

Flows: Building Blocks of Reasoning and Collaborating AI

Anonymous ACL submission

Abstract

Recent advances in artificial intelligence (AI) have produced highly capable and controllable systems. This creates unprecedented opportunities for structured reasoning as well as collaboration among multiple AI systems and humans. To fully realize this potential, it is essential to develop a principled way of designing and studying such structured interactions. For this purpose, we introduce the conceptual framework *Flows*. Flows are self-contained building blocks of computation, with an isolated state, communicating through a standardized message-based interface. This modular design allows Flows to be recursively composed into arbitrarily nested interactions, with a substantial reduction of complexity. Crucially, any interaction can be implemented using this framework, including prior work on AI–AI and human–AI interactions, prompt engineering schemes, and tool augmentation. We demonstrate the potential of *Flows* on competitive coding, a challenging task on which even GPT-4 struggles. Our results suggest that structured reasoning and collaboration substantially improve generalization, with AI-only Flows adding +21 and human–AI Flows adding +54 absolute points in terms of solve rate. To support rapid and rigorous research, we introduce the `aiFlows` library embodying *Flows*.

1 Introduction

The success of large language models (LLMs) largely lies in their remarkable emergent ability to adapt to information within their context (i.e., prompt) (Brown et al., 2020; Wei et al., 2022; Kojima et al., 2022). By strategically crafting the context, LLMs can be conditioned to perform complex reasoning (Wei et al., 2022; Nye et al., 2021) and effectively utilize external tools (Schick et al., 2023), significantly enhancing their capabilities. Some of the most exciting recent developments involve defining *control flows*, wherein an LLM

controls a set of tools, orchestrated to solve increasingly complex tasks. Examples of such control flows include ReAct (Yao et al., 2023b), AutoGPT (Richards, 2023), or BabyAGI (Nakajima, 2023). However, these represent but a few of the many conceivable control flows, offering only a glimpse into the vast potential of structured LLM interactions. To realize this potential, we need to develop ways for systematically studying such interactions.

Currently, no general yet efficient abstraction exists for effectively modeling structured interactions. Previous work and existing frameworks, such as LangChain (Chase, 2022), Chameleon (Lu et al., 2023), and HuggingGPT (Shen et al., 2023), have converged on modeling agents as entities that use LLMs to select and execute actions towards specific tasks, where the set of possible actions is pre-defined by the available tools. In this view, tools serve a narrow, well-defined goal and can perform sophisticated tasks (e.g., querying a search engine or executing code). However, their behavior is limited to a single interaction. To highlight the implications of this limitation, consider the following scenario: Alice wants to apply for a job at HappyCorp. If Alice is an agent, she would need to explicitly plan the entire process, including preparing the application, sending it, and evaluating it, which may involve a background check, organizing an interview, and more. Alice would need the knowledge and the “computational” ability to plan every detail. Furthermore, unforeseen events may arise (e.g., the interviewer being on parental leave), requiring Alice to adapt. In reality, most of the complexity is hidden from Alice behind an interface to HappyCorp’s hiring process that might itself be composed of sub-processes involving many other *agents* and *tools*. The hiring process, carefully designed by experts, can be reused by many agents, and its sub-processes can be modified or improved with minimal or no impact on the other components. This makes it evident that agents and tools should

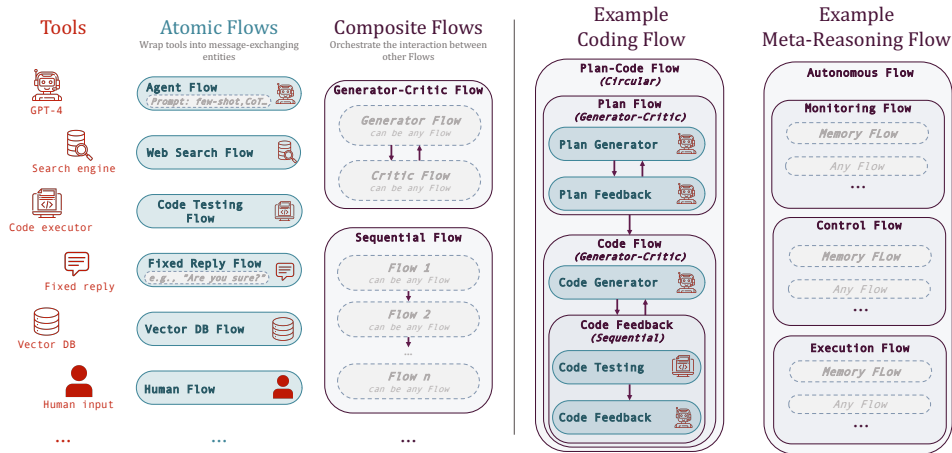


Figure 1: **Flows framework exemplified.** The first column depicts examples of tools. The second column depicts Atomic Flows constructed from the example tools. The third column depicts examples of Composite Flows defining structured interaction between Atomic or Composite Flows. The fourth column illustrates a specific Composite competitive coding Flow as those used in the experiments. The fifth column outlines the structure of a hypothetical Flow, defining a meta-reasoning process that could support autonomous behavior.¹

083 be able to interact in complex, dynamic or static, 114
 084 ways as parts of nested, modular processes, and 115
 085 the distinction between the two becomes blurred as 116
 086 they both serve as computational units in a complex 117
 087 computational process. 118

