# ControlLLM: Augment Language Models with Tools by Searching on Graphs

Zhaoyang Liu[*1,2]    Zeqiang Lai[*2]    Zhangwei Gao[2]    Erfei Cui[2]    Xizhou Zhu[2,3]
Lewei Lu[4]    Qifeng Chen[1 ✉]    Yu Qiao[2]    Jifeng Dai[2,3]    Wenhai Wang[2 ✉]

[1]The Hong Kong University of Science and Technology
[2]OpenGVLab, Shanghai AI Laboratory    [3]Tsinghua University    [4]SenseTime

## Abstract

*We present ControlLLM, a novel framework that enables large language models (LLMs) to utilize multi-modal tools for solving complex real-world tasks. Despite the remarkable performance of LLMs, they still struggle with tool invocation due to ambiguous user prompts, inaccurate tool selection and parameterization, and inefficient tool scheduling. To overcome these challenges, our framework comprises three key components: (1) a task decomposer that breaks down a complex task into clear subtasks with well-defined inputs and outputs; (2) a Thoughts-on-Graph (ToG) paradigm that searches the optimal solution path on a pre-built tool graph, which specifies the parameter and dependency relations among different tools; and (3) an execution engine with a rich toolbox that interprets the solution path and runs the tools efficiently on different computational devices. We evaluate our framework on diverse tasks involving image, audio, and video processing, demonstrating its superior accuracy, efficiency, and versatility compared to existing methods.*

## 1. Introduction

Large-scale language models (LLMs) like ChatGPT [18] and LLaMA series [29, 30] have demonstrated impressive capability in understanding and generating natural language. Their skill at interaction, planning, and reasoning are also rapidly extended beyond NLP, propelling the advancement of the studies in multi-modal interaction [1, 13, 14, 17, 31, 32, 40].

An emerging example of multi-modal interaction is tool-augmented language models [16, 25, 26, 36, 37], which strive to enhance the capabilities of language models beyond text to include varying modalities such as images, videos, au-

dio, *etc*. These models employ LLMs as primary controllers and incorporate tools with diverse functionalities as plugins. This helps to solve a wide range of multi-modal tasks and opens the door for innovative applications in multi-modal interaction. However, challenges persist in this field such as task decomposition, tool selection and invocation, argument and return value assignment to tools, and efficient tool execution scheduling.

To address this challenge, recent methods [16, 26, 36, 38] use ChatGPT with techniques like Chain-of-Thought (CoT) [34], Tree-of-Thought (ToT) [38], self-consistency [33] and in-context learning [7] to solve complex tasks by breaking them into a chain or tree of sub-tasks But these methods are based on the assumption that each sub-task has at most one preceding task. This presumption is inadequate for real-world applications, particularly multi-modal tasks, which usually require multiple inputs. For instance, to remove an object in an image, one needs to input the image path, the object's position, and the removal command. Therefore, a chain-shaped or tree-shaped paradigm may struggle to invoke tools in complex situations accurately and fail to manage intricate topological relationships among tools (see Fig. 1).

In this work, we introduce ControlLLM, a new framework designed to assist Large Language Models (LLMs) in accurately and efficiently controlling multi-modal tools and providing comprehensive solutions for complex real-world tasks involving multi-modal inputs. Our framework places particular emphasis on three aspects as follows:

**Task Decomposition.** A task decomposer is introduced to analyze the user prompt and breaks it down into well-defined subtasks, each with clear fields such as task description, task domain, arguments, and returned output. By breaking down complex tasks into manageable subtasks, the task decomposer significantly enhances the system's ability to handle intricate user prompts. It paves the way for follow-up task planning and solution execution.

**Task Planning.** This part handles tool selection and

---

\* Equal contribution.
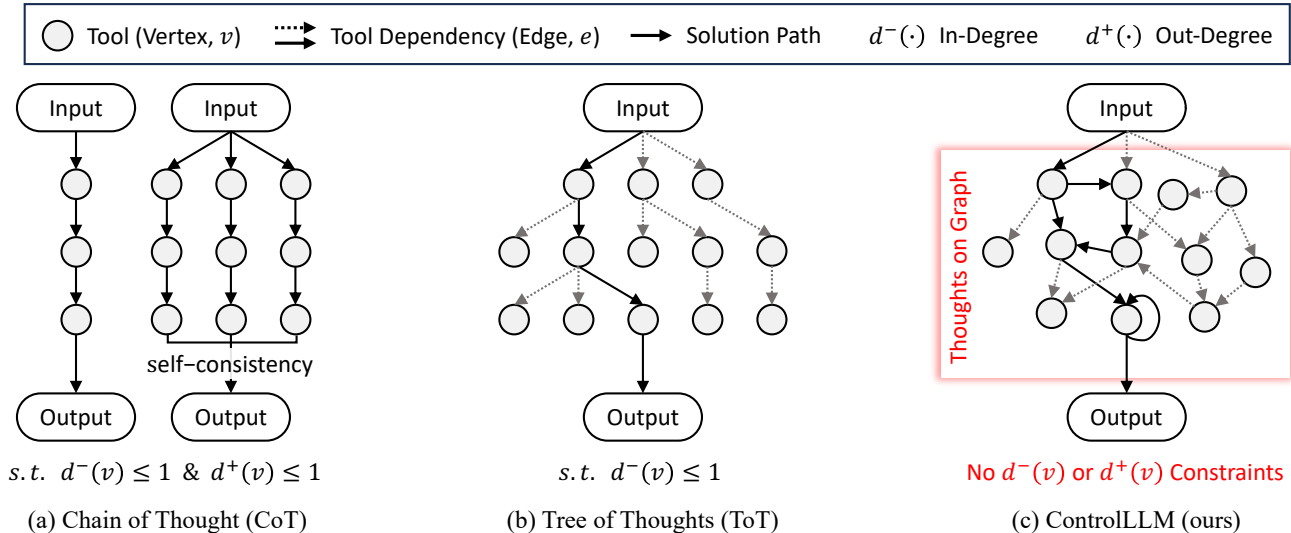✉ Corresponding authors (wangwenhai@pjlab.org.cn, cqf@ust.hk).

Figure 1. **Comparison of different paradigms for task planning.** (a) Chain of Thought (CoT) [34], CoT with self-consistency [33] and (b) Tree of Thoughts [38] (ToT) essentially rely on the LLMs to perform task planning, where the edge is actually formed by LLMs at run time. (c) The Thoughts-on-Graph (ToG) paradigm in our method searches for solutions on a pre-built graph that captures the dependencies of tools, which avoids the hallucination problem in tool invocation.

tool argument assignment. We propose a thoughts-on-graph (ToG) paradigm that traverses a topological tool graph to search for solutions. The nodes of the graph are tools that are connected based on their dependencies and relationships. ToG orchestrates the selected tools and controls the flow of resources among them to form possible solutions. ToG can find the optimal solution for each sub-task by applying diverse search strategies on the graph. Due to the concrete definition in subtask and explicit tool dependencies in a tool graph, ToG can effectively search all feasible solution paths in cases where the selected optimal solution fails to meet users' preferences.

**Solution Execution.** We design an execution engine that can execute the solution generated by ToG and craft informative and well-formatted responses. The engine has access to a versatile toolbox consisting of various tools from different sources, such as locally deployed APIs or cloud services. The engine can also parallelize the tool executions according to the topology of the solution path to reduce the latency and provide feedback during the execution process.

Our ControlLLM offers several advantages. (1) It can accurately handle complex real-world tasks that involve multimodal inputs and outputs, while previous methods [4,15,16, 26,36,37] are usually limited by the modality or complexity of the tasks; (2) It can reduce the impact of ambiguity problems in natural language. With well-defined objectives of subtask, our ToG can provide multiple alternative solutions; (3) It can overcome the token limitation of LLMs during task planning. It is because our method searches the optimal solution path on the tool graph, instead of asking LLMs to

generate solutions to meet user requirements. (4) It is easy to scale up the toolbox in our method. Since the optimal solution lies in the tool graph, when the toolbox changes, we only need to rebuild the graph without re-training LLMs or updating in-context prompts.

In summary, the main contributions are as follows:

(1) We propose ControlLLM, a framework that lets LLMs use various tools across different modalities to solve complex tasks in the real world. With a powerful toolbox, ControlLLM can be easily extended to tasks with natural language, images, audio, video, or any mix of them.

(2) We design three tailored components in ControlLLM: Task decomposition, which breaks down the user prompt into subtasks with well-defined inputs and outputs; The ToG paradigm for task planning, which searches the optimal solution path on a graph that depicts tool dependencies; And an execution engine with a powerful toolbox, which efficiently schedules and executes the solution path.

(3) We construct a benchmark to assess the efficacy of ControlLLM on tasks with different complexity levels. The experimental results demonstrate significant improvements in tool usage. Notably, ControlLLM achieves a success rate of 98% in the metric of overall solution evaluation on challenging tasks, while the best baseline only reaches 59%.

