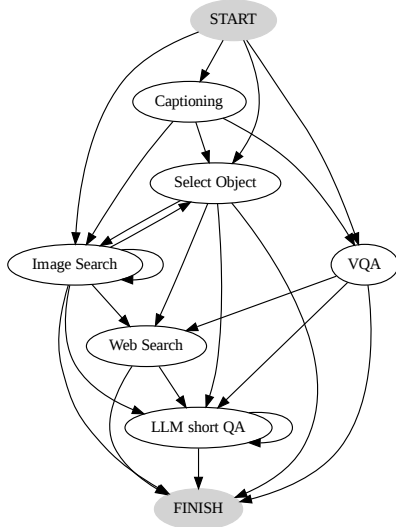
152 based on the present *state*, the planner provides a range of potential subsequent actions \mathcal{A}_s . The
 153 potential action space \mathcal{A}_s may be large, making the search space intractable. To address this issue, the
 154 planner refers to the human decisions from the transition graph \mathcal{G} to eliminate irrelevant actions. The
 155 planner also excludes the actions that have already been taken before and are stored in the working
 156 memory \mathcal{M} . Formally, this procedure is $\mathcal{A}_s \leftarrow \phi(\text{state}, \mathcal{G}, \mathcal{M})$.

157 Next, it collects a set of relevant in-context examples \mathcal{E}_s that are assembled from the decisions
 158 previously made by humans during the user study relevant to actions \mathcal{A}_s , that is $\mathcal{E}_s \leftarrow \theta(\mathcal{E}, \mathcal{A}_s)$. With
 159 the gathered in-context examples \mathcal{E}_s and the working memory \mathcal{M} that holds data collected from past
 160 tool interactions, the planner formulates a prompt, denoted by $p_s \leftarrow \psi(\mathcal{E}_s, \mathcal{M})$. The prompt p_s
 161 is then sent to the LLM which returns a structured answer, determining the next tool t_s to be activated
 162 and the query q_s to be dispatched to it. We denote this action by $t_s, q_s \leftarrow LLM(p_s)$. This design
 163 allows the planner to be invoked multiple times throughout the process, thereby facilitating dynamic
 164 decision-making that gradually leads to answering the input query.

165 The Algorithm 2 shows the overall decision-making workflow of AVIS. The entire process repeats
 166 until a satisfactory answer is produced. Initially, the working memory is populated only with the input
 167 visual question I , and the initial *state* is set to START. At each iteration, we first invoke the planner
 168 \mathcal{P} to determine the next tool and the query to employ, as outlined in Algorithm 1. Subsequently, the
 169 selected external tool executes and delivers its output o_s . The output from the tools can be quite
 170 diverse, ranging from a list of identified objects, to a collection of similar images with their captions,
 171 to snippets of search results or knowledge graph entities.

172 Therefore, we employ a reasoner \mathcal{R} to analyze the output o_s , extract the useful information and decide
 173 into which category the tool output falls: informative, uninformative, or final answer. Our method
 174 utilizes the LLM with appropriate prompting and in-context examples to perform the reasoning. If
 175 the reasoner concludes that it’s ready to provide an answer, it will output the final response, thus
 176 concluding the task. If it determines that the tool output is uninformative, it will revert back to the
 177 planner to select another action based on the current state. If it finds the tool output to be useful, it
 178 will modify the state and transfer control back to the planner to make a new decision at the new state.

179 To illustrate with a tangible example, we can refer to the output that the model would receive as
 180 depicted in Figure 4(c). There are several entities within the answer. The role of the reasoner is
 181 twofold: to determine which entity is pertinent for responding to the question and to assess whether
 182 the model has obtained the necessary information to transition to the next state.



183

Figure 3: Transition graph \mathcal{G} defines feasible actions the planner can take. This graph is induced by our user study introduced in Sec. 3.3.

Algorithm 2 AVIS Decision Making Workflow

- 1: $\mathcal{M} \leftarrow \{\text{input}\}, \text{state} \leftarrow \text{START}$
 - 2: $t_s, q_s \leftarrow \mathcal{P}(\text{state}, \mathcal{G}, \mathcal{E}, \mathcal{M})$ ▷ Call the planner \mathcal{P} to decide the next tool to use t_s and the query to pass to it q_s
 - 3: $o_s \leftarrow \text{Exec}(t_s, q_s)$ ▷ Call tool t_s with query q_s and get output o_s
 - 4: $\hat{o}_s \leftarrow \mathcal{R}(o_s, \mathcal{M})$ ▷ Process the output and extract the key info \hat{o}_s using the reasoner \mathcal{R}
 - 5: $\mathcal{M}.\text{add}(\hat{o}_s)$ ▷ Update the working memory
 - 6: **switch** \hat{o}_s **do**
 - 7: **case** \hat{o}_s is not informative
 - 8: goto(2) ▷ Go to line 2 to make decision at the same state, excluding t_s .
 - 9: **case** \hat{o}_s has useful information
 - 10: $\text{state} \leftarrow t_s$ ▷ Update state
 - 11: goto(2) ▷ Go to line 2 to make decision for the next state.
 - 12: **case** \hat{o}_s is ready as final answer
 - 13: $\text{ans} \leftarrow \hat{o}_s$ ▷ Output answer
-

184 Our approach, which employs dynamic decision-making coupled with backtracking, differs from
 185 previous methods [23, 36] that follow a plan-then-execute paradigm. Our system is structured to

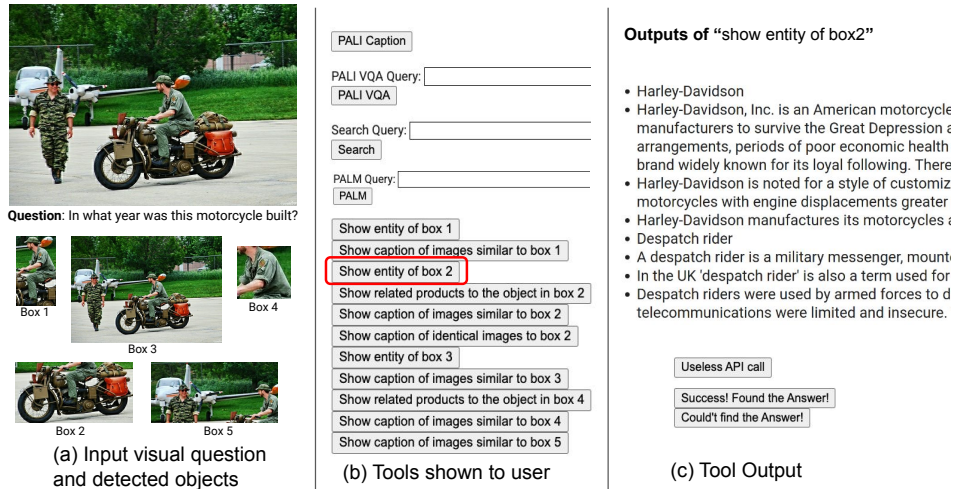


Figure 4: We conduct a user study to gather examples of user decision-making when responding to visual information-seeking questions. Given a visual question as depicted in (a), the user makes a series of tool calls using the available APIs shown in (b). Each tool call yields an output which the user reviews whether it is useful and determines the subsequent action, illustrated in (c).

186 make decisions grounded to the results of current executions and to conduct iterative searches for
 187 tool combinations. This process eventually yields the most effective strategy to accomplish the task.

188 3.2 Tools and their APIs

189 To respond effectively to visual queries that necessitate in-depth information retrieval, it’s important
 190 to equip AVIS with a comprehensive suite of tools. In this section, we describe these tools.

191 **Image Captioning Model:** We employ the PALI 17B [8] captioning model, which obtains state-of-
 192 the-art results for image captioning. This tool has the capability to generate captions for either the
 193 entire image or for a cropped image corresponding to the bounding box of a detected object.

194 **Visual Question Answering Model:** We utilize the PALI 17B [8] VQA model, which has been
 195 fine-tuned on the VQA-v2 [11] dataset. This tool takes an image and a question as inputs and provides
 196 a text-based answer as the output.

197 **Object Detection:** We use an object detector trained on a super-set of Open Images dataset [17]
 198 categories that is provided by Google Lens API [1]. We use high confidence threshold to only keep
 199 the top-ranked detected boxes for the input image.

200 **Image Search:** We utilize Google Image Search to obtain a broad range of information related to the
 201 image crop of a detected box as provided in Google Lens API [1]. This information encompasses
 202 various details, such as knowledge graph entities, titles of associated products, and captions of
 203 analogous or identical images. When it comes to decision-making, our planner considers the
 204 utilization of each piece of information as a separate action.

205 **OCR:** To identify text available in input image, we take advantage of the Optical Character Recognition
 206 (OCR) feature available in the Google Lens API [1].

207 **Web Search:** We employ the Google Web Search API [2]. It accepts a text-based query as input and
 208 produces the following outputs: (i) related document links and snippets, (ii) in certain instances, a
 209 knowledge panel providing a direct answer to the query.

210 **LLM short QA:** We incorporate a Language Model (LLM) powered question-answering component
 211 as another tool. This tool accepts a query in text form and produces an answer also in text form. It is
 212 important to note that the use of the LLM here as a question-answering tool is distinct from its role in
 213 the planner or reasoner as outlined in Alg. 1 and Alg. 2.

214 3.3 Gathering User Behavior to Inform LLM Decision Making

215 Many of the visual questions in datasets such as Infoseek [7], Oven [14] and OK-VQA [26] ask for
 216 fine-grained answers, which poses a challenge even for humans, often requiring the assistance of

217 various APIs and web searches for answers. Figure 4(a) illustrates an example visual question taken
218 from the OK-VQA [26] dataset. In order to gather insights into human decision-making process, we
219 carried out a user study. More specifically, our goal is to understand how humans utilize external
220 tools to answer visual queries that involve seeking information.

221 The user is equipped with an identical set of tools as our method. They are presented with the input
222 image and question, along with image crops for each detected object. Based on the information
223 obtained through image search for each cropped image, the user is offered one or multiple buttons
224 associated with each box. These buttons provide the user with the ability to access diverse information
225 pertaining to the image crop of the box. This includes details such as corresponding knowledge graph
226 entities, captions of similar images, titles of associated related products, and captions of identical
227 images. An example set of tools and APIs are shown in Figure 4(b).