## 2. Related Work

**Planning, Reasoning, and Decision Making.** It is a longstanding vision to empower autonomous agents with the abilities of planning, reasoning, and decision-

making [12, 27, 35]. Despite progressive development, it was recent advancements in large language models (LLM) [3, 5] that have taken a breakthrough step in addressing these problems on the broad user requests. Nevertheless, it is shown that LLMs still suffer from difficulties in dealing with knowledge-heavy and complex tasks [23]. To overcome these issues, Chain of Thoughts (CoT) [34] is introduced as a simple yet effective prompting technique to elite the complex reasoning capabilities of LLMs. Following this line of work, CoT with self consistency [33], Tree of Thoughts (ToT) [34], and other techniques [6, 10, 41], have been proposed to improve the reasoning abilities further. There are also several works [2, 39] that introduce techniques called Graph-of-Thought (GoT). They all share a common insight that excessively relies on LLMs to generate thoughts for solving complicated NLP problems. In contrast, our ToG aims to endow the language model with the ability to use tools for a multi-modal dialogue system. It designs the search algorithm to form a complicated thought network of task planning. This paradigm evidently improves the performance of task planning.

**Tool-Augmented LLM.** Drawing inspiration from the evolving planning and decision-making capabilities observed in Large Language Model (LLM) systems, a new wave of research [21] starts to enhance LLMs with external tools for accessing up-to-date information, reducing hallucination, multi-modal interactions, *etc*. Prominent examples include VisProg [8], Visual ChatGPT [36], Hugging-GPT [26], InternGPT [16], AutoGPT[1], and Transformers Agent[2]. A distinctive trait of this line of research is its reliance on the zero-shot or few-shot in-context learning capabilities inherent in LLMs [3]. These capabilities enable task decomposition, tool selection, and parameter completion without requiring explicit fine-tuning. However, due to the inherent limitations of LLMs, issues such as hallucination and challenges in effective decomposition and deduction can arise with substantial frequency. In conjunction with this line of work, there are also instruction-tuning methods [9,19,20,22,25,37]. Whereas alleviating the aforementioned issues after being tuned on the conversations involved tools, these methods are still limited at expanding the toolset, i.e., additional training is required to add tools.

**Multi-Modal LLMs.** Developing LLMs that inherently possess multi-modal capabilities is another approach to extending the usage boundary of LLMs for more complex real-world scenarios. For instance, BLIP-2 [13] binds

a frozen image encoder and LLM to enable the vision-language understanding and generation. Similarly, Vision-LLM [32] and LISA [11] empower the LLMs with the visual perception capabilities such as object detection and segmentation. Nevertheless, these methods could only cover a limited range of modalities and often require huge effects on model fine-tuning.

## 3. ControlLLM

The prevalence of LLMs has unprecedentedly boosted the development of human-computer interaction, and we can empower the LLMs with abilities to interact with broader modalities via tools. In response, we present an innovative framework, namely **ControlLLM**, characterized by its flexibility, extensibility, and analytical precision. As depicted in Fig. 2, our framework consists of four sequential stages, *i.e.*, task decomposition, task planning, solution execution, and response generation. In the following, we illustrate the design of each stage in detail.

### 3.1. Task Decomposition

ControlLLM starts with task decomposition – a stage for decomposing the user request $r$ into a list of parallel subtasks. We here can utilize a language model $\mathcal{M}$, *e.g.*, Chat-GPT or instruction-tuned LLaMA, to automatically decompose the user request as follows:

$$\{s_0, ..., s_i, ..., s_n\} = \mathcal{M}(r), \tag{1}$$

where $s_i$ is the i-th subtask, $n$ is the number of all subtasks. We will elaborate on the different choices of language model $\mathcal{M}$ in Sec. 3.4 and discuss their impacts in experiments in Sec. 4.4. The result of task decomposition is JSON format, and the protocol is in Table 1.

Task decomposition is different from task planning. It only breaks down the user's request into several parallel subtasks and summarizes the input resources for each subtask from the user request. It does not need to know what tools to use or how to use them. The objective of this stage is to achieve three aims. Firstly, it splits user requests into smaller and more manageable units, *i.e.*, subtasks, thereby accelerating task planning. Secondly, it seeks to determine the search domain that is most relevant and appropriate for the given problem, thus further narrowing down the search space. Thirdly, it endeavors to infer the input and output resource types from the context of the instructions, which identifies the start and end nodes for ToG to search.

### 3.2. Task Planning with Thoughts-on-Graph

This stage constitutes the crux of the entire system. Given the results of task decomposition, we design a Thoughts-on-Graph (ToG) paradigm to heuristically find the solution on the graph.

---

[1]https://github.com/Significant-Gravitas/Auto-GPT
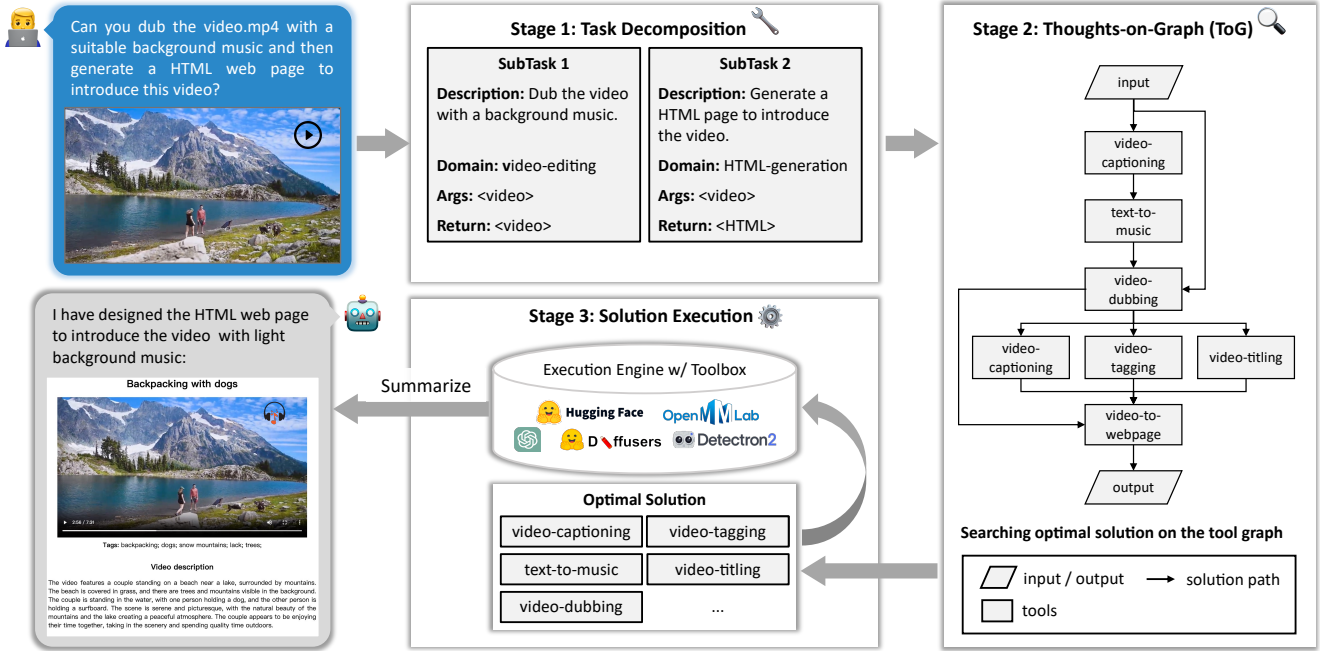[2]https://huggingface.co/docs/transformers/transformers_agents

Figure 2. **System Design of ControlLLM**. The framework consists of three stages. The first stage is task decomposition, which parses the user input into several subtasks. Then, in Stage 2, ToG utilizes a depth-first search algorithm to find the optimal solution for each subtask. The execution engine in the last stage executes the solution and returns the output to users. We here use the example of generating a web page for the video to illustrate our method.

| Field | Description |
|---|---|
| description | a brief summary of what the subtask wants to achieve. It gives some guidance on how to approach the problem for ToG. |
| task_domain | the domain that this task belongs to. It helps ToG narrow down the search space and find the most relevant and suitable tools for the subtask. |
| args | the inputs that the user provides for this subtask. It is usually in the form of key-value pairs, where the key is the type of the argument, and the value is the resource path or text you want to use. For example, [{"type": "image", "value": "image_1.png"}, {"type": "text", "value": "remove the dog in the picture"}]. |
| return | the expected output of the subtask. For example, the return is {"image": "⟨GEN⟩-0"}, which means the expected output is an image and "⟨GEN⟩-0" is just a temporary placeholder. |

Table 1. We elaborate on each field in the output protocol of task decomposition.

### 3.2.1 Building the Tool Graph

In this stage, we embark on constructing a Tool Graph $G$ by simply using an adjacency matrix, which serves as a fundamental module for analyzing and optimizing the interactions between tools and resources. Our endeavor is driven by observing a discernible topological structure that inherently exists between the input and output of diverse tools, as demonstrated in Fig. 2. This compelling insight propels us to craft a comprehensive tool graph that encapsulates this inherent relationship.

*Resource* node can be formally defined as one-tuple: ⟨type⟩, where "type" represents the specific type of re-

source. *Tool* node can be expressed as a three-tuple: ⟨desc, args, ret⟩, where each element carries significant implications for comprehending the functionalities of a tool. The desc field encapsulates the description of the tool, elucidating its purpose, methodology, and intended applications. The args field represents a list of input resource types that the tool accepts, thereby giving the prerequisites for its effective utilization. Finally, the ret field designates the resource type that the tool generates upon completion of its operations.

**Edge Definitions.** Edges in the tool graph intricately connect the nodes, highlighting the relationships between

4

different tools. We define two types of edges in the graph.