228 When the user initiates an action, such as clicking on a button or submitting a query to web search,
229 PALM, or PALI VQA, the corresponding tool is invoked, and the resulting output is displayed to the
230 user. We record the sequence of actions taken by the user and the outputs that they receive at each
231 step. For instance, in Figure 4, we show an example of how a user needs to perform four actions to
232 answer the question: *i) display entities in box 2, ii) show the caption of similar images to box 2, iii)*
233 *conduct a search for "In what year was Harley-Davidson XA built?", and iv) utilize PALM using the*
234 *combination of the search output and the question "In what year was Harley-Davidson XA built?".*

235 The collected user behavior serves as a guide for our system in two key ways. Firstly, we construct a
236 transition graph by analyzing the sequence of decisions made by users. This graph defines distinct
237 states and restricts the available set of actions at each state. For example, at the START state, the
238 system can take only one of these three actions: PALI caption, PALI VQA, or object detection.
239 Figure 3 illustrates the transition graph that has been constructed based on the decision-making
240 process of the users. Secondly, we utilize the examples of user decision-making to guide our planner
241 and reasoner with relevant contextual instances. These in-context examples aid in enhancing the
242 performance and effectiveness of our system.

243 We conducted a user study involving 10 participants who collectively answered a total of 644 visual
244 questions. During the study, we presented users with visual questions that were randomly selected
245 from both the Infoseek [7] and OK-VQA [26] datasets. This approach allowed us to provide the
246 participants with a varied and diverse set of visual questions to assess and respond to. We show the
247 details for this study as well as example prompts in the Appendix.

248 4 Experiments

249 We evaluate AVIS on two visual question answering datasets: *i) OK-VQA [26]*, which requires
250 common-sense knowledge not observed in given image; and *ii) Infoseek_{wikidata} [7]*, which further
251 necessitates more fine-grained information that cannot be covered by common sense knowledge.

252 **Experimental Setup.** We follow the decision-making workflow in Alg. 2 to implement AVIS to solve
253 visual questions. For the Planner, we write the basic instructions for describing each tool, and keep a
254 pool of real user behavior when they select each tool, which we collected in the user study. At each
255 step s , we prepare the prompt based on the feasible action lists \mathcal{A}_s . For the Reasoner, we write the
256 prompt for all APIs that return a long list of results, including *Object Detection*, *Product Detection*,
257 *Web Image Search* and *Web Text Search*, that guides reasoner to extract the relevant information. Note
258 that we design the reasoner in a way such that the “uninformative” answers can be detected. In order
259 to support this, we manually prepare several bad examples that do not provide any useful information,
260 pass it to the reasoner as a part of the prompt. We show the detailed prompts for these two modules
261 in the Appendix.

262 We use the frozen PALM 540B language model [9] for both the planner and the reasoner, with
263 deterministic generation ensured by setting the temperature parameter to zero. We use 10 examples
264 as in-context prompts for each dataset, and report the VQA accuracy [11] as the evaluation metric.

265 **Baselines.** A significant novelty of AVIS is the ability to dynamically determine the relevant
266 tools according to different states. To show that this design choice is useful, we add a number of
267 baselines that do not contain a LLM-planner for dynamic decision making. Instead, they follow a
268 pre-determined sequence to call a list of tools. We propose the following baselines:

Model	Unseen Entity	Unseen Question
PALM [9] (Q-only, few-shot)	3.7	5.1
OFA [22] (fine-tune)	9.7	14.8
PALI [6] (VQA, zero-shot)	1.8	2.2
PALI [6] (fine-tune)	16.0	20.7
PALM [9] w/ CLIP [32] (few-shot + external knowledge)	21.9	18.6
FiD [44] w/ CLIP [32] (fine-tune + external knowledge)	20.7	18.1
—baselines without dynamic decision making, sequentially execute the tools—		
baseline-PALM w/ (PALI*, few-shot)	12.8	14.9
baseline-PALM w/ (PALI* + Object, few-shot)	31.3	36.1
baseline-PALM w/ (PALI* + Object + Search, few-shot)	36.1	38.2
AVIS (ours, few-shot)	50.7	56.4
w/o PALI*	47.9	54.2
w/o Object	41.2	48.4
w/o Search	42.5	49.6

Table 1: **Visual Question Answering** results (accuracy) on Infoseek_{wikidata}. The first four rows are results from their paper that do not use external knowledge, and the next two are from their paper that use CLIP as knowledge source. The tool PALI* denotes the frozen multi-task PALI-17B model for both visual question answering and image captioning. Object means object detection, and search means image and text search.

	Model	Accuracy (%)
Supervised	KRISP [25]	38.4
	KAT [12]	54.4
	ReVIVE [20]	58.0
	REVEAL [15]	59.1
	PALI [6] (OK-VQA, finetune)	<u>64.5</u>
Zero-shot	PALI [6] (VQA, zero-shot)	41.6
	PICa-Full [41]	48.0
	Flamingo (zero-shot) [4]	50.6
	BLIP-2 [18]	45.9
	ViperGPT [36]	51.9
	Flamingo (few-shot) [4]	57.8
Few-shot	(baselines without dynamic decision making, sequentially executing the tools)	
	baseline-PALM w/ (PALI*)	44.3
	baseline-PALM w/ (PALI*+Object)	38.2
	baseline-PALM w/ (PALI*+Object + Search)	47.9
	AVIS (ours)	60.2
	w/o PALI*	47.1
w/o Object	58.3	
w/o Search	55.0	

Table 2: **Visual Question Answering** results (accuracy) on OK-VQA. The tool PALI* denotes the frozen multi-task PALI-17B model for both visual question answering and image captioning. Object means object detection, and search means image and text search.

- 269
- 270
- 271
- 272
- 273
- 274
- 275
- **baseline-PALM w/ PALI***, which integrates the captions generated by PALI and the visual answers from PALI VQA. PALI* denotes the combination of both VQA and captioning tool.
 - **baseline-PALM w/ (PALI* + Object)**, which in addition calls the object detection tool, and then integrates all object data, including products and text detected by OCR.
 - **baseline-PALM w/ (PALI* + Object + Search)**, a model which first selects a relevant object with the help of PALM, then sequentially executes the image search and Google search with the object name. It then calls PALM again to answer the question.

276 For each of the three baselines, we prepare a few-shot Chain-Of-Thought (COT) prompting [39], in
277 which the COT prompt guides the model to explain why predictions are made based on the provided
278 information. Note that these baselines utilize a set of tools in a fixed order, without the capacity for
279 dynamic decision making.

280 We also evaluate the usefulness of each tool group (i.e., PALI*, Object, and Search) through an
281 ablation study. This involves removing each tool group from our framework individually, and
282 assessing the impact on performance.

283 **Experimental Results.** Table 5 presents the results of AVIS and other baselines on the
284 Infoseek_{wikidata} dataset. Infoseek_{wikidata} is a challenging dataset that requires identifying highly
285 specific entities. Even robust visual-language models, such as OFA [22] and PALI [6], fail to yield



Figure 5: Examples of AVIS’s dynamic planning and reasoning procedure for solving visual questions.

286 high accuracy when fine-tuned on this dataset. However, our AVIS, without fine-tuning and by
 287 leveraging a complete set of tools guided by 10 in-context examples, achieves the accuracy of 50.7
 288 and 56.4 on the unseen entity and question splits, respectively. This significantly outperforms the
 289 fine-tuned results of PALI-17B, which are 16.0 and 20.7, as well as the PALM model augmented
 290 with CLIP knowledge, which are 21.9 and 18.6, respectively.

291 Table 5 also illustrates that our improvements are not solely due to the additional information provided
 292 by the external tools, but due to our dynamic decision-making pipeline. We compare the results of
 293 AVIS with the three baselines that conduct sequential execution. While these baselines do improve
 294 the performance, our AVIS framework outperforms the best baseline model by up to 17.3 accuracy.
 295 Note that AVIS and the baselines use exactly the same set of tools. This considerable performance
 296 gap clearly shows the clear advantage of our dynamic decision-making design. Furthermore, we
 297 show the importance of each tool in the last block of Table 5. Removal of any of the tools degrades
 298 the overall accuracy. Among the three tool groups, Object and Search are more important than PALI,
 299 as they provide more fine-grained information crucial for the Infoseek dataset.

300 We report the OK-VQA experiments in Table 2. AVIS with few-shot in-context examples achieves
 301 an accuracy of 60.2, higher than most of the existing methods tailored for this dataset, including
 302 KAT [12], ReVIVE [20] and REVEAL [15]. AVIS achieves lower but comparable performance
 303 compared to PALI model fine-tuned on OK-VQA. This difference, compared to Infoseek, may be
 304 attributed to the fact that most QA examples in OK-VQA rely more on commonsense knowledge
 305 than on fine-grained knowledge. Therefore, it is feasible to encode such generic knowledge in the
 306 model parameters and requires less external knowledge. Note that PALI zero-shot VQA model itself
 307 achieves 41.6 accuracy, which is significantly higher than in Infoseek, which supports this hypothesis.
 308 Table 2 also shows that the object detection is less crucial as a tool on this data set, compared to PALI
 309 captioning and VQA. Nonetheless, AVIS equipped with all tools achieves the best performance.

310 **Case studies for dynamic decision making.** One of the key features of AVIS is its ability to
 311 dynamically make decisions instead of executing a fixed sequence. Figure 5 presents three examples
 312 of AVIS’s dynamic planning and reasoning process. They demonstrate the flexibility of AVIS to use
 313 different tools at various stages. It is also worth noting that our reasoner design enables AVIS to
 314 identify irrelevant information, backtrack to a previous state, and repeat the search. For instance, in
 315 the second example concerning the taxonomy of fungi, AVIS initially makes an incorrect decision
 316 by selecting a leaf object. However, the reasoner identifies that this is not relevant to the question,
 317 prompting AVIS to plan again. This time, it successfully selects the object related to false turkey-tail
 318 fungi, leading to the correct answer, Stereum.

319 5 Conclusion

320 In this paper, we propose a novel approach that equips the Large Language Models (LLM) with the
 321 ability to use a variety of tools for answering knowledge-intensive visual questions. Our methodology,
 322 anchored in human decision-making data collected from a user study, employs a structured framework
 323 that uses an LLM-powered planner to dynamically decide on tool selection and query formation. An
 324 LLM-powered reasoner is tasked with processing and extracting key information from the output of
 325 the selected tool. Our method iteratively employs the planner and reasoner to leverage different tools
 326 until all necessary information required to answer the visual question is amassed.

327 **References**

- 328 [1] Google lens. Web interface available at <https://images.google.com>.
- 329 [2] Google search. Web interface available at <https://www.google.com>.
- 330 [3] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, K. Gopalakrishnan, K. Hausman,
331 A. Herzog, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint*
332 *arXiv:2204.01691*, 2022.
- 333 [4] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican,
334 M. Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural*
335 *Information Processing Systems*, 35:23716–23736, 2022.
- 336 [5] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry,
337 A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing*
338 *systems*, 33:1877–1901, 2020.
- 339 [6] X. Chen, X. Wang, S. Changpinyo, A. Piergiovanni, P. Padlewski, D. Salz, S. Goodman, A. Grycner,
340 B. Mustafa, L. Beyer, et al. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint*
341 *arXiv:2209.06794*, 2022.
- 342 [7] Y. Chen, H. Hu, Y. Luan, H. Sun, S. Changpinyo, A. Ritter, and M.-W. Chang. Can pre-trained vision and
343 language models answer visual information-seeking questions? In *arXiv preprint arXiv:2302.11713*, 2023.
- 344 [8] F. Chern, B. Hechtman, A. Davis, R. Guo, D. Majnemer, and S. Kumar. TPU-KNN: K nearest neighbor
345 search at peak flop/s. *CoRR*, abs/2206.14286, 2022.
- 346 [9] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton,
347 S. Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*,
348 2022.
- 349 [10] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong,
350 T. Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- 351 [11] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh. Making the V in VQA matter: Elevating
352 the role of image understanding in visual question answering. In *2017 IEEE Conference on Computer*
353 *Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 6325–6334.
354 IEEE Computer Society, 2017.
- 355 [12] L. Gui, B. Wang, Q. Huang, A. Hauptmann, Y. Bisk, and J. Gao. Kat: A knowledge augmented transformer
356 for vision-and-language. *arXiv preprint arXiv:2112.08614*, 2021.
- 357 [13] T. Gupta and A. Kembhavi. Visual programming: Compositional visual reasoning without training. In
358 *arXiv preprint arXiv:2211.11559*, 2022.
- 359 [14] H. Hu, Y. Luan, Y. Chen, U. Khandelwal, M. Joshi, K. Lee, K. Toutanova, and M.-W. Chang. Open-
360 domain visual entity recognition: Towards recognizing millions of wikipedia entities. In *arXiv preprint*
361 *arXiv:2302.11154*, 2023.
- 362 [15] Z. Hu, A. Iscen, C. Sun, Z. Wang, K.-W. Chang, Y. Sun, C. Schmid, D. A. Ross, and A. Fathi. Reveal:
363 Retrieval-augmented visual-language pre-training with multi-source multimodal knowledge memory. In
364 *CVPR*, 2023.
- 365 [16] A. Kulshreshtha, D. D. F. Adiwardana, D. R. So, G. Nemade, J. Hall, N. Fiedel, Q. V. Le, R. Thop-
366 pilan, T. Luong, Y. Lu, and Z. Yang. Towards a human-like open-domain chatbot. In *arXiv preprint*
367 *arXiv:2001.09977*, 2020.
- 368 [17] A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Mallocci,
369 A. Kolesnikov, T. Duerig, and V. Ferrari. The open images dataset v4: Unified image classification, object
370 detection, and visual relationship detection at scale. *IJCV*, 2020.
- 371 [18] J. Li, D. Li, S. Savarese, and S. C. H. Hoi. BLIP-2: bootstrapping language-image pre-training with frozen
372 image encoders and large language models. *CoRR*, abs/2301.12597, 2023.
- 373 [19] Y. Li, D. H. Choi, J. Chung, N. Kushman, J. Schrittwieser, R. Leblond, T. Eccles, J. Keeling, F. Gi-
374 meno, A. D. Lago, T. Hubert, P. Choy, C. de Masson d’Autume, I. Babuschkin, X. Chen, P. Huang,
375 J. Welbl, S. Gowal, A. Cherepanov, J. Molloy, D. J. Mankowitz, E. S. Robson, P. Kohli, N. de Fre-
376 itas, K. Kavukcuoglu, and O. Vinyals. Competition-level code generation with alphacode. *CoRR*,
377 abs/2203.07814, 2022.

- 378 [20] Y. Lin, Y. Xie, D. Chen, Y. Xu, C. Zhu, and L. Yuan. Revive: Regional visual representation matters in
379 knowledge-based visual question answering. *arXiv preprint arXiv:2206.01201*, 2022.
- 380 [21] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning. In *arXiv preprint arXiv:2304.08485*, 2023.
- 381 [22] J. Lu, C. Clark, R. Zellers, R. Mottaghi, and A. Kembhavi. Unified-io: A unified model for vision, language,
382 and multi-modal tasks. *CoRR*, abs/2206.08916, 2022.
- 383 [23] P. Lu, B. Peng, H. Cheng, M. Galley, K.-W. Chang, Y. N. Wu, S.-C. Zhu, and J. Gao. Chameleon:
384 Plug-and-play compositional reasoning with large language models. In *arXiv preprint arXiv:2304.09842*,
385 2023.
- 386 [24] A. Madaan, S. Zhou, U. Alon, Y. Yang, and G. Neubig. Language models of code are few-shot common-
387 sense learners. In Y. Goldberg, Z. Kozareva, and Y. Zhang, editors, *Proceedings of the 2022 Conference on*
388 *Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates,*
389 *December 7-11, 2022*, pages 1384–1403. Association for Computational Linguistics, 2022.
- 390 [25] K. Marino, X. Chen, D. Parikh, A. Gupta, and M. Rohrbach. Krisp: Integrating implicit and symbolic
391 knowledge for open-domain knowledge-based vqa. In *Proceedings of the IEEE/CVF Conference on*
392 *Computer Vision and Pattern Recognition*, pages 14111–14121, 2021.
- 393 [26] K. Marino, M. Rastegari, A. Farhadi, and R. Mottaghi. OK-VQA: A visual question answering benchmark
394 requiring external knowledge. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*
395 *2019, Long Beach, CA, USA, June 16-20, 2019*, pages 3195–3204. Computer Vision Foundation / IEEE,
396 2019.
- 397 [27] G. Mialon, R. Dessi, M. Lomeli, C. Nalmpantis, R. Pasunuru, R. Raileanu, B. Rozière, T. Schick,
398 J. Dwivedi-Yu, A. Celikyilmaz, E. Grave, Y. LeCun, and T. Scialom. Augmented language models: a
399 survey. In *arXiv preprint arXiv:2302.07842*, 2023.
- 400 [28] R. Nakano, J. Hilton, S. Balaji, J. Wu, L. Ouyang, C. Kim, C. Hesse, S. Jain, V. Kosaraju, W. Saunders, et al.
401 Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*,
402 2021.
- 403 [29] OpenAI. Gpt-4 technical report. In *arXiv preprint arXiv:2303.08774*, 2023.
- 404 [30] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama,
405 A. Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural*
406 *Information Processing Systems*, 35:27730–27744, 2022.
- 407 [31] A. Parisi, Y. Zhao, and N. Fiedel. Talm: Tool augmented language models. In *arXiv preprint*
408 *arXiv:2205.12255*, 2022.
- 409 [32] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin,
410 J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from natural language
411 supervision. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24*
412 *July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763.
413 PMLR, 2021.
- 414 [33] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, et al. Improving language understanding by
415 generative pre-training. 2018.
- 416 [34] T. L. Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilic, D. Hesslow, R. Castagné, A. S. Luccioni, F. Yvon,
417 M. Gallé, J. Tow, A. M. Rush, S. Biderman, A. Webson, P. S. Ammanamanchi, T. Wang, B. Sagot,
418 N. Muennighoff, A. V. del Moral, O. Ruwase, R. Bawden, S. Bekman, A. McMillan-Major, I. Beltagy,
419 H. Nguyen, L. Saulnier, S. Tan, P. O. Suarez, V. Sanh, H. Laurençon, Y. Jernite, J. Launay, M. Mitchell,
420 C. Raffel, A. Gokaslan, A. Simhi, A. Soroa, A. F. Aji, A. Alfassy, A. Rogers, A. K. Nitzav, C. Xu, C. Mou,
421 C. Emezue, C. Klamm, C. Leong, D. van Strien, D. I. Adelani, and et al. BLOOM: A 176b-parameter
422 open-access multilingual language model. *CoRR*, abs/2211.05100, 2022.
- 423 [35] T. Schick, J. Dwivedi-Yu, R. Dessi, R. Raileanu, M. Lomeli, L. Zettlemoyer, N. Cancedda, and T. Scialom.
424 Toolformer: Language models can teach themselves to use tools. In *arXiv preprint arXiv:2302.04761*,
425 2023.
- 426 [36] D. Surís, S. Menon, and C. Vondrick. Vipergpt: Visual inference via python execution for reasoning. In
427 *arXiv preprint arXiv:2303.08128*, 2023.

- 428 [37] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro,
429 F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*,
430 2023.
- 431 [38] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou,
432 D. Metzler, E. H. Chi, T. Hashimoto, O. Vinyals, P. Liang, J. Dean, and W. Fedus. Emergent abilities of
433 large language models. In *arXiv preprint arXiv:2206.07682*, 2022.
- 434 [39] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. H. Chi, Q. V. Le, and D. Zhou.
435 Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*, 2022.
- 436 [40] C. Wu, S. Yin, W. Qi, X. Wang, Z. Tang, and N. Duan. Visual chatgpt: Talking, drawing and editing with
437 visual foundation models. In *arXiv preprint arXiv:2303.04671*, 2023.
- 438 [41] Z. Yang, Z. Gan, J. Wang, X. Hu, Y. Lu, Z. Liu, and L. Wang. An empirical study of GPT-3 for few-shot
439 knowledge-based VQA. *ArXiv preprint*, abs/2109.05014, 2021.
- 440 [42] Z. Yang, L. Li, J. Wang, K. Lin, E. Azarnasab, F. Ahmed, Z. Liu, C. Liu, M. Zeng, and L. Wang. Mm-react:
441 Prompting chatgpt for multimodal reasoning and action. In *arXiv preprint arXiv:2303.11381*, 2023.
- 442 [43] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao. React: Synergizing reasoning and
443 acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- 444 [44] D. Yu, C. Zhu, Y. Fang, W. Yu, S. Wang, Y. Xu, X. Ren, Y. Yang, and M. Zeng. KG-FiD: Infusing
445 knowledge graph in fusion-in-decoder for open-domain question answering. In *Proceedings of the 60th*
446 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4961–
447 4974, Dublin, Ireland, 2022. Association for Computational Linguistics.