(1) *Tool-resource edge* is established from the tool to its returned resource type. This signifies that the tool is capable of generating resources of the corresponding type. Mathematically, a tool-resource edge is represented as:

$$G(T_j, R_i) = \begin{cases} \text{true}, & \text{if } R_i \text{ equals to } \texttt{ret} \text{ of } T_j, \\ \text{false}, & \text{otherwise}, \end{cases} \quad (2)$$

where $T_j$ is $j$-th tool node, $R_i$ is $i$-th resource node, true denotes two nodes are connected, and false denotes two nodes are disconnected.

### 3.2.2 Resource-Tool Edge

(2) *Resource-tool edge* represents this resource type that can be accepted as input arguments for its adjacent tool. This connection indicates how the resources flow to the tool. The resource-tool edge is mathematically defined as:

$$G(R_i, T_j) = \begin{cases} \text{true}, & \text{if } R_i \text{ belongs to } \texttt{args} \text{ of } T_j, \\ \text{false}, & \text{otherwise}. \end{cases} \quad (3)$$

Through the establishment of this graph, we can use diverse search strategies to make informed decisions regarding tool selection, input resource assignment, and workflow optimization.

### 3.2.3 Searching on the Graph

The pseudocode of the solution search algorithm is shown in Algorithm 1. Our ToG is built upon a depth-first search (*DFS*) algorithm where the tool selection function $\mathcal{F}$ is used to sample the tool nodes on the tool graph. The searched solution is a sequence of tools that take input resources as input and return the output resource to finish the user request. The algorithm starts from the input resource nodes and explores all possible paths to the output resource node while keeping track of the intermediate resources and tools along the way. The algorithm stops when it reaches the expected output node or when it exceeds a maximum length limit. In the end, the algorithm returns all the solutions it finds as a list of tool sequences. Each step from *resource node* to *tool node* represents a thought process, as it involves a decision that determines whether to use this tool and how to specify its input arguments from available resources.

To find a trade-off between time and space complexity, we develop a tool assessment module in which the language model is leveraged to score the tools in each search step and then filter out some irrelevant tools. With this assessment module, we design four search strategies for the function $\mathcal{F}$ to determine which tool nodes within "task_domain" to visit among all adjacent nodes when searching on the graph:

---

**Algorithm 1** The Python pseudocode of depth-first solution search in Thoughts-on-Graph

---

**Input:**
    t: subtask obtained by Eq. 1
    g: tool graph $G$ constructed in Sec. 3.2.1
    r: available resources, initialized with subtask["args"]
    s: recorded tools during searching
**Output:**
    solutions: all potential solutions for the subtask

1: **function** DFS_SEARCH(t, g, r, s)
2:     **if** len(s) > $m$:
3:         **return** []
    # $\mathcal{F}$ finds all tool candidates, explained in Sec. 3.2.3
4:     available_tools = $\mathcal{F}$(g, r)
5:     solutions = []
6:     **for** tool **in** available_tools:
7:         solution.append(tool)
8:         resources.append(tool["args"])
9:         **if** tool["returns"] == subtask["returns"]:
10:           results.append(solution)
11:         results = DFS_SEARCH(t, g, r, s)
12:         solutions.extend(results)
13:         resources.remove(tool["args"])
14:         solution.remove(tool)
15:     **return** solutions         ▷ Return
16: **end function**

---

**Greedy Strategy**. This strategy selects the tool node with the highest score at each step, where the score indicates the relevance of the tool to the task. A higher score indicates that the tool is more helpful for solving the task. Greedy search is fast and simple, but it may not find the optimal solution or even any solution at all.

**Beam Strategy**. This strategy limits the number of nodes selected at each level. It only keeps the $k$ best nodes according to their assessment score and discards the rest. Beam search can expand the search space but reduce the search efficiency.

**Adaptive Strategy**. This is a variant of beam search where it dynamically adjusts the beam size by choosing the tools with scores higher than a fixed threshold, which is a trade-off between exploration and exploitation. It can widen the search space when there are many available choices and narrow down when there are few confident choices.

**Exhaustive Strategy**. This strategy explores all possible paths from the start node to the goal node. An exhaustive search is guaranteed to find an optimal solution if one exists, but it may be very slow and consume a lot of memory during the search.

The impacts of different search strategies are studied in Sec. 4.4. By initiating a systematic traversal of this graph, commencing at the "args" nodes and culminating at the "re-

turn" node, a diverse list of conceivable solutions is meticulously synthesized. This diverse list, akin to a brainstorm or mind map, represents the spectrum of potential solutions.

### 3.2.4 Solution Expert

In this section, we delve into the core concept of the solution expert that streamlines the process of evaluating and selecting optimal solutions from all possible candidates. By systematically converting each solution into a formatted string description and harnessing the capabilities of prompt engineering, the solution expert enables us to make informed decisions based on evaluated scores.

**Solution Description Formatting.** To facilitate the solution expert to comprehend the solution, we need to generate the description for each solution candidate. This involves transforming raw solution data into structured, formatted string descriptions. These descriptions encapsulate the essence of each solution, highlighting its key features, including inputs, output, and docstring.

**Solution Evaluation.** In the subsequent phase, the solution expert capitalizes on prompt engineering techniques to assess each solution based on subtask descriptions and formatted solution descriptions. The designed prompts serve as a bridge, guiding language model $\mathcal{M}$ to evaluate the feasibility of each solution against the objective of the subtask. Through this process, we can assign scores to solutions, gauging their effectiveness and relevance to the task. Prompt engineering ensures that the evaluation process is focused, targeted, and aligned with the subtask. The prompt template is shown in the Table 8.

**Solution Ranking.** The final aim of this module is to select the top-performing solutions. The optimal solution is identified as the highest score assessed in the last step. Given that sometimes the selected optimal solution may not meet the user requirements, we also provide several alternative solutions by setting a threshold score of 3. These solutions, which exhibit a higher degree of alignment with the subtask's requirements, emerge as the most promising candidates for user preference.

Through collaborative efforts, the solution expert ensures that solutions are appropriately tailored, optimized, and well-adapted to the task.

### 3.2.5 Resource Expert

In the Algorithm 1, we encounter a challenge stemming from the potential presence of multiple instances of the same resource type within the available resource list. This challenge introduces complexity, making it difficult to straightforwardly deduce certain arguments for tools using predefined rules. To address this intricacy, we design a solution expert.

This module transforms the task of completing missing arguments into a fill-in-the-blank exercise. To achieve this, the resource expert crafts prompts that not only incorporate the task description but also include the available resource list. In this manner, a language model $\mathcal{M}$ is employed to dynamically complete the missing parameters within a solution by interacting with the contextual information presented. We put the prompt template in the Table 9.

### 3.3. Solution Execution

After the task solutions are generated, they are passed to a tool engine for execution, as shown in Fig. 2. During this stage, the execution engine initially parses the solutions into a sequence of *Actions*. Each action is associated with particular tool services, which could be implemented via either handcrafted mapping tables or an automatic scheduler based on some strategies. Different from previous works [16, 36, 37] that adopt static tool mapping, our design empowers the system with the flexibility to schedule diverse tools (such as different object detection models) based on users' preference for computational resources and execution duration trade-offs.

The parsed actions are executed by an interpreter that automatically dispatches the action to the local, remote, or hybrid endpoints. Similar to HuggingGPT [26], multiple independent subtasks would be executed in parallel to improve efficiency. Besides, our interpreter maintains a state memory storing all the intermediate results, including their values and types. This enables the running-time automatic correction for the action parameters.

**Response Generation** With all the execution results in hand, we could respond to the user requests. The unprocessed results may lack comprehensiveness and clarity, potentially making it difficult for user to understand. To this end, we introduce an additional stage to aggregate all the execution results and generate user-friendly responses. This is achieved by prompting the LLMs, such as ChatGPT, with the user request, action list, and execution results and asking them to summarize the answers intelligently. The prompt can be found in Table 10.

### 3.4. The Choices of Language Model

One feasible yet direct choice is to use off-the-shelf **large language models** (LLMs) such as ChatGPT and Llama 2, which are pre-trained on large-scale text corpora and can handle various NLP tasks. These LLMs are readily available. We design a series of elaborate prompts with in-context learning as shown in Appendix 6.3 for task decomposition, tool assessment, solution expert, and resource expert. The advantage of using off-the-shelf LLMs is that they have strong zero-shot capabilities and do not need to train a language model from scratch. However, the disadvantage is that they may lead to low performance as they are not

| Features | ControlLLM (our work) | HuggingGPT [26] | Visual ChatGPT [36] | InternGPT [16] | GPT4Tools [37] |
|---|---|---|---|---|---|
| Image Perception | ✓ | ✓ | ✓ | ✓ | ✓ |
| Image Editing | ✓ | ✓ | ✓ | ✓ | ✓ |
| Image Generation | ✓ | ✓ | ✓ | ✓ | ✓ |
| Video Perception | ✓ | ✓ | ✗ | ✓ | ✗ |
| Video Editing | ✓ | ✓ | ✗ | ✓ | ✗ |
| Video Generation | ✓ | ✓ | ✗ | ✓ | ✗ |
| Audio Perception | ✓ | ✗ | ✗ | ✗ | ✗ |
| Audio Generation | ✓ | ✓ | ✗ | ✗ | ✗ |
| Multi-Solution | ✓ | ✗ | ✗ | ✗ | ✗ |
| Pointing Inputs | ✓ | ✗ | ✗ | ✓ | ✗ |
| Resource Type Awareness | ✓ | ✗ | ✗ | ✗ | ✗ |

Table 2. **A comparison of features between different methods.** The table shows that our framework supports more features that facilitate the user experience of multi-modal interaction. It proves the high scalability of our framework.

trained for our requirements.