448 **A Implementation of AVIS workflow**

449 We implemented AVIS using the code snippet referenced in Code 7. Throughout our experiments,
450 we employed the APIs of Google Search, LENS, PALI, and PALM directly, without the need for
451 additional GPU/TPU computational resources. Tools that didn't require input queries, such as object
452 detection, captioning, and image search, had their results pre-calculated over the two datasets to
453 reduce the time cost. Other services like VQA, text search, and LLM QA were called during runtime.

454 **B Comparison to pure Autonomous baseline without Transition Graph**

455 One of the significant contributions of this paper lies in the use of a transition graph, synthesized
456 from an authentic user study. To underscore the importance of this graph, along with user prompts in
457 facilitating the efficacy of AVIS, we devised a baseline that operates independently of the transition
458 graph. In this scenario, the model, at each timestep, is presented with a comprehensive list of all
459 tools, each paired with a task description. This baseline shares similarities with the recently launched
460 AutoGPT¹, BabyAGI² projects, which attempted to utilize LLMs as autonomous agents to select all
461 possible actions available in the web.

462 The results are show in Table 3 on Infoseek WikiData unseen entity set and OKVQA. Note that this
463 baseline doesn't achieve the number as high as AVIS with the transition graph and user prompts.
464 The key reason for this discrepancy is the global characteristics inherent in the tool list we have.
465 For instance, we typically first address the visual sub-question through object detection and image
466 search, followed by resolving the knowledge component via Google Search and LLM. However,
467 solely relying on the task description, devoid of human behavior as guidance, can result in the model
468 generating unrealistic tools. We will discuss this intuition more in the following sections.

Model	Infoseek	OKVQA
AVIS w.o/ Transition Graph	38.2	47.3
AVIS w/ Transition Graph	50.7	60.2

Table 3: Ablation of AVIS with or without the guidance of Transition Graph

469 **C Analysis of AVIS's generated tool execution sequence**

470 We have also conducted an analysis to determine whether common patterns exist within the generated
471 programs of AVIS's predictions.

472 We gathered the tool execution traces for all samples within the Infoseek unseen entity dataset.
473 Initially, we display the frequency of each tool being invoked in Figure 6, followed by a more detailed
474 analysis of the first to fourth most commonly called tools in Figures 7-10. As illustrated, the AVIS
475 model, guided by the transition graph and prompts, does not utilize all possible combination of tools,
476 but favors some certain combinations. For instance, as depicted in Fig 7, "object select" is utilized
477 more frequently than other tools at the outset. Similarly, as demonstrated in Fig 9, during the third
478 step, when the model accumulates the visual answer, it is likely to invoke "web search" to gather
479 additional information.

480 We have also calculated the transition probability of the induced graph in Fig 11. The structure
481 of this graph differs slightly from the guided transition graph because during actual runtime, the
482 model will not predict some of the edges. Overall, it reveals a clear two-step question-solving pattern.
483 Initially, AVIS gathers sufficient visual information through the use of visual tools such as "object
484 detection," "VQA," or "identical image search," and then employs "LLM QA" to obtain the visual
485 answer. Subsequently, it iteratively calls "web search" and "LLM QA" post-search with a prompt,
486 eventually deriving the final answer. We also present the distribution of the lengths of generated
487 sequences in Figure 13. As illustrated, the lengths vary considerably, rather than maintaining a fixed
488 value, with a length of 5 being most common for the generated sequences.

¹<https://github.com/Significant-Gravitas/Auto-GPT>

²<https://github.com/yoheinakajima/babyagi>

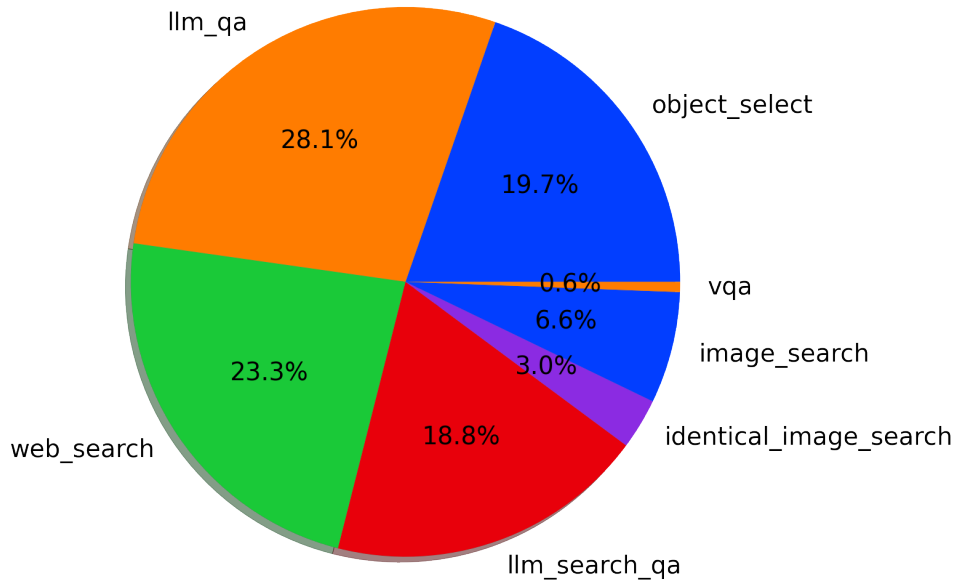


Figure 6: Overall frequency of tool usage on Infoseek dataset.

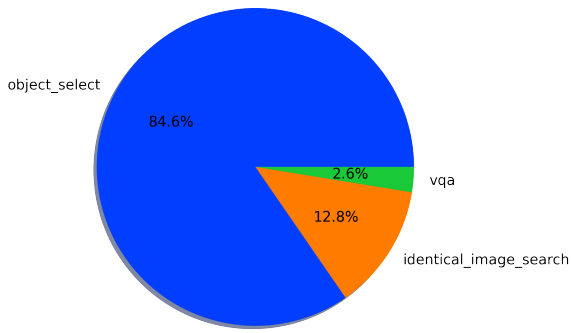


Figure 7: Frequency of the first used tool.

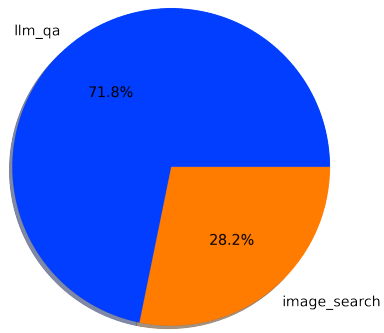


Figure 8: Frequency of the second used tool.

489 Another intriguing aspect worth exploring is our reasoner component. As explained in the paper, the
 490 reasoner evaluates whether the output of each tool is "informative," "not informative," or "answerable".
 491 We exhibit the overall frequency of these predictions in Figure 12. As shown, the model tends to
 492 classify most of the outputs as either informative or answerable. However, approximately 8.1%
 493 of returned entries are deemed "not informative," in which case AVIS would backtrack to select
 494 alternative actions. We further demonstrate a few examples of different choices in Table 4.

495 D Dataset Details

496 **Infoseek**³ is a Visual Question Answering (VQA) dataset, specifically geared towards information-
 497 seeking questions that cannot be answered merely through common sense knowledge. This dataset
 498 was curated by initially gathering human-annotated questions, which were then automatically inte-
 499 grated with existing visual entity recognition datasets and Wikidata to generate complex question-

³<https://open-vision-language.github.io/infoseek/>

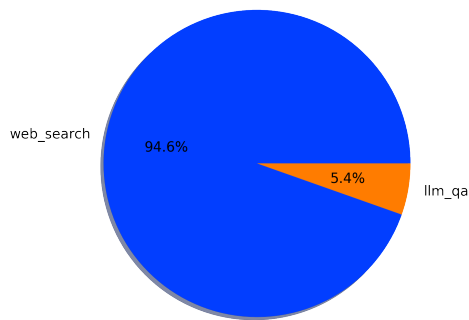


Figure 9: Frequency of the third used tool.

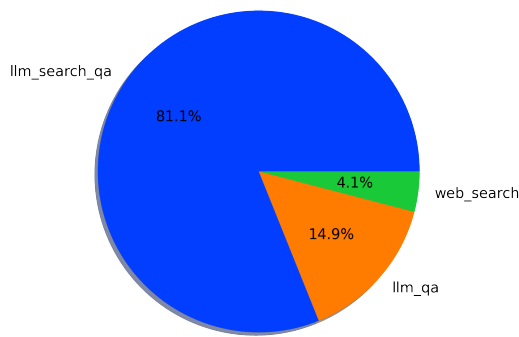


Figure 10: Frequency of the fourth used tool.

Question	Action	Returned Content	Output of Reasoner
where is the island located in?	image search	Retrievals: [willow - Students Britannica Kids Homework Help (score=57.2), Silky Willow Live Stakes For Sale Buy Live Stakes Wholesale (score=55.3), ...]	not informative
Who or what maintains this place?	LLM QA	the history of trafilgar square. the most relevant knowledge entry about trafilgar square is "Trafalgar Square is the focal point for protests, revelry and victory celebrations, however it was only in 1835 that the area became known by this name". From the result, we know the history of trafilgar square is trafilgar square is the focal point for protests, revelry and victory celebrations, however it was only in 1835 that the area became known by this name. Terefore, the predicted answer is trafilgar square is the focal point for protests , revelry and victory celebrations , however it was only in 1835 that the area became known by this name.	not informative
What is the total quantity of produced items for this type of aircraft?	image search	Retrievals: [high confidence: Condor bietet neue Langstrecke nach Brasilien an - AERO International (score=97.3), high confidence: Air safety incidents for aircraft registration D-ABOE - AeroInside (score=95.0), ...]	yes, answerable
what is the name of this mountain?	image search	Object: [a view of a castle in the distance . There are trees and buildings present at the bottom of this image. We can see a hill and a tower in the middle of this image. We can see the sky in the background. (Caption, whole image), Mount of Olives (ridge): The Mount of Olives or Mount Olivet is a mountain ridge east of and adjacent to Jerusalem's Old City (score=88.6), Mount Zion (peak): Mount Zion is a hill in Jerusalem, located just outside the walls of the Old City (score=79.0)]	yes, informative

Table 4: Several examples of API execution results and the reasoner's justification.

500 answer pairs. At the time of submission, we only have access to its wikidata split. Here we also
 501 report the results on human split in Table 5.