Another alternative choice is to finetune a language model, *e.g.*, LLaMA [29], by using self-instruct method [33]. More details of optimizing $\mathcal{M}$ can be referred to Appendix 6.1. The advantage of finetuning a language model is that it can achieve high performance by adapting to the data and the task. However, the disadvantage is that it requires computational resources to train the model and may suffer from overfitting or data scarcity issues.

Regarding this issue, it is essential to carefully consider the trade-offs between readily available off-the-shelf LLMs with zero-shot capabilities and the potential for finetuning a model to achieve superior performance at the cost of computational resources. We will thus further discuss the impacts of different language models $\mathcal{M}$ in 4.4 and explore the optimal settings for our framework.

## 4. Experiments

### 4.1. Benchmark

In this section, we build a benchmark that is used to evaluate our proposed framework compared with other state-of-the-art methods. In order to make fair comparisons, we only evaluate and test on the intersection of toolsets from different methods [16, 26, 36, 37]. As a result, the benchmark consists of a set of tasks that require various tools to solve complex problems. It is designed to cover different task domains, such as question answering, image generation, image editing, image perception, visual question answering, *etc*. In this benchmark, the tasks involve more than 20 tools across different modalities.

This benchmark includes more than 100 instructions that specify the user's goals and preferences for each task. The instructions are classified into three levels of difficulty: easy, medium, and hard. The difficulty level reflects the complexity and specificity of the instructions, as well as the

number and diversity of tools required to complete the task. For example, an easy instruction may only require one or two simple tools, while a hard instruction may require several complex tools with fine-grained control. We believe that this benchmark can provide a comprehensive comparison of the tool control capabilities of different methods.

### 4.2. Evaluation Protocol

Effectively evaluating the performance of tool-augmented Language Models (LLMs) remains a challenging task. This challenge stems from several factors, including the inherent ambiguities in defining standard answers, the absence of shared benchmarks, and formatted solutions for systematically assessing different methods. Consequently, existing evaluation methods either provide limited case studies [36], quantitative results for specific, deterministic problems [20, 25], or comparisons against various LLM backbones [22, 24, 26].

To address this, we compare our method against similar models such as VisualChatGPT [36], HuggingGPT [26], GPT4Tools [37], and InternGPT [16], all of which share comparable toolsets. Since the APIs of tools in different methods are slightly inconsistent, it is hard to annotate all feasible solutions for each method. As such, we adopt a crowd-sourcing evaluation approach with three annotation experts via a straightforward protocol. This protocol breaks down the evaluation into three main aspects: tool selection, argument inference, and overall assessment. Please note that the evaluation protocol is independent of the tools' capabilities. That is to say, when the tools and their input arguments are correct, we do not account for the case where the output fails to satisfy the user's expectations due to the limitations of tools.

### 4.2.1 Metrics for Tool Selection

A) Irrelevant Tool Inclusion Rate (*abbr. $IR$*): This metric gauges the performance of method in excluding irrelevant tools. It measures the proportion of the predicted solutions that contain the irrelevant tools. A higher $IR$ indicates that the method tends to include more unnecessary tools, potentially hindering effective task planning.

B) Necessary Tool Inclusion Rate (*abbr. $NR$*): This metric assesses the inclusion of necessary tools in the predicted solution $S^p$ but without considering whether the arguments of tools are correct. If $NR$ is high, it indicates the method has strong capabilities in tool selection.

### 4.2.2 Metrics for Arguments Inference

A) Resource Hallucination Rate (*abbr. $HR$*): This indicator reveals the extent of hallucination in the method's responses when inferring the arguments for tools. It measures whether all arguments of the tools used in the predicted solution exist physically. A lower $HR$ suggests that the method is less prone to generating hallucinated content.

B) Resource Type Consistency Rate (*abbr. $CR$*): This metric examines whether the types of resources used as inputs in the predicted solution match those of the corresponding tools. It evaluates the method's ability to ensure consistency between argument types and tools.

### 4.2.3 Solution Evaluation

The Solution Evaluation (*abbr. $SE$*) measures the success rate of all generated solutions on our benchmark. It only considers whether the output solution can effectively address the user's problem, irrespective of whether it contains irrelevant tools. It focuses on whether the tool call chain can resolve the user's request. A higher score in the solution evaluation indicates that the method is able to provide an effective solution to user requests.

In summary, these intuitive metrics together provide a comprehensive assessment of tool-augmented LLMs in terms of tool selection, argument inference, and overall effectiveness in addressing user queries. The formal definition of these metrics can refer to Appendix 6.2

### 4.3. Comparisons with the stage-of-the-art methods

### 4.3.1 Feature Comparison

Table 2 presents a comprehensive feature comparison among various methods [16, 26, 36, 37], highlighting ControlLLM's distinct advantages in the landscape of multi-modal interaction. Notably, "Multi-Solution" signifies the method's ability to provide multiple feasible solutions, granting users more options. "Pointing Inputs" signifies support for pointing devices during interaction, enhancing precision. "Resource Type Awareness" indicates the method's capability to discern the type of resource used in interaction, ensuring more context-aware responses. In summary, ControlLLM emerges as the standout choice, excelling in various features. It offers a comprehensive set of tools in domains of image, video, and audio. Moreover, its support for resource type awareness, multiple solutions, and pointing inputs demonstrates its adaptability and scalability, making it the highly versatile framework for diverse multi-modal interaction scenarios.

### 4.3.2 Quantitative Comparison

In this section, we provide a comprehensive analysis for ControlLLM to compare with state-of-the-art methods, as summarized in Table 3. We here provide three implementations for our method: a) ControlLLM-ChatGPT leverages the ChatGPT-3.5 as language model $\mathcal{M}$; b) ControlLLM-LLaMA that finetunes a LLaMA-7B as a language model $\mathcal{M}$; c) ControlLLM-Mix is regarded as our default setting, which finetunes LLaMA-7B as a task decomposer in stage 1 while the rest modules employ the ChatGPT to finish the tasks. ControlLLM-Mix combines the advantages of the other two variants.

Our evaluation is based on a set of metrics assessing tool selection and argument inference, as well as the overall effectiveness of the solutions. ControlLLM excels in several key aspects. Notably, it achieves the lowest Irrelevant Tool Inclusion Rate ($IR$) at a mere 0.03, indicating its exceptional ability to exclude irrelevant tools during task planning. This is a significant advantage, as minimizing irrelevant tools is crucial for efficient task execution. ControlLLM also boasts the highest Necessary Tool Inclusion Rate ($NR$) at 0.93, signifying its proficiency in selecting essential tools. Furthermore, ControlLLM demonstrates superior performance in argument inference, with the lowest Argument Hallucination Rate ($HR$) of 0.02 and the highest Argument Type Consistency Rate ($CR$) of 0.98. These results underscore its ability to generate accurate and consistent arguments, addressing a common challenge in language models. In the metric of solution evaluation, ControlLLM maintains its lead with a score of 0.94, indicating its effectiveness in resolving user requests. In summary, ControlLLM exhibits remarkable performance in tool selection, argument inference, and overall solution effectiveness, evidently outperforming the state-of-the-art methods in this field.

### 4.4. Ablation Studies

In this section, we delve into ablation studies to gain deeper insights for the performance of our method.

### 4.4.1 Ablation Studies on Different LLMs

In this section, we conduct ablation studies to evaluate the impact of different LLMs on task planning. We also in-

| Methods | Tool | | Argument | | Solution Evaluation ↑ | | | |
|---|---|---|---|---|---|---|---|---|
| | $IR\downarrow$ | $NR\uparrow$ | $HR\downarrow$ | $CR\uparrow$ | All | Easy | Medium | Hard |
| HuggingGPT [26] | 0.45 | 0.64 | 0.16 | 0.69 | 0.59 | 0.73 | 0.50 | 0.33 |
| Visual ChatGPT [36] | 0.26 | 0.58 | 0.09 | 0.76 | 0.57 | 0.73 | 0.63 | 0.10 |
| InternGPT [16] | 0.12 | 0.51 | 0.49 | 0.43 | 0.44 | 0.60 | 0.46 | 0.00 |
| GPT4Tools [37] | 0.19 | 0.44 | 0.28 | 0.72 | 0.43 | 0.64 | 0.33 | 0.00 |
| ControlLLM-ChatGPT | 0.16 | 0.63 | 0.83 | 0.83 | 0.64 | 0.71 | 0.67 | 0.43 |
| ControlLLM-LLaMA | 0.06 | 0.95 | 0.99 | 0.99 | 0.91 | 0.99 | 0.86 | 0.76 |
| ControlLLM-Mix* | **0.03** | **0.93** | **0.02** | **0.98** | **0.94** | **0.99** | **0.96** | **0.81** |

Table 3. **Compared with other state-of-the-art methods**. ↓ means the smaller the better, ↑ means the larger the better. The results of state-of-the-art methods [16, 26, 36, 37] are reproduced on our own benchmark. ∗ denotes the default setting of ControlLLM if not stated.