502 **OK-VQA**⁴ is another VQA dataset, unique in its requirement for the application of external knowl-
 503 edge that transcends the information directly visible in the input images. The creation of this dataset
 504 involved crowdsourced workers who were tasked with annotating complex questions, drawing upon
 505 the extensive knowledge resources available on Wikipedia.

506 E Prompt Examples

507 Below we show different prompt examples to support our AVIS workflow. First is the prompts for
 508 planning, which selects which tool to use and what query to send. It is consists of a overall task
 509 descriptions and many real examples showing at which circumstances real users select this tool.

```

510 1planner_prompt =
511 2"""You goal is to answer the following query: %s.
512 3
513 4To answer it, you will be provided with the following tools:
514 5%s
515 6
516 7Please make the decision based on the current context.
517 8
518 9%s
  
```

⁴<https://okvqa.allenai.org/>

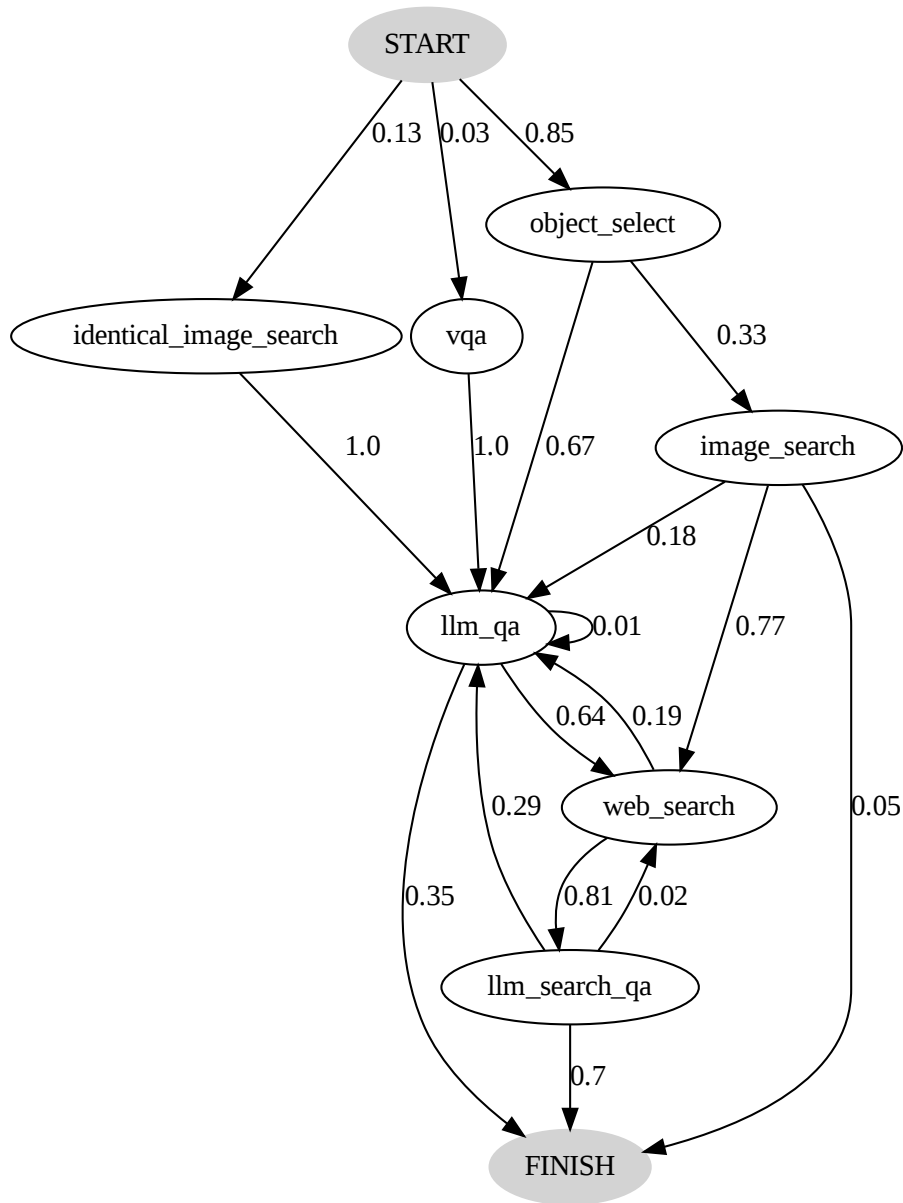


Figure 11: Induced transition frequency graph of AVIS over Infoseek dataset.

```

519 10Query: %s
520 11Context: %s
521 12Action: \n
522 13""
523 14
524 15task_instructions = {
525 16'vqa':
526 17    'You will ask simple question about this image to a external QA module. Please use this when the input
527    query is very straightforward and simple.',\
528 18'object_select':
529 19    'You will select one of the object we detect to dig further. Please use when the question asks about a
530    specific object.',\
531 20'identical_image_search':

```

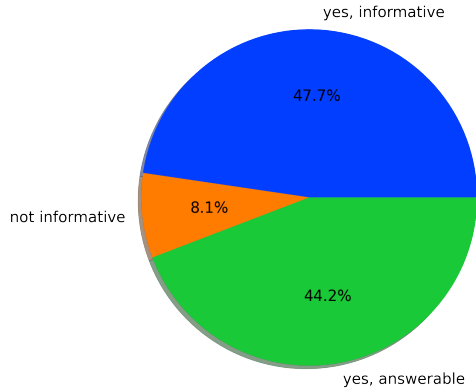



Figure 12: Overall frequency of judgement by reasoner of AVIS.

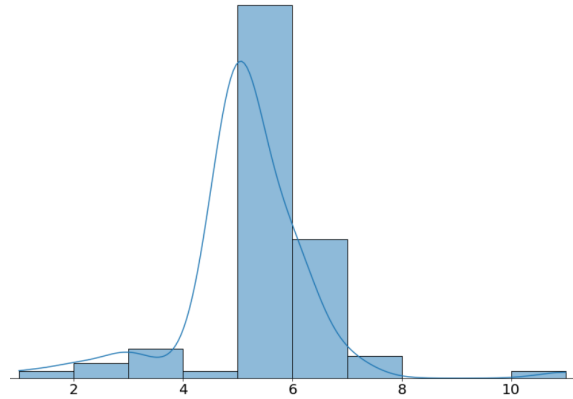


Figure 13: Length distribution of AVIS's generated action sequences.

Model	Unseen Entity	Unseen Question
PALM (Q-only, few-shot)	6.6	4.8
OFA (fine-tune)	2.9	6.2
PALI (fine-tune)	5.9	13.3
PALM w/ CLIP (few-shot + external knowledge)	14.9	15.6
FiD w/ CLIP (fine-tune + external knowledge)	17.6	18.9
AVIS (ours, few-shot)	31.4	33.6

Table 5: **Visual Question Answering** results (accuracy) on $\text{Infoseek}_{\text{human}}$. The first four rows are results from their paper that do not use external knowledge, and the next two are from their paper that use CLIP as knowledge source.

```

532 21 'You will see captions of all images identical to the given image. Please use when the question asks
533     about the whole image instead of a part.',\
534 22 'image_search':
535 23 'You will see captions of all images similar to this object. Please use when you need more information.',\
536 24 'web_search':
537 25 'You will send question to Google Search to get knowledge. Please use when the current query requires
538     extra knowledge',\
539 26 'llm_qa':
540 27 'You will send question to a QA module. Please use this when the input query is simple and contain
541     common-sense knowledge'
542 28 }

```

Listing 1: Planner prompt skeleton and Task instructions

```

543 1vqa_plan_prompts = [
544 2""Query: what is the train carrying?
545 3Context: [
546 4 a train traveling down train tracks next to a forest . There are four trains on the railway track. In the
547 background there are trees,poles and sky. (Caption, whole image)
548 5 Extracted Text: BNSF (score=100.0),
549 6 BNSF Railway: BNSF Railway is one of the largest freight railroads in North America (score=89.3),
550 7]
551 8Action: vqa
552 9""",\
553 10""Query: What is the girl wearing on her legs?
554 11Context: [
555 12 a woman standing in a field putting on a coat . There is a woman standing on the ground. This is grass and
556 there are plants. In the background we can see some trees and this is sky. (Caption, whole image)
557 13]
558 14Action: vqa
559 15""",\
560 16""Query: what color is the bus?
561 17Context: [
562 18 a double decker bus parked in front of a building . There is a double decker bus on the road and this is
563 snow. Here we can see a pole, light, trees, and houses. In the background there is sky. (Caption, whole
564 image)
565 19 Extracted Text: ENVIRO400 (score=100.0),
566 20 Extracted Text: Les Miserables (score=100.0),
567 21 Query Suggestion: les miserables (score=100.0),
568 22 Volvo Olympian: The Volvo Olympian was a rear-engined 2-axle and 3-axle double decker bus chassis
569 manufactured by Volvo at its Irvine, Scotland factory (score=88.5),
570 23 Alexander Dennis Enviro400: The Alexander Dennis Enviro400 is a twin-axle low-floor double-decker bus that
571 was built by the British bus manufacturer Alexander Dennis between 2005 and 2018 (score=85.4),
572 24]
573 25Action: vqa
574 26""",\
575 27""Query: what is the person doing?
576 28Context: [
577 29 two people sitting on the floor opening presents . There are sofas on the sofas there are pillows, here
578 there is table, on the table there are plants and other objects, here there are two persons sitting on
579 the ground, gift boxes, dog and this is floor. (Caption, whole image)
580 30]
581 31Action: vqa
582 32""
583 33]
584 34object_select_plan_prompts = [
585 35""Query: what is the name of this building?
586 36Context: [
587 37 a group of people that are standing in front of a building . There is a building in the left corner which
588 has few people standing in front of it and there is a fire hydrant in the right corner and there is a
589 street light pole beside it. (Caption, whole image)
590 38 Query Suggestion: Alcatraz Warden's House San Francisco (score=100.0),
591 39 Alcatraz Island (historic_site): Alcatraz Island is a small island 1 (score=91.9),
592 40 Warden's House: The Warden's House was the home of the wardens of the federal penitentiary on Alcatraz
593 Island, off San Francisco (score=78.1),
594 41]
595 42Action: object_select
596 43""
597 44""Query: what is the island?
598 45Context: [
599 46 a view of a mountain from a cable car . There is a ropeway. Behind that there are trees and hills.
600 (Caption, whole image)
601 47 Ngong Ping 360 (gondola_lift_station): Ngong Ping 360 is a bicable gondola lift on Lantau Island in Hong
602 Kong (score=91.8),
603 48 Tian Tan Buddha (monument): The Big Buddha is a large bronze statue of Buddha, completed in 1993, and
604 located at Ngong Ping, Lantau Island, in Hong Kong (score=79.0),
605 49]
606 50Action: object_select
607 51""
608 52""Query: what is the name of this place?
609 53Context: [
610 54 a cemetery with a building in the background . There is a road and there are many atoms and trees beside it
611 and there is a building in the right corner. (Caption, whole image)
612 55]
613 56Action: object_select
614 57""
615 58""Query: what is the name of this bird?
616 59Context: [
617 60 a bird sitting on top of a lush green hillside . There is a bird on the grassland in the foreground area of
618 the image and the background is blurry. (Caption, whole image)
619 61 Atlantic puffin (type_of_bird): The Atlantic puffin, also known as the common puffin, is a species of
620 seabird in the auk family (score=73.2),
621 62 Horned puffin (type_of_bird): The horned puffin is an auk found in the North Pacific Ocean, including the
622 coasts of Alaska, Siberia and British Columbia (score=73.2),
623 63 Puffins (type_of_bird): Puffins are any of three species of small alcid in the bird genus Fratercula
624 (score=73.2),
625 64 Fraterculini (score=48.8),
626 65 Auk (type_of_bird): An auk or alcid is a bird of the family Alcidae in the order Charadriiformes
627 (score=11.8),
628 66]
629 67Action: object_select
630 68""
631 69]
632 70identical_image_search_plan_prompts = [
633 71""Query: what is the name of this building?
634 72Context: [