| Task Decomp. | LLMs | Tool | | Argument | | Solution Evaluation ↑ | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | $IT\downarrow$ | $NR\uparrow$ | $HR\downarrow$ | $CR\uparrow$ | All | Easy | Meduim | Hard |
| *w/o* priors | Llama2-13B | 0.28 | 0.71 | 0.01 | 0.99 | 0.68 | 0.87 | 0.50 | 0.38 |
| | ChatGPT-3.5 | 0.13 | 0.84 | 0.01 | 0.99 | 0.83 | 0.99 | 0.67 | 0.57 |
| | ChatGPT-4 | 0.06 | 0.91 | 0.03 | 0.97 | 0.91 | 0.98 | 0.83 | 0.81 |
| *w/* priors | Llama2-13B | 0.12 | 0.83 | 0.04 | 0.95 | 0.82 | 0.95 | 0.71 | 0.62 |
| | ChatGPT-3.5 | 0.03 | 0.93 | 0.02 | 0.98 | 0.94 | 0.98 | 0.96 | 0.81 |
| | ChatGPT-4 | **0.01** | **0.98** | **0.02** | **0.98** | **0.98** | **1.00** | **1.00** | **0.91** |

Table 4. **The effects of task decomposition with regard to different LLMs**. We find, if adding prior knowledge, such as which tools might be used, into the subtasks description in task decomposition, the performance of task planning can be evidently improved.

vestigate the effects of incorporating prior knowledge into the subtask descriptions during task decomposition. The method without prior knowledge usually directly uses the user's request as a subtask description and does not offer any hints or suggestions on tool selections in the subtask description. In contrast, In the method with prior knowledge, we managed to add prior knowledge into the subtask description, which is expected to guide the tool assessment. The results are presented in Table 4. The prior knowledge indeed improves the necessary tool inclusion rate ($NR$) and reduces the chance of selecting irrelevant tools ($IT$) when using the same large language model. Furthermore, we found the ability of language models plays a decisive role in tool selection. The more powerful the language model, the higher the score of solution evaluation.

### 4.4.2 Impact of Search Strategies

Table 5 investigates the impact of different search strategies within our Thoughts-on-Graph. We observe that the exhaustive search strategy outperforms the others in most metrics, but this strategy is time-consuming. On the other hand, the greedy search strategy achieves the lowest performance. This is because sometimes it can not search for a feasible path based on the tool with a high score due to the inaccurate tool assessment. It is thus efficient but usually fails to find the solution, especially in hard cases. In addition, the adaptive search strategy strikes a balance between performance metrics, offering competitive results in most aspects. To trade-off between time and accuracy, we thus choose the adaptive strategy as our default search strategy.

### 4.5. Qualitative analyses

In this section, we present extensive case studies to qualitatively assess our method.

Fig. 3 shows two simple cases to illustrate the capabilities of our ControlLLM in task planning. In contrast to HuggingGPT [26], we found our method is able to generate more diverse solutions to meet users' expectations, thanks to our proposed Thoughts-on-Graph paradigm.

Then, we provide more cases across different modalities to validate the user experience for our method in practice. In Fig. 4, we show some cases of image perception, which involves analyzing and understanding the content of an image, such as detecting objects, counting objects, finding objects, segmenting objects, answering questions about the image, *etc*. These tasks require the system to invoke tools to process visual information and extract relevant fea-

| Search | Tool | | Argument | | Solution Evaluation ↑ | | | |
|---|---|---|---|---|---|---|---|---|
| Strategies | $IT\downarrow$ | $NR\uparrow$ | $HR\downarrow$ | $CR\uparrow$ | All | Easy | Meduim | Hard |
| Greedy | 0.19 | 0.49 | 0.24 | 0.76 | 0.49 | 0.56 | 0.58 | 0.19 |
| Beam | 0.14 | 0.88 | **0.00** | 1.00 | 0.88 | 0.96 | 0.79 | 0.76 |
| Adaptive | **0.03** | 0.93 | 0.02 | 0.98 | 0.94 | 0.98 | 0.96 | 0.81 |
| Exhaustive | 0.06 | **0.97** | 0.01 | **0.99** | **0.97** | **1.00** | 0.96 | **0.91** |

Table 5. **The evaluation for different search strategies**. As introduced in Sec. 3.2.3, although exhaustive strategy achieves the best performance on most metrics, the adaptive strategy strikes a good balance between efficiency and effectiveness.



Figure 3. **Comparisons of task planning**. We here use two simple cases to illustrate the differences between two different methods in task planning. Here, each output node is generated by different solution paths.

tures and labels from the image. Fig. 5 gives examples of image processing and image editing, which assist users in processing or editing the image according to some criteria or instructions. Fig. 6 mainly focuses on image question answering and image generation, showing the graphic dialogue capability. In Fig. 7, we provide some multi-modal interaction cases on image, video, and audio domains. In addition, we also illustrate the capabilities of complicated scenarios with solutions searched by ToG during task planning in Fig. 8 and Fig. 9. These complex tasks involve combining multiple tools to find a more advanced and creative solution path that can solve more challenging problems. It requires a system that can integrate different types of information and outputs from tools and generate comprehensive and meaningful responses based on execution results. These figures demonstrate the strong capabilities of ControlLLM in task planning for both simple and complicated scenarios. It thus leads to a better user experience.

Figure 4. **Case study on different task domains** (a).

**image-processing**

**User:** Extract the edge information for this image.

**AI**: I have successfully extracted the edge information from the image you provided. To achieve this, I used an image processing tool called "*image_to_edge*". I passed your input image as input to the tool. The tool processed the image and generated a new image file containing only the edge information.
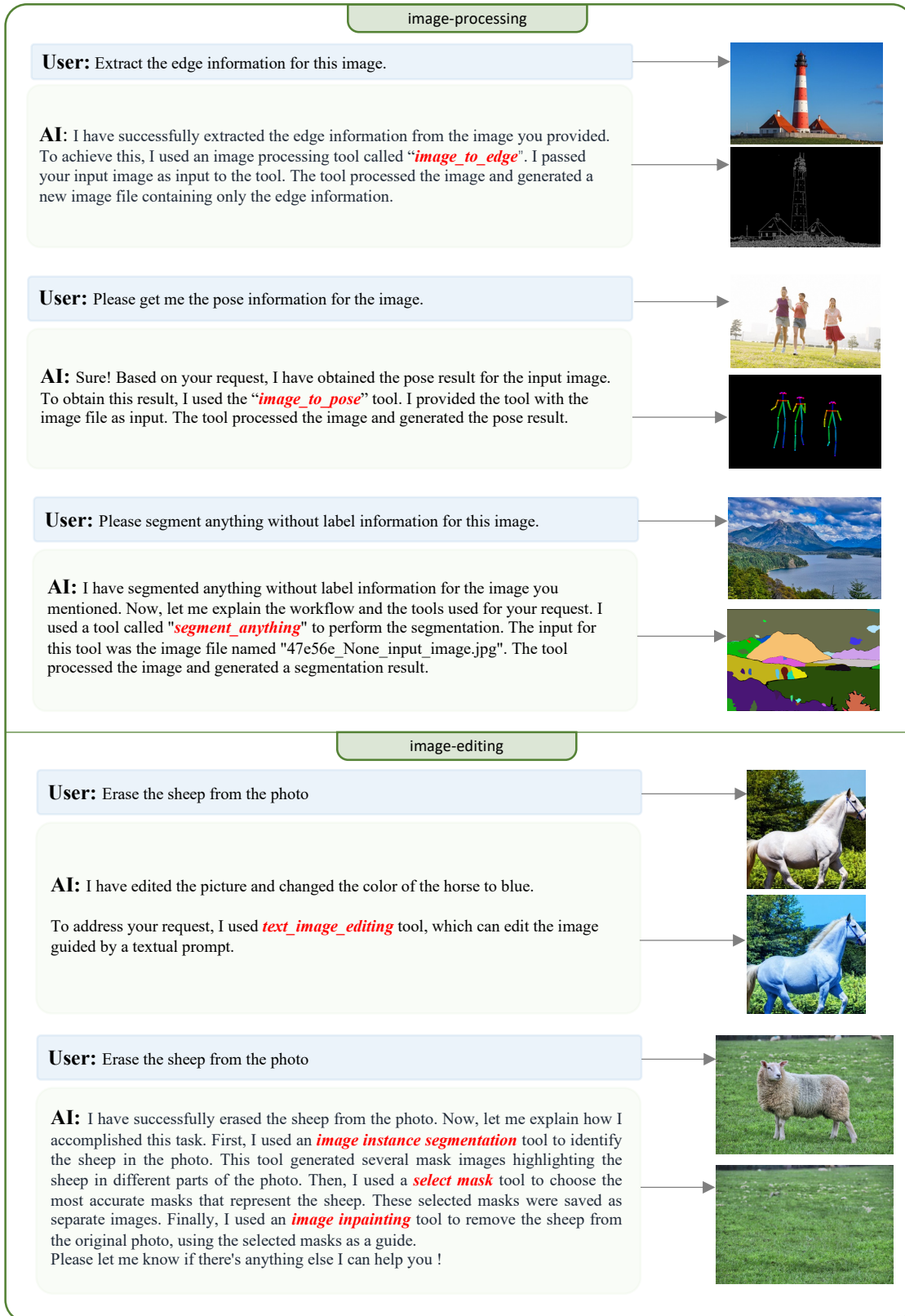
**User:** Please get me the pose information for the image.

**AI:** Sure! Based on your request, I have obtained the pose result for the input image. To obtain this result, I used the "*image_to_pose*" tool. I provided the tool with the image file as input. The tool processed the image and generated the pose result.

**User:** Please segment anything without label information for this image.

**AI:** I have segmented anything without label information for the image you mentioned. Now, let me explain the workflow and the tools used for your request. I used a tool called "*segment_anything*" to perform the segmentation. The input for this tool was the image file named "47e56e_None_input_image.jpg". The tool processed the image and generated a segmentation result.

**image-editing**

**User:** Erase the sheep from the photo

**AI:** I have edited the picture and changed the color of the horse to blue.

To address your request, I used *text_image_editing* tool, which can edit the image guided by a textual prompt.

**User:** Erase the sheep from the photo

**AI:** I have successfully erased the sheep from the photo. Now, let me explain how I accomplished this task. First, I used an *image instance segmentation* tool to identify the sheep in the photo. This tool generated several mask images highlighting the sheep in different parts of the photo. Then, I used a *select mask* tool to choose the most accurate masks that represent the sheep. These selected masks were saved as separate images. Finally, I used an *image inpainting* tool to remove the sheep from the original photo, using the selected masks as a guide.
Please let me know if there's anything else I can help you !

Figure 5. **Case study on different task domains** (b).

**image-question-answering**

**User:** Why is this image so funny?

**AI**: The image is funny because it features two toy frogs, one of which is holding a camera and taking a picture of the other. The scene is set in a grassy field, and the toy frogs are posed in a way that mimics human behavior, specifically a photographer capturing a subject. The humor comes from the unexpected and anthropomorphic portrayal of the toy frogs, as it is not common to see inanimate objects, such as toys, engaging in activities like photography. This playful and whimsical scene can evoke amusement in the viewer.
I used an *image question answering* tool. The tool analyzed the image and determined why this image is so funny.

**User:** Given the image, is there going to be a potentially dangerous situation in the picture?

**AI:** There is a potentially dangerous situation as a man is riding on the back of an elephant, while another elephant is walking through the water. The presence of two elephants in close proximity to each other and the man could pose a risk of accidents or injuries, especially if the elephants become agitated or frightened. It is essential for the man to maintain control over the elephant he is riding and ensure that the other elephant is calm and well-behaved to minimize any potential dangers. To solve your question, I used an *image question answering* tool.

**User:** how many people can be seen riding the elephant in the image?

**AI:** In the image, there is one person riding on the back of the elephant. I used an *image question answering* tool to answer your question.

**image-generation**

**User:** I want a new image that has mountains and sunset. Can you use the depth result from input image to create it?

**AI**: The new image is generated. The workflow involved two tools: "*image_to_depth*" and "*depth_text_to_image*".
The first tool, "*image_to_depth*", used an input image to generate a depth result.
The second tool, "*depth_text_to_image*", took a text input of "The new image should have mountains and sunset." and a depth map (generated by the previous tool)

**User:** please generate a new image based the edge from this image. The new image shows a beautiful tower in the rain.

**AI:** The new image generated.
Workflow:
The image provided was processed using an *image-to-edge* tool, which generated an edge representation of the image.
The edge representation was then used along with the text "The new image shows a beautiful tower in the rain" in an *edge-text-to-image* tool. This tool combined the edge representation and the text to generate a new image.

Figure 6. **Case study on different task domains** (c).

**image & video & audio tasks**

**User:** Convert the following sentence into speech: "Hope is the thing with feathers That perches in the soul, And sings the tune without the words, And never stops at all"

**AI**: I have converted the following sentence into speech for you: "Hope is the thing with feathers That perches in the soul, And sings the tune without the words, And never stops at all". To do this, I used a *text_to_speech* tool that can generate speech from text. I gave it the input sentence and it returned an audio file with a human-like voice reading the sentence. The tool can also adjust the speed, pitch, and accent of the voice to suit different preferences. I hope you like the results.



**User:** Generate an image where a woman is surfing. Then convert this image to a dynamic video.

**AI**: I have generated an image where a woman is surfing and converted it to a dynamic video for you.
To do this, I used two tools that can create and manipulate visual content. The first tool called *text_to_image* can generate an image based on a text description. I gave it the input 'a woman is surfing' and it returned an image of a woman surfing on a blue ocean wave. The second tool is *image_to_video* that can convert an image to a video by adding motion and sound effects. I gave it the image from the first tool, and it returned a video of the woman surfing with realistic water and wind sounds. I hope you like the results.



**User:** What happened in this video?

**AI**: In this video, a basketball game is taking place, with a player attempting to shoot the ball. The player is holding the ball and getting ready to make a shot, while other players are on the court, possibly trying to defend or score.
To solve your question, I used two tools: "*action_recognition*" and "*video_captioning*", which both take the video as input.



**User:** Please generate music according to the following prompt: An 80s driving pop song with heavy drums and synth pads in the background

**AI:** Sure! I have generated a piece of music based on your prompt.

To create this music, I used a *text-to-music* tool. I provided the tool with the prompt you gave me, which was "An 80s driving pop song with heavy drums and synth pads in the background". The tool then generated a music composition that matches the description.



**User:** Can you help me dub this video with uploaded audio?

**AI:** Sure! Based on your request, I have dubbed the video with the audio for you to dub a video. I used a "*dub_video*" tool that can add or replace the audio track of a video. I gave it the input video file and the input audio file, and it returned a new video file with the audio synchronized with the video.
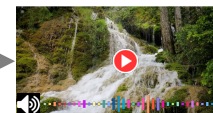


Figure 7. **Case study on different task domains** (d).

**User:** Please dub the video.mp4 with a suitable background music and then generate a HTML web page to introduce this video.

Task planning

[Action(tool_name='*video_captioning*', inputs={'video': 'video.mp4'}, outputs=['<TOOL-GEN>-text-0']),
Action(tool_name='*text_to_music*', inputs={'text': '<TOOL-GEN>-text-0'}, outputs=['<TOOL-GEN>-text_to_music-audio-0']), Action(tool_name='*dub_video*', inputs={'video': "video.mp4 ", 'audio': "<TOOL-GEN>-text_to_music-audio-0",}, outputs=['<TOOL-GEN>-dub_video-video-0'])]
[Action(tool_name='*title_generation*', inputs={'text': '<TOOL-GEN>-text-0'}, outputs=['<TOOL-GEN>-text-1']),
Action(tool_name='*text_to_tags*', inputs={'text': '<TOOL-GEN>-text-0'}, outputs=['<TOOL-GEN>-tags-0']),
Action(tool_name='*video_to_webpage*', inputs={'video': '<TOOL-GEN>-dub_video-video-0', 'title': '<TOOL-GEN>-text-1',
'tags': '<TOOL-GEN>-tags-0', 'description': '<TOOL-GEN>-text-0'}, outputs=['<GEN>-0'])]

**AI:** Sure, I can generate a web page for your input video.
I have executed a workflow using several tools to generate the desired web page. The *video_captioning* tool was used to generate a text description for the input video. The *dub_video* tool combined the input video with the generated audio to create a dubbed video.
The *title_generation* tool was used to generate a title for the video based on the generated text.
The *text_to_tags* tool generated a set of tags related to the video content.
Finally, the *video_to_webpage* tool was used to generate a web page for the video, including the video itself, the generated title, tags, and description.
The resulting web page includes the following elements:
1) A title prompt asking for a title to be generated within 10 words. 2) A video player to play the input video. 3) Tags related to the video content. 4) A video description.



**User:** Generate a new image based on prompt. Prompt: "A serene and beautiful landscape with a calm lake reflecting the blue sky and white clouds. There are some green hills and trees on the shore". Then help me generate a promotional video for this generated image.

Task planning

[Action(tool_name='*text_to_image*', inputs={'text': 'A serene and beautiful landscape with a calm lake reflecting the blue sky and white clouds. There are some green hills and trees on the shore, and a small wooden boat floating on the water.'}, outputs=['<TOOL-GEN>-text_to_image-image-0']), Action(tool_name='*image_captioning*', inputs={'image': '<TOOL-GEN>-text_to_image-image-0'}, outputs=['<TOOL-GEN>-text-0']), Action(tool_name='*text_to_speech*', inputs={'text': '<TOOL-GEN>-text-0'}, outputs=['<TOOL-GEN>-text_to_speech-audio-0']),
Action(tool_name='*image_audio_to_video*', inputs={'image': "<TOOL-GEN>-text_to_image-image-0", 'audio': "<TOOL-GEN>-text_to_speech-audio-0"}, outputs=['<GEN>-0'])]

**AI:** Sure! I can help you generate a video to introduce an image. Based on the information provided, I have used a combination of tools to fulfill your request.
First, I used the "*text_to_image*" tool to convert your description into an image. The tool created an image of a serene and beautiful landscape with a calm lake reflecting the blue sky and white clouds. There are also green hills and trees on the shore, as well as a small wooden boat floating on the water.
Next, I used the "*image_captioning*" tool to generate a textual description of the image. The tool analyzed the image and produced a caption that best describes it.
Then, I used the "*text_to_speech*" tool to convert the generated textual description into an audio file. This tool transformed the text into speech, allowing for a more engaging experience.
Finally, I used the "*image_audio_to_video*" tool to combine the generated image and audio into a video. The tool merged the image and audio together, creating a video that introduces the image.
The result of this workflow is a video file. This video showcases the image described in your initial request and includes the audio narration of its description.
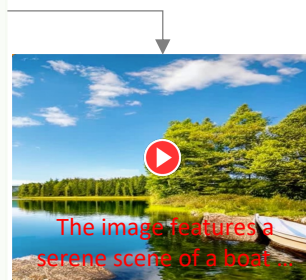


Figure 8. **Case study on different task domains** (e).

Figure 9. **Case study on different task domains** (f).