```

```

635 73 a group of people that are standing in front of a building . There is a building in the left corner which
636     has few people standing in front of it and there is a fire hydrant in the right corner and there is a
637     street light pole beside it. (Caption, whole image)
638 74 Query Suggestion: Alcatraz Warden's House San Francisco (score=100.0),
639 75 Alcatraz Island (historic_site): Alcatraz Island is a small island 1 (score=91.9),
640 76 Warden's House: The Warden's House was the home of the wardens of the federal penitentiary on Alcatraz
641     Island, off San Francisco (score=78.1),
642 77]
643 78Action: identical_image_search
644 79""",
645 80""Query: what is the aircraft?
646 81Context: [
647 82 a fighter jet sitting on top of an airport tarmac . There is a plane and missiles on the ground. At the
648     left a person is standing wearing a cap. (Caption, whole image)
649 83 Extracted Text: AIRLINERS.NET (score=100.0),
650 84 Query Suggestion: airliners.net (score=100.0),
651 85 Airliners: Airliners (score=74.8),
652 86 British Aerospace Hawk 200: The British Aerospace Hawk 200 is a single-seat, single engine light multirole
653     fighter designed for air defence, air denial, anti-shipping, interdiction, close air support, and
654     ground attack (score=74.8),
655 87 product: Airfix BAE Hawk T1 1:72 (score=0.0),
656 88 product: Rolls-royce Adour In The Hawk / Bae Hawk 200 . Pdf/download (score=0.0),
657 89]
658 90Action: identical_image_search
659 91""",
660 92""Query: what is the name of this place?
661 93Context: [
662 94 a row of pillars sitting next to a dirt road . There is a building and this is plant. Here we can see
663     pillars and a sky. (Caption, whole image)
664 95 Query Suggestion: Palmyra Archaeology (score=100.0),
665 96 Great Colonnade at Palmyra (ancient_roman_architecture_structure): The Great Colonnade at Palmyra was the
666     main colonnaded avenue in the ancient city of Palmyra in the Syrian Desert (score=90.3),
667 97]
668 98Action: identical_image_search
669 99""",
670 100""Query: what is the name of this lake?
671 101Context: [
672 102 a view of a river surrounded by mountains . There are trees in the right corner and there is a river and
673     mountains in front of it. (Caption, whole image)
674 103 Monte Bre (peak): Monte Bre is a small mountain east of Lugano on the flank of Monte Boglia with a view of
675     the bay of Lugano and the Pennine Alps and the Bernese Alps (score=85.5),
676 104 product: Top Searched (score=0.0),
677 105]
678 106Action: identical_image_search
679 107""",
680 108]
681 109action_prompt_dict = {'vqa': vqa_plan_prompts, 'object_select': object_select_plan_prompts,
682     'identical_image_search': identical_image_search_plan_prompts, 'image_search':
683     image_search_plan_prompts, 'web_search': web_search_plan_prompts,
684 110'llm_qa': llm_qa_plan_prompts}

```

Listing 2: Planning Prompts Example

685 We then show how AVIS decompose question into a visual sub-question and a knowledge sub-question.
686 This is done at beginning to guide later tool usage.

```

687 1question_decomposition_prompt = """
688 2     Read the following question for a given image. Decompose the question into two sub-questions.
689 3
690 4     The first will ask information about the image, and the second requires reasoning over the textual
691     knowledge.
692 5     In the second question, we use # to denote the answer of the first question.
693 6
694 7
695 8     Question: what chemical makes the vegetable orange?
696 9     Visual: which orange vegetable is shown?
697 10    Knowledge: chemical makes # orange?
698 11
699 12
700 13    Question: How long can their horns grow?
701 14    Visual: which animals are shown?
702 15    Knowledge: How long can #'s horns grow?
703 16
704 17
705 18    Question: What is a competition for these animals called?
706 19    Visual: which animals are shown?
707 20    Knowledge: competition for #?
708 21
709 22
710 23    Question: What is the name of the ancient greek sport that evolved into the sport featured above?
711 24    Visual: which sport is played?
712 25    Knowledge: name of the ancient greek sport that evolved into #?
713 26
714 27
715 28    Question: Which food item here has the most protein?
716 29    Visual: what are the food items shown?
717 30    Knowledge: Which food item of # has the most protein?
718 31
719 32
720 33    Question: How many calories are in this meal?

```

```

721 34 Visual: what are the food items shown?
722 35 Knowledge: calories in #?
723 36
724 37
725 38 Question: What type of sandwich is this?
726 39 Visual: which type of sandwich is shown?
727 40 Knowledge: #
728 41
729 42 Question: What is the name of the restaurant where this was served?
730 43 Visual: which food items are served?
731 44 Knowledge: restaurant where # was served?
732 45
733 46
734 47 Question: What genus of bird is flying here?
735 48 Visual: what genus of bird is flying?
736 49 Knowledge: #
737 50
738 51
739 52 Question: What is the main ingredient in this food?
740 53 Visual: which food is shown?
741 54 Knowledge: main ingredient in #?
742 55 ""

```

Listing 3: Question Decomposition Prompts

743 Below are several examples to help AVIS learn how to select the most suitable object ID.

```

744 1 object_select_prompt = ""
745 2 Please think step by step. In the following, you will be given a "Query", a list of "Objects".
746 3
747 4 Your task is to predict the object #ID that is mostly relevant to answer the queries. Please generate the
748 5 detailed explanation why you select this object, and then output ID in "Object #ID".
749 6
750 7
751 7 Query: which city is this place?
752 8 Object #0 [
753 9 a row of pillars sitting next to a dirt road . There is a building and this is plant. Here we can see
754 10 pillars and a sky. (Caption, whole image)
755 11 Query Suggestion: Palmyra Archaeology (score=100.0),
756 12 Great Colonnade at Palmyra (ancient_roman_architecture_structure): The Great Colonnade at Palmyra was the
757 13 main colonnaded avenue in the ancient city of Palmyra in the Syrian Desert (score=90.3),
758 14 ]
759 15 Object #1 [
760 16 a green plant sitting next to a brick wall . There is a plant and this is wall. And there is a sky.
761 17 (Caption, center)
762 18 Date palm (type_of_palm_trees): Phoenix dactylifera, commonly known as date palm, is a flowering plant
763 19 species in the palm family, Arecaceae, cultivated for its edible sweet fruit called dates (score=81.7),
764 20 ]
765 21 Object #2 [
766 22 a wicker basket sitting on top of a rock . There is a blur image of a rock. (Caption, lower right)
767 23 ]
768 24 Output: The query asks about the city of the place. Only Object #0 (whole image) mentions city name Palmyra,
769 25 which is an ancient city. Also, Object #0 contains Query Suggestion "Palmyra Archaeology".
770 26 Therefore, the predicted Object #ID is 0.
771 27
772 28
773 28 Query: where is this place?
774 29 Object #0 [
775 30 a view of a valley surrounded by mountains . There are hills and this is grass. Here we can see trees and
776 31 this is sky. (Caption, whole image)
777 32 ]
778 33 Object #1 [
779 34 a view of a lush green hillside with trees . There is a house on the rock and there are few plants beside
780 35 it and there is a greenery ground in the background. (Caption, center)
781 36 Monterey Pine (type_of_conifers): Pinus radiata, the Monterey pine, insignis pine or radiata pine, is a
782 37 species of pine native to the Central Coast of California and Mexico (score=49.1),
783 38 European rabbit (type_of_leporids): The European rabbit or coney is a species of rabbit native to the
784 39 Iberian Peninsula, western France, and the northern Atlas Mountains in northwest Africa (score=31.3),
785 40 ]
786 41 Object #2 [
787 42 a green plant growing on a rocky surface . There is a blur image of trees and rocks. (Caption, lower center)
788 43 product: GreenView Fairway Formula Seed Success Paillis biodegradable avec engrais Sac de 4,5 kg Couvre 200
789 44 m2 (score=0.0),
790 45 ]
791 46 Object #3 [
792 47 a rocky hillside with lots of green vegetation . There are trees and this is rock. (Caption, lower left)
793 48 Willow: Willows, also called salallows and osiers, of the genus Salix, comprise around 350 species of
794 49 typically deciduous trees and shrubs, found primarily on moist soils in cold and temperate regions
795 50 (score=31.3),
796 51 Tamarisk: The genus Tamarix is composed of about 50-60 species of flowering plants in the family
797 52 Tamaricaceae, native to drier areas of Eurasia and Africa (score=26.8),
798 53 ]
799 54 Output: The query asks about the location of this place. Although these entries doesn't explicitly contain
800 55 location name, but Object #1 (center) contains Monterey Pine and European rabbit, which might hint the
801 56 location later.
802 57 Therefore, the predicted Object #ID is 1.
803 58 ""

```

Listing 4: Object Select Prompts

804 Below are the prompts to extract answer from objects and extracted captions of similar images.

```
805 1reason_vqa_prompt = ""
806 2Please think step by step. In the following, you will be given:
807 3
808 4- Query: The query to be asked.
809 5- Think: Why the following knowledge is retrieved.
810 6- Entity: A list of entities that describe the object.
811 7- Retrievals: A list of web documents that are similar to the object. If there's "high confidence", it's very
812 important.
813 8
814 9Your task is to predict a short answer to the query based on the provided information. You need to first
815 identify which knowledge entry is mostly relevant, and then extract the answer from the knowledge.
816 10Rely on Object information more, and if there contains "Query Suggestion", try to use it. Otherwise, if a
817 information appears lots of time, there's a higher chance it's the answer.
818 11After explaining your decision choice, saying "Answer is" and appending your predicted short answer. Please
819 also generate the type of the answer after a comma.
820 12If you are uncertain about the answer, especially when the knowledge is irrelevant to the query, say "cannot
821 be answered". Do not generate the answer not inside the provided knowledge.
822 13
823 14
824 15
825 16Query: what is this building?
826 17Think: object (whole image) contains stockholm city hall, which is the seat of stockholm municipality in
827 stockholm, sweden.
828 18Object: [
829 19 Stockholm City Hall (city_hall): Stockholm City Hall is the seat of Stockholm Municipality in Stockholm,
830 Sweden (score=96.1),
831 20 Bla Hallen (banquet_hall): The Blue Hall is the main hall of the Stockholm City Hall best known as the
832 banquet hall for the annual Nobel Banquet, and also used for state visits, student balls, jubilees and
833 other large events (score=79.0),
834 21]
835 22Retrievals: [
836 23 high confidence: City Hall - Blue Hall (1) | Stockholm (2) | Pictures | Sweden in Global-Geography
837 (score=47.8),
838 24 high confidence: le salon bleu a city hall (salle de remise des prix nobel) - Picture of Stockholm,
839 Stockholm County - Tripadvisor (score=47.7),
840 25]
841 26
842 27Output: The query asks about the building. From both Object and Retrievals, there are mentions about
843 Stockholm City Hall and Blue Hall. As Stockholm City Hall contains Blue Hall, the answer shall be
844 Stockholm City Hall.
845 28Therefore, the predicted answer is Stockholm City Hall.
846 29
847 30
848 31Query: which sport is played?
849 32Think: Object shows a snail sitting on top of a tennis ball.
850 33Object: [
851 34 Cantareus apertus (type_of_gastropods): Cantareus apertus, commonly known as the green garden snail, is a
852 species of air-breathing land snail, a terrestrial pulmonate gastropod mollusc in the family Helicidae,
853 the typical snails,
854 35 Garden snail (type_of_gastropods): Cornu aspersum, known by the common name garden snail, is a species of
855 land snail in the family Helicidae, which includes some of the most familiar land snails,
856 36 Helix aspersa aspersa (type_of_gastropods),
857 37 Slug; Slug, or land slug, is a common name for any apparently shell-less terrestrial gastropod mollusc,
858 38 Snail: A snail is a shelled gastropod,
859 39]
860 40Retrievals: [
861 41 2019 NEWBIE Competition Winner Steven Ryan, Snail Farming - YouTube,
862 42 Alive specimens. a. Megalobulimus ovatus (CMIOC 11136), b. Thaumastus... | Download Scientific Diagram,
863 43 Brown garden snail > Manaaki Whenua,
864 44 Common garden snail and baby,
865 45 Easy Everyday Food for Garden Snails - Ask the plantician,
866 46 Green Life Soil: Natural pest & disease control in a winter garden,
867 47 Helminthoglyptinae - Wikipedia,
868 48 Hydrosalpingitis in broilers - Veterinaria Digital,
869 49 Master Gardener: Protecting squash and cucumbers from slugs and snails - Press Enterprise,
870 50 Mother Baby Blue Snails On Phalaenopsis Stock Photo 530400856 | Shutterstock,
871 51]
872 52
873 53Output: The query asks about sport. From both entities and retrievals, they only talk about snail, and there
874 is no information about which sport is played.
875 54Therefore, given the provided information, this query cannot be answered.
876 55
877 56
878 57Query: which sport is played?
879 58Think: object , object , and object all contain people playing basketball. however, object is the only one
880 that contains a group of women playing basketball.
881 59therefore, the predicted object #id is 0.
882 60Retrievals: [
883 61 08.07.2011 Zanele Mdoana of South Africa in action during the Quarter-finals between New Zealand and South
884 Africa, Mission Foods World Netball Championships 2011 from the Singapore Indoor Stadium in Singapore
885 Stock Photo - Alamy,
886 62 55 Brazilian Handball Team Images, Stock Photos & Vectors | Shutterstock,
887 63 ::Malawi High Commission::,
888 64 Amanda Mynhardt Photostream | Netball, Netball singapore, Netball south africa,
889 65 Australia pass Malawi test with flying colours at Netball World Cup | Netball World Cup 2019 | The Guardian,
890 66 Australia's Jo Weston (second left) and Barbados' Latonia Blackman in action during the Netball World Cup
891 match at the M&S Bank Arena, Liverpool Stock Photo - Alamy,
892 67 Birmingham 29795 World Netball Championships Final Editorial Stock Photo - Stock Image | Shutterstock,
893 68 Bridget kumwenda malawi netball hi-res stock photography and images - Alamy,
```

```

894 69 England V Australia International Netball Series Photos and Premium High Res Pictures | Netball, Netball
895 quotes, Inspirational women,
896 70 File:XX0992 - Madrid basketball Donna Burns - 3b - Scan.jpg - Wikimedia Commons,
897 71]
898 72
899 73Output: The query asks about which sport is played. From retrievals, there exist many mentions about netball,
900 and mentions that they are played by women.
901 74therefore, the predicted answer is women netball.
902 75
903 76
904 77
905 78Query: what is the name of the insect?
906 79Think: only object (while image) mentions the name of the insect, western tiger swallowtail.
907 80Object: [
908 81 Query Suggestion: Western Tiger Swallowtail (score=100.0),
909 82 Canadian tiger swallowtail (type_of_lepidoptera): Papilio canadensis, the Canadian tiger swallowtail, is a
910 species of butterfly in the family Papilionidae (score=78.4),
911 83 Eastern tiger swallowtail (us_state_butterfly): Papilio glaucus, the eastern tiger swallowtail, is a
912 species of butterfly native to eastern North America (score=78.4),
913 84]
914 85Retrievals: [
915 86 high confidence: kupu-kupu - Wiktionary (score=100.0),
916 87 high confidence: Top Spots for Nature Watching and Birding | VisitMaryland.org (score=100.0),
917 88 high confidence: File:Eastern Tiger Swallowtail Papilio glaucus on Milkweed 2800px.jpg - Wikimedia Commons
918 (score=99.8),
919 89 high confidence: Photographing Butterflies - Life in the Finger Lakes (score=97.8),
920 90]
921 91
922 92Output: The query asks about the name of the insect. From Object, it contains a very informative "Query
923 Suggestion: Western Tiger Swallowtail".
924 93Therefore, the predicted answer is Western Tiger Swallowtail.
925 94
926 95""

```

Listing 5: Reason Prompt (Visual Question)

927 Below are prompts AVIS extract answer from search results:

```

928 1reason_ga_prompt = ""
929 2 Please think step by step. In the following, you will be given a "Query", and a list of "Knowledge" from
930 Google Search related to this query.
931 3
932 4 Your task is to predict a short answer to the query based on the provided information. You need to first
933 identify the most relevant knowledge entry, and then predict a short answer based on the knowledge. If
934 a information appears lots of time, there's a higher chance it's the answer.
935 5
936 6 After explaining your decision choice, saying "Answer is" and appending your predicted answer.
937 7 If you are uncertain about the answer, especially when the knowledge is irrelevant to the query, say
938 "cannot be answered". Do not generate the answer not inside the provided knowledge.
939 8
940 9
941 10Query: What chemical makes carrot orange?
942 11Knowledge: [
943 12Title: How did carrots become orange? - The Economist
944 13Content: High Confidence Response: carotenoids.
945 14
946 15Context: The chemical compounds that give carrots their vivid colour, carotenoids, are usually used by plants
947 that grow above ground to assist in the process of photosynthesis.
948 16
949 17Title:
950 18Content: carotenoids
951 19
952 20The chemical compounds that give carrots their vivid colour, carotenoids, are usually used by plants that
953 grow above ground to assist in the process of photosynthesis.
954 21
955 22Title: Can Eating Too Many Carrots Make Your Skin Turn Orange? | Britannica - Encyclopedia Britannica
956 23Content: Maybe not! Carrots and other orange fruits and vegetables are rich in a pigment known as
957 beta-carotene. In humans, this pigment is converted to vitamin A by specialized cells in the small
958 intestine. When high levels of beta-carotene are consumed, not all of the pigment is converted to
959 vitamin A.
960 24Fortunately, the skin discoloration fades when the diet is changed and the levels of beta-carotene in the
961 blood decline.
962 25
963 26Title: Why are carrots orange? | Ask Dr. Universe | Washington State University
964 27Content: Orange carrots are packed with chemicals called carotenoids—specifically, beta-carotene. Your body
965 turns beta-carotene into vitamin A, which helps you grow and protects you from getting sick.
966 Beta-carotene isn't just nutritious. It's also loaded with orange pigment.
967 28That's why vegetables with lots of beta-carotene-like sweet potatoes, squash, and pumpkins-share the same
968 color. But what about that rainbow of other carrot colors? They have their own special qualities, too.
969 Purple carrots get their color from
970 29]
971 30Output: The query asks about chemical that makes carrot orange. Because there's one high confidence result,
972 the most relevant knowledge entries about such chemical is "High Confidence Response: carotenoids."
973 31From this result we know the chemical shall be carotene.
974 32Therefore, the predicted answer is carotene.
975 33
976 34
977 35
978 36Query: What is the name of the drainage basin of ounasjoki?
979 37Knowledge: [
980 38Title: Ounasjoki - Wikipedia

```

981 39Content: It is also the largest river entirely within its borders. Ounasjoki is approximately 299.6
982 kilometres (186.2 mi) in length, and the catchment area is 13,968 square kilometres (5,393 sq mi), 27%
983 of the Kemijoki catchment area.

984 40Tributaries

985 41

986 42- Nakkalajoki.

987 43- Kakkalojoki.

988 44- Syva Tepastojoki.

989 45- Loukinen.