## 5. Conclusion

In this paper, we propose **ControlLLM**, a multi-modal interaction framework that can accurately control tool usage across various domains, including text, image, audio, video, *etc*. The proposed framework consists of three key stages: (1) *task decomposition* to concrete the objective of the task, (2) a *Thoughts-on-Graph* (ToG) paradigm to search the optimal solution path on the constructed tool graph, (3) and an *execution engine* with a versatile toolbox to execute solution efficiently. We conduct extensive experiments and demonstrate that our ControlLLM achieves superior performance regarding tool selection, argument inference, and overall solution effectiveness compared to existing methods.

Nevertheless, this work still has some limitations. Since the goal of this work is to improve the accuracy of tool usage, even if the solution is feasible, we cannot guarantee that the output of the tools will meet the human's expectations. On the other hand, due to the inherent ambiguity in natural language, it is difficult to ensure that the optimal solution selected is consistent with the user's goal. In this case, we can only provide more alternative solutions searched by ToG for users to choose from if the optimal solution fails.

## References

[1] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022. 1

[2] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michał Podstawski, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. Graph of Thoughts: Solving Elaborate Problems with Large Language Models, 2023. 3

[3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 3

[4] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. 2

[5] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022. 3

[6] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022. 3

[7] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*, 2022. 1

[8] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14953–14962, 2023. 3

[9] Shibo Hao, Tianyang Liu, Zhen Wang, and Zhiting Hu. Toolkengpt: Augmenting frozen language models with massive tools via tool embeddings. *arXiv preprint arXiv:2305.11554*, 2023. 3

[10] Namgyu Ho, Laura Schmid, and Se-Young Yun. Large language models are reasoning teachers. *arXiv preprint arXiv:2212.10071*, 2022. 3

[11] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. *arXiv preprint arXiv:2308.00692*, 2023. 3

[12] Jean-Claude Latombe. *Robot motion planning*, volume 124. Springer Science & Business Media, 2012. 3

[13] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. 1, 3

[14] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022. 1

[15] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 2

[16] Zhaoyang Liu, Yinan He, Wenhai Wang, Weiyun Wang, Yi Wang, Shoufa Chen, Qinglong Zhang, Zeqiang Lai, Yang Yang, Qingyun Li, Jiashuo Yu, et al. Interngpt: Solving vision-centric tasks by interacting with chatbots beyond language. *arXiv preprint arXiv:2305.05662*, 2023. 1, 2, 3, 6, 7, 8, 9

[17] Lei Ma, Jincong Han, Zhaoxin Wang, and Dian Zhang. Cephgpt-4: An interactive multimodal cephalometric measurement and diagnostic system with visual large language model. *arXiv preprint arXiv:2307.07518*, 2023. 1

[18] OpenAI. Chatgpt (Mar 14 version) [large language model]. 6, 2023. 1

[19] Aaron Parisi, Yao Zhao, and Noah Fiedel. Talm: Tool augmented language models. *arXiv preprint arXiv:2205.12255*, 2022. 3

[20] Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*, 2023. 3, 7

[21] Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, et al. Tool learning with foundation models. *arXiv preprint arXiv:2304.08354*, 2023. 3

[22] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian,

et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*, 2023. 3, 7

[23] Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021. 3

[24] Jingqing Ruan, Yihong Chen, Bin Zhang, Zhiwei Xu, Tianpeng Bao, Guoqing Du, Shiwei Shi, Hangyu Mao, Xingyu Zeng, and Rui Zhao. Tptu: Task planning and tool usage of large language model-based ai agents. *arXiv preprint arXiv:2308.03427*, 2023. 7

[25] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023. 1, 3, 7

[26] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface. *arXiv preprint arXiv:2303.17580*, 2023. 1, 2, 3, 6, 7, 8, 9, 20, 22

[27] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017. 3

[28] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023. 19

[29] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 1, 7, 19

[30] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 1

[31] Sai Vemprala, Rogerio Bonatti, Arthur Bucker, and Ashish Kapoor. Chatgpt for robotics: Design principles and model abilities. *Microsoft Auton. Syst. Robot. Res*, 2:20, 2023. 1

[32] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *arXiv preprint arXiv:2305.11175*, 2023. 1, 3

[33] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022. 1, 2, 3, 7

[34] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022. 1, 2, 3

[35] Sholom M Weiss, Casimir A Kulikowski, Saul Amarel, and Aran Safir. A model-based method for computer-aided medical decision-making. *Artificial intelligence*, 11(1-2):145–172, 1978. 3

[36] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023. 1, 2, 3, 6, 7, 8, 9

[37] Rui Yang, Lin Song, Yanwei Li, Sijie Zhao, Yixiao Ge, Xiu Li, and Ying Shan. Gpt4tools: Teaching large language model to use tools via self-instruction, 2023. 1, 2, 3, 6, 7, 8, 9

[38] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*, 2023. 1, 2

[39] Yao Yao, Zuchao Li, and Hai Zhao. Beyond chain-of-thought, effective graph-of-thought reasoning in large language models. *arXiv preprint arXiv:2305.16582*, 2023. 3

[40] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023. 1

[41] Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*, 2022. 3

# 6. Appendix

## 6.1. Training the language model

### 6.1.1 Instruction generation.

The first step to train $\mathcal{M}$ is to construct the instruction corpus. We here opt for ChatGPT-3.5 to generate the training corpus. The following steps will elaborate on the details of instructions generation for task decomposition, tool assessment, solution expert, and resource expert, respectively.

For task decomposition, we generate two different types of instructions as follows: 1) Basic instructions, where they only contain one subtask after task decomposition. We set some seed instructions with ground-truth results of task decomposition, which serve as initial templates for generating more diverse instructions. Then, we use ChatGPT to generate more diverse instructions based on the pre-defined seed instructions. During the generation process, we center on the seed instructions and produce more instructions using more diverse expressions and styles. These instructions

need to share the task decomposition results with the seed instructions as ground truth. 2) Compound instructions, which involve multiple tools and intermediate resources. We devise a method to assemble the basic instructions into the compound instructions in a coherent and logical manner. This method aims to enhance the model's modeling ability and improve the system's complex interaction capability, by enabling the model to handle user requests that span multiple domains and require multiple steps of processing. We here generate almost 100k instructions for training. The instructions generated in this step will be used in the following task as well.

For the tool assessment, solution expert, and resource expert, we use prompts in 6.3 to collect the output from ChatGPT by running GoT on the instructions generated in task decomposition. Unlike directly generating the solution, these tasks only involve making a decision, like scoring the tools or solutions based on the input, so they are relatively simple, and ChatGPT with strong zero-shot capabilities can easily solve them. Therefore, we decide to directly distill the knowledge of ChatGPT by using prompt techniques. Through the experiments, we verify the feasibility of this strategy.

### 6.1.2  Training Recipes

We follow the training protocol in [28][3], where LLaMA [29] is used as our default language model $\mathcal{M}$. The language model $\mathcal{M}$ is finetune for 3 epochs with the initial learning rate 2e-5 and consine decay. We fixed training batch size as 128 by adaptively setting "gradient_accumulation_steps". The whole training is on 8xA100 GPUs.

### 6.2. Metric Definitions

A) Irrelevant Tool Inclusion Rate (*abbr.* $IR$):

$$F(s^p) = \begin{cases} \text{true}, & \text{The solution } s^p \text{ contains the irrelevant tools} \\ \text{false}, & \text{otherwise} \end{cases}$$
(4)

$$IR = \frac{\sum_i^{||S^p||} \mathbb{I}(F(S_i^p))}{||S^p||}$$
(5)

Where $\mathbb{I}$ is indicator function, $S^p$ is predicted solution. This formula measures the proportion of tools that are irrelevant or unnecessary for solving the subtask. It counts the number of solutions that contains the useless tools. A high $IR$ value means that the method produces more useless tools for solving the problem.

B) Necessary Tool Inclusion Rate (*abbr.* $NR$):

$$H(s^p) = \begin{cases} \text{true}, & \text{Solution } s^p \text{ contains necessary tools} \\ \text{false}, & \text{otherwise} \end{cases},$$
(6)

$$NR = \frac{\sum_i^{||S^p||} \mathbb{I}(H(S_i^p))}{||S^p||}.$$
(7)

The necessary tools play critical role in solving the user request. For example, if you want to know the position of specific object, the object detection tool is necessary. This metric measures the proportion of solutions that contain the necessary tools for solving the task. It checks whether the solution has all the tools that can produce the expected output. A high NR value means that the method has a strong ability in task planning and finding the right tools for the user's request.

#### 6.2.1  Arguments Inference

A) Resource Hallucination Rate (*abbr.* $HR$):

$$P(s^p) = \begin{cases} \text{true}, & s^p \text{contains false resources} \\ \text{false}, & \text{otherwise} \end{cases},$$
(8)

$$HR = \frac{\sum_i^{||S^p||} \mathbb{I}(P(S_i^p))}{||S^p||},$$
(9)

which measures the proportion of hallucination when the method works on our benchmark. It checks whether all the arguments of the tools exist in the input resources or not. A low HR value means that the method rarely leads to hallucinations that are common in LLMs.

B) Resource Type Consistency Rate (*abbr.* $CR$):

$$Q(s^p) = \begin{cases} \text{true}, & \text{No resource type conflict in } s^p \\ \text{false}, & \text{otherwise} \end{cases},$$
(10)

$$CR = \frac{\sum_i^{||S^p||} \mathbb{I}(Q(S_i^p))}{||S^p||},$$
(11)

which checks whether the resource type of each argument is compatible with the tool specification. A high CR value means that the method can correctly infer and assign arguments for each tool.