990 46- Meltausjoki.

991 47Course. The Ounasjoki originates at Ounasjarvi lake in Enontekio. It flows first eastwards through
992 Perilajarvi lake and turns south after some seven kilometres. The river then follows southern-sou
993 48

994 49Title: DRAINAGE BASIN OF THE BALTIC SEA - UNECE

995 50Content: Vistula. 194,424. Baltic Sea. BY, PL, SK, UA. - Bug. 39,400. Vistula. BY, PL, UA. - Dunajec. 4726.7.
996 Vistula. PL, SK. -Poprad. 2,077. Dunajec. PL, SK. Oder. 118,861. Baltic Sea. CZ, DE, PL. - Neisse ...
997 Oder. CZ, DE, PL. - Olse ... Oder. CZ, PL. 1 The assessment of water bodies in italics was not included
998 in the present publication. 2 For the Venta River Basin District, which includes the basins of the
999 Barta/Bartuva and Sventoji rivers. Oulu. Lulea. Rovaniemi. Lake. Oulujarvi. Lake. Tornetrask. Torne.
1000 Oulujoki.

1001 51]

1002 52Output: The query asks about drainage basin of ounasjoki. The most relevant knowledge entry that contain
1003 basin is "Venta River Basin District, which includes the basins of the Barta/Bartuva and Sventoji
1004 rivers."

1005 53From this result we know the drainage basin shall be Venta River Basin.

1006 54Therefore, the predicted answer is Venta River Basin.

1007 55

1008 56

1009 57Query: What is the typical diameter (in centimetre) of tennis?

1010 58Knowledge: [

1011 59Title: What Size Is A Tennis Ball In Cm? - Metro League

1012 60Content: To Recap. A tennis ball is typically about 2 cm in diameter. Similar Posts: What Is A Junk Ball In
1013 Tennis?

1014 61How tall is a tennis ball? Tennis Balls come in different sizes, some as small as 2.575"-2.7" (6.54-6.86 cm)
1015 and others up to 8 inches (20 cm). The mass of a tennis ball must be between 1.975-2.095 oz (56-59 g).
1016 62

1017 63Title: Tennis Ball Dimensions & Drawings | Dimensions.com

1018 64Content: Tennis Balls have a diameter of 2.575"-2.7" (6.54-6.86 cm) and circumference of 8.09"-8.48"
1019 (20.6-21.5 cm). The mass of a Tennis Ball must be between 1.975-2.095 oz (56-59.4 g).

1020 65Tennis Balls have a diameter of 2.575"-2.7" (6.54-6.86 cm) and circumference of 8.09"-8.48" (20.6-21.5 cm).
1021 The mass of a Tennis Ball must be between 1.975-2.095 oz (56-59.4 g). A Tennis Ball is a ball designed
1022 for the sport of tennis.

1023 66

1024 67Title: Tennis ball - Wikipedia

1025 68Content: Modern tennis balls must conform to certain criteria for size, weight, deformation, and bounce
1026 criteria to be approved for regulation play. The International Tennis Federation (ITF) defines the
1027 official diameter as 6.54-6.86 cm (2.57-2.70 inches). Balls must have masses in the range 56.0-59.4 g
1028 (1.98-2.10 ounces).

1029 69]

1030 70Output: The query asks about diameter of tennis (in centimetre). the most relevant knowledge entry about
1031 diameter of tennis is "tennis balls have a diameter of 2.575"-2.7" (6.54-6.86 cm) and circumference of
1032 8.09"-8.48" (20.6-21.5 cm)".

1033 71As the query ask about centimetre, cm. From this result we know the diameter shall be 6.54 - 6.86.

1034 72Therefore, the predicted answer is 6.54 - 6.86.

1035 73

1036 74

1037 75

1038 76Query: Who is the inventor of women netball, sport?

1039 77Knowledge: [

1040 78Title:

1041 79Content: History of netball - Wikipedia

1042 80

1043 81In 1893, Martina Bergman-osterberg informally introduced one version of basketball to her female physical
1044 training students at the Hampstead Physical Training College in London, after having seen the game
1045 being played in the United States.

1046 82

1047 83Title: History of netball - Wikipedia

1048 84Content: In 1893, Martina Bergman-osterberg informally introduced one version of basketball to her female
1049 physical training students at the Hampstead Physical Training College in London, after having seen the
1050 game being played in the United States. Madame osterberg advocated physical fitness for women to better
1051 prepare them for motherhood and in the wider context of women's emancipation.

1052 85

1053 86Title: Netball - Wikipedia

1054 87Content: A common misunderstanding of netball's origins has resulted in the mistaken belief that netball was
1055 created to prevent women from playing basketball. However, netball's development traces back to
1056 American sports teacher Clara Gregory Baer's misinterpretation of the basketball rule book in 1895.

1057 88History. Netball's early development emerged from Clara Baer's misinterpretation of the early rules of James
1058 Naismith's new sport of basketball (which he developed while studying in Massachusetts) and eventually
1059 evol

1060 89

1061 90Title: History of Netball - World Netball

1062 91Content: Women's indoor basketball began exactly two days later when female teachers to the gym were
1063 captivated by the game but it wasn't until 1895 that the current game of netball was well and truly
1064 shaped. When Clara Baer, a sports teacher in New Orleans, wrote to Naismith asking for a copy of the
1065 rules, the subsequent rules package contained a drawing of
1066 92]

1067 93Output: The query asks about inventor of women netball. The most relevant knowledge entry about women netball
1068 inventor is "In 1893, Martina Bergman-Osterberg informally introduced one version of basketball to her
1069 female physical training students".

1070 94From the result, we know the inventor shall be Martina Bergman-Osterberg.

1071 95Therefore, the predicted answer is Martina Bergman-Osterberg.

1072 96

```

1073 97
1074 98 Query: How many elevators does torre picasso have?
1075 99 Knowledge: [
1076 100 Title:
1077 101 Content: Torre Picasso | Turismo Madrid
1078 102
1079 103 The interior of the Picasso Tower houses offices designed as intelligent spaces equipped with the highest
1080 104 technology, comfort and use of space. It has 18 lifts, divided into three groups of six.
1081 105
1082 106 Title: Torre Picasso - Wikipedia
1083 107 Content: 26 elevators; 18 serve office floors divided into three zones:
1084 108 - 1st-18th floors at 2.5 m/s (8.20 ft/s)
1085 109 - 18th-32nd floors at 4 m/s (13.12 ft/s)
1086 110 - 32nd-43rd floors at 6 m/s (19.69 ft/s) (fastest in Spain)
1087 111
1088 112 Title: Torre Picasso - Field Trip
1089 113 Content: 26 elevators, of which 18 to office floors in 3 groups of 6:
1090 114
1091 115 - 1st-18th floors at 2.5 m/s (8.20 ft/s)
1092 116 - 18th-32nd floors at 4 m/s (13.12 ft/s)
1093 117 - 32nd-43rd floors at 6 m/s (19.69 ft/s) (apparently fastest in Spain)
1094 118
1095 119 Title: Torre Picasso - Wikiwand
1096 120 Content: The building as seen from the junction of the Paseo de la Castellana and the Plaza de Pablo Ruiz
1097 121 Picasso. 26 elevators; 18 serve office floors divided into three zones: 1st-18th floors at 2.5 m/s
1098 122 (8.20 ft/s) 18th-32nd floors at 4 m/s (13.12 ft/s)
1099 123
1100 124 Output: The query asks about number of elevators in torre picasso. the most relevant knowledge entry about
1101 125 number of elevators in torre picasso is "26 elevators; 18 serve office floors divided into three
1102 126 zones:".
1103 127
1104 128 From the result, we know the number of elevators shall be 26.
1105 129 therefore, the predicted answer is 26.
1106 130
1107 131 ""

```

Listing 6: Reason Prompt (Knowledge Question)

```

1108 1 class MemoryState:
1109 2     state: str = ''
1110 3     traversed_actions: list = []
1111 4     query: str = ''
1112 5     context: str = ''
1113 6
1114 7     def __init__(self, state, query = '', context = ''):
1115 8         self.state = state
1116 9         self.query = query
1117 10        self.context = context
1118 11
1119 12 def plan(transition_graph, cur_memory, lens_res, retr_res):
1120 13     action_list = [a for a in transition_graph[cur_memory.state] if a not in cur_memory.traversed_actions]
1121 14     action_prompt = ''
1122 15     for a in action_list:
1123 16         action_prompt += ' --' + a + ': ' + task_instructions[a] + '\n'
1124 17     prompt_example = ""
1125 18     for a in action_list:
1126 19         prompt_example += action_prompt_dict[a] + "\n"
1127 20     action_prompt = planner_prompt % (cur_memory.query, action_prompt, prompt_example, cur_memory.query,
1128 21     cur_memory.context)
1129 22     action = api_utils.call_palm(action_prompt)[0]
1130 23
1131 24     instruction = []
1132 25     if action in require_instruction:
1133 26         exclude_ids = cur_memory.traversed_actions:
1134 27         prompt = instruction_prompt(cur_memory.query, lens_res, exclude_ids)
1135 28         res = api_utils.call_palm(prompt)[0]
1136 29         reason = parse_reason('the query asks about ' + reason)
1137 30         instruction = [reason, res]
1138 31     return action, instruction
1139 32
1140 33 def avis_execution(d):
1141 34     state = 'START'
1142 35     answer = None
1143 36
1144 37     prompt = question_decomposition_prompt + 'Question: ' + q + '\n'
1145 38     res = api_utils.call_palm(prompt)[0]
1146 39
1147 40     vqi = res.find('Visual: ')
1148 41     kqi = res.find('Knowledge: ')
1149 42     vq = res[vqi + 8: kqi-1]
1150 43     kq = res[kqi+11:]
1151 44
1152 45     working_memory = [MemoryState(state = 'START', query = vq, context = lens_res[0])]
1153 46     while not answer:
1154 47         cur_memory = working_memory[-1]
1155 48         action, instruction = plan(transition_graph, cur_memory, lens_res, retr_res)
1156 49         exec_res = execute(action, instruction, lens_res, retr_res)
1157 50         res = reason(exec_res)
1158 51         if 'not informative' in res:
1159 52             cur_memory.traversed_actions += [action]
1160 53         elif 'answer is' in res:

```



```
1161 53     answer = res[10:]
1162 54     else:
1163 55         working_memory += [MemoryState(state = action, query = kq, context = res)]
1164 56     return answer
```

Listing 7: Workflow of AVIS (code snippets)