#### 6.2.2  Solution Evaluation

$$W(s^p) = \begin{cases} \text{true}, & s^p \text{ can solve the task} \\ \text{false}, & \text{otherwise} \end{cases},$$
(12)

$$SE = \frac{\sum_i^{||S^p||} \mathbb{I}(W(S_i^p))}{||S^p||}.$$
(13)

This metric focuses on the rate of feasible solutions that can solve the tasks over our benchmark, regardless of whether it contains irrelevant tools, as long as the tool call chain output the information that is able to solve the user's request.

## 6.3. Prompt Engineering

**Task decomposition**. The prompt in Table 6 is designed for ControlLLM-ChatGPT in Table 4. It guides the ChatGPT to decompose the user request into several subtasks.

**Tool Assessment**. Table 7 outlines a prompt design for the Tool Assessment task, where the AI assistant evaluates tools' suitability for a given task. The assistant's evaluation is captured in the JSON format, including reasoning and a score. The scoring criteria range from 1 to 5, reflecting the tool's relevance to the task. The prompt emphasizes the connection between tool descriptions and task requirements. This prompt guides AI in making informed decisions when assessing tools' utility for a specific task.

**Solution Expert**. The Table 8 presents a prompt design for the Solution Expert task, wherein an AI assistant evaluates and scores solutions based on their relevance to a given task. The assistant's reasoning process is captured in the JSON format, providing insights into the scoring decision. The prompt emphasizes the connection between solution descriptions and task requirements. This prompt guides AI in making thoughtful decisions when scoring solutions.

**Resource Expert**. As shown in Table 9, we outline a prompt structure for the Resource Expert task. In this task, the AI assistant is required to identify the appropriate inputs for various tools based on the provided context. The AI should refrain from inventing non-existent resources. It emphasizes providing explanations for choices and formatting solutions correctly. This prompt assists AI in accurately selecting and explaining resource inputs for diverse tools and tasks within a constrained context.

**Response Generation.**. We design a prompt template for the Response Generation task in Table 10. In this task, the AI assistant is tasked with explaining the process and outcomes using input and inference results. The AI is instructed to respond directly to the user's request, followed by describing the task's procedure, offering analysis, and presenting model inference results using a first-person perspective. If the inference results involve file paths, the complete path should be provided, or if there are no results, the AI should communicate its inability. The prompt sets the context for generating informative and user-understandable responses in the Response Generation task.

---

| The Prompt of Tool Assessment |
|---|
| The following is a friendly conversation between a human and an AI. The AI is professional and parse user input to several tasks with lots of specific details from its context. If the AI does not know the answer to a question, it truthfully says it does not know. The AI assistant can parse user input to several tasks with JSON format as follows: <Solution>["description": task_description, "task": [task_domain_1, task_domain_2], "id": task_id, "dep": dependency_task_id, "args": ["type": "text", "image" or audio, "value": text, image_url or <GEN>-dep_id], "returns":["type": "segmentation", "value": "<GEN>-task_id"]]</Solution>. The "description" should describe the task in details and AI assistant can add some details to improve the user's request without changing the user's original intention. The special tag "<GEN>-dep_id" refers to the one generated text/image/audio/video/segmentation/mask in the dependency task (Please consider whether the dependency task generates resources of this type.) and "dep_id" must be in "dep" list. The special tag "<GEN>-task_id" refers to the one generated text/image/audio/video/segmentation/mask in this task and "task_id" should be in line with field "id" of this task. The "dep" field denotes the ids of the previous prerequisite tasks which generate a new resource that the current task relies on. The "args" field and the "returns" field denote the input resources and output resources of this task, respectively. The type of resource must be in ["text", "image", "line", "normal", "hed", "scribble", "pose", "edge", "bbox", "category", "segmentation", "audio", "video", "segmentation", "mask"], nothing else. The "task" MUST be selected from the following options: "question-answering", "visual-question-answering", "image-generation", "image-editing", "image-perception", "image-processing", "audio-perception", "audio-generation", "audio-editing", "video-question-answering", "video-perception", "video-generation", "video-editing", nothing else. Think step by step about all the tasks that can resolve the user's request. Parse out as few tasks as possible while ensuring that the user request can be resolved. Pay attention to the dependencies and order among tasks. If some args is not found, you cannot assume that they already exist. You can think a new task to generate those args that do not exist or ask for user's help. If the user input can't be parsed, you need to reply empty JSON []. You should always respond in the following format: <br> <Solution><YOUR_SOLUTION></Solution> <br> <YOUR_SOLUTION>should be strictly with JSON format described above. |

Table 6. Our prompt inspired by [26] for task descomposition.

| The Prompt of Tool Assessment |
| --- |
| Given a task and a tool, the AI assistant helps the system to decide whether this tool can process the task. The assistant should focus more on the description of the model and give a score to each tool. The AI assistant respond with JSON format as follows: <Solution>"Thought": thought, "Score": score </Solution>. "Thought" filed records the model's thinking process step by step within 80 words, which give the reasons why assistant gives this score. "Score" filed denotes score that uses to assess whether this tool is useful for this task. Score is in [1, 2, 3, 4, 5]. Here is the scoring criteria: "Score"=1: the tool is totally not related to the task and does not provide any useful output for solving the task. "Score"=2: the tool is somewhat not related to the task and may not provide any useful output for solving the task. "Score"=3: the tool is probably related to the task and provides some intermediate output that is partially helpful for solving the task but it may not be the optimal one. "Score">3: the tool is closely or directly related to the task and provides output that is mostly helpful for solving the task, or that matches the returns of the task with regard to the type. In a nut shell, for the given task, the higher the score, the more useful the tool is. You should always respond in the following format: <Solution>SOLUTION </Solution>'SOLUTION' should strictly comply with JSON format described above. Task description: "{{task}}".\n\n Here is the description of the tool "{{tool_name}}": \n{{tool_name}}: {{tool_description}}\nArgs: \n{{arguments}}\nReturns: \n{{returns}}\n\nThe above information may be useful for AI to make decision. Please refer to the scoring criteria and score the tool {{tool_name}} for this task. Notice that If the tool description contains keywords from the task description, the score of this tool should be greater than or equal to 3. |

Table 7. Our prompt for tool assessment.

| The Prompt of Solution Expert |
| --- |
| Given a task and a solution, The AI assistant needs to score the solution and respond in json format. Please notice that AI assistant should think. The AI assistant should pay more attention to relevance between the description of each tool in the solution and task. The AI assistant respond with JSON format as follows: <Solution>{"Thought": "thought", "Score": score}</Solution>. "Thought" filed records the model's thinking process step by step within 80 words, which give the reasons why assistant gives this score. "Score" filed denotes score that uses to assess whether this tool is useful for this task. "Score" is in [1, 2, 3, 4, 5]. Here is the scoring criteria: "Score"=1: The solution is totally not related to user's request and can not solve the task. "Score"=2: The solution is somewhat not related to user's request and may not solve the task. "Score"=3: The solution is probably related to the user's intention and may solve the task but it may not be the optimal one. "Score">3: The solution is closely or directly related to what the user wants and could satisfactorily solve the task. In a nut shell, the higher the score, the greater the likelihood of the solution solving the given task. You should always respond in the following format: <Solution>'SOLUTION' </Solution>'SOLUTION' should strictly comply with JSON format described above. User's request: "{{request}}" Task description: "{{task}}". Here is the description of the solution: {{solution}} Please refer to the scoring criteria and score this solution based on the task description. You should think carefully before scoring the solution. Notice that If the keywords in the solution are close in meaning to the keywords in the task description, then the score of this solution is at least 3. |

Table 8. Our prompt for solution expert.

| The Prompt of Resource Expert |
|---|
| The AI assistant needs to find the inputs corresponding to each tool from the context and respond in json format. Please notice that AI assistant should never fake the resources that do not exist. The AI assistant can infer the absent input parameters from the context and respond with JSON format as follows: [{"image": "xxx.png"}, {"bbox": "<GEN>-detr-bbox-0"}, "text": "<summarize the text input from context>"]\n AI assistant should always respond in the following format: "<Explanation>[briefly explain your choice here]</Explanation><Solution>'SOLUTION' </Solution>" 'SOLUTION' should be strictly with JSON format described above. User's request: "{{request}}" Task: "{{task_description}}". <Resources>: {{resources}} We use {{tool_name}} to solve this task: '{{tool_name}}': {{tool_description}} Args: {{arguments}} Returns: {{returns}}\n For the type of "text", AI assistant should summarize the content from the context based on the task and the tool's description. For other types of input, AI assistant need to select the inputs from <Resources>. Now we prepare the inputs for {{tool_name}}: {{input}}. Please complete this inputs and return the completed inputs with the format described above like: <Solution>'SOLUTION' </Solution>. |

Table 9. Our prompt for resource expert.

| The Prompt of Response Generation |
|---|
| With the input and the inference results, the AI assistant needs to describe the process and results. The previous stages can be formed as - User Input: {{ *User Input* }}, Task Planning: {{ *Tasks* }}, Model Selection: {{ *Model Assignment* }}, Task Execution: {{ *Predictions* }}. You must first answer the user's request in a straightforward manner. Then describe the task process and show your analysis and model inference results to the user in the first person. If inference results contain a file path, must tell the user the complete file path. If there is nothing in the results, please tell me you can't make it. |

Table 10. Our Prompt Design in Response Generation. We here refer to the prompts from [26] to generate a user-friendly response.