088 Starting from the observation that everything 119
 089 is a (control) flow defining a potentially complex 120
 090 interaction between many diverse components, we 121
 091 introduce a conceptual framework where Flows are 122
 092 the fundamental building blocks of computation. 123
 093 Flows are independent, self-contained, goal-driven 124
 094 entities able to complete semantically meaningful 125
 095 units of work. To exchange information, Flows 126
 096 communicate via a standardized message-based 127
 097 interface. The framework is depicted in Fig. 1. 128

098 The *Flows* abstraction ensures modularity. Alice, 129
 099 a higher-level meta-reasoning Flow that can sup- 130
 100 port autonomous behavior, does not need to know 131
 101 anything beyond how to interface with Happy- 132
 102 Corp’s hiring Flow. This substantially reduces com- 133
 103 plexity (Alice is interacting with a deeply nested, 134
 104 compositional structured interaction through a sim- 135
 105 ple interface) and provides flexibility, allowing sub- 136
 106 Flows to be swapped without consequences as long 137
 107 as they have the same interface. Indeed, Happy- 138
 108 Corp’s pre-filtering Flow can be swapped from a 139
 109 rule-based system to an AI model or even a human 140
 110 Flow without affecting the structure of the overall 141
 111 process. The abstraction also enables reusability 142
 112 and the composition of sub-Flows into new Flows 143
 113 for different tasks. Furthermore, the framework 144

114 shares key design choices with the Actor model, 115
 116 one of the most prominent models of concurrent 117
 118 computation (cf. Sec. 3). Certainly, once Alice sub- 119
 120 mits her application to HappyCorp, she does not 121
 122 need to wait for the response; she can move to her 123
 124 next goal while the other Flows run concurrently. 125

126 We showcase the potential of the proposed 127
 128 framework and library by investigating complex 129
 130 collaborative and structured reasoning patterns on 131
 132 the challenging task of competitive coding, a mind 133
 134 sport involving participants trying to solve prob- 135
 136 lems defined by a natural language description. 137

138 **Contributions.** (i) We propose *Flows*, a conceptual 139
 140 framework providing an abstraction that enables 141
 142 the design and implementation of arbitrarily nested 143
 144 interactions with a substantial reduction of com- 145
 146 plexity and increase in flexibility in comparison 147
 148 to existing frameworks. *Flows* can represent *any* 149
 150 interaction and provides a common framework for 151
 152 reasoning about interaction patterns, specifying hy- 153
 154 potheses, and structuring research, more broadly. 154
 155 (ii) We open-source the aiFlows library, which 156
 157 embodies *Flows*, together with the visualization 158
 159 toolkit FlowViz and FlowVerse, a repository of 160
 161 Flows that can be readily used, extended, and com- 162
 163 posed into novel, more complex Flows. (iii) We 164
 165 leverage *Flows* and the accompanying library to 166
 167 systematically investigate the benefits of complex 168
 169 interactions for solving competitive coding prob- 169
 170 lems and develop AI-only Flows adding +21 and 170
 171 human–AI Flows adding +54 absolute points in 171
 172 terms of solve rate. 172

¹For more details on meta-reasoning Flows see Sec. 7

2 Related Work

Existing libraries for modeling structured interactions. LangChain (Chase, 2022) has become the go-to library for creating applications using large language models. However, most recent works involving structured interaction, such as Cameleon (Lu et al., 2023), Camel (Li et al., 2023), HuggingGPT (Shen et al., 2023), AutoGPT (Richards, 2023), and BabyAGI (Nakajima, 2023) come with their own library. We argue that the reason why researchers opt to implement bespoke solutions is the lack of a general yet efficient abstraction for modeling structured interactions that makes it easy to explore novel ideas. *Flows*, with its modular design, provides such an abstraction (cf. Appendix A.3).

Competitive coding (CC). With the advent of transformers, Li et al. (2022) finetuned an LLM on GitHub code repositories and a dataset scraped from Codeforces. Recently, Zelikman et al. (2022) proposed decomposing CC problems into function descriptions and, for each function description, using an LLM to generate the implementation in a modular way. While these methods yield promising results, CC remains a challenging task far from being solved (OpenAI, 2023). This is why it presents itself as an ideal testbed for studying collaborative and structured reasoning interactions.

3 Flows

This section introduces *Flows* as a conceptual framework, describes its benefits, and presents the aiFlows library, which embodies the framework.

3.1 Flows as a Conceptual Framework

The framework is centered around *Flows* and *messages*. Flows represent the fundamental building block of computation. They are independent, self-contained, goal-driven entities able to complete a semantically meaningful unit of work. To exchange information, Flows communicate via a standardized message-based interface. Messages can be of any type the recipient Flow can process.

We differentiate between two types of Flows: Atomic and Composite.² Atomic Flows complete the work directly by leveraging *tools*. Tools can be as simple as a textual sequence specifying a (simple) Flow’s fixed response or as complex as a compiler, a search engine, powerful AI systems

like LLaMA (Touvron et al., 2023a,b), Stable Diffusion (Rombach et al., 2021), and GPT-4; or even a human. Notably, in the *Flows* framework, AI systems correspond to tools. An Atomic Flow is effectively a minimal wrapper around a tool and achieves two things: (i) it fully specifies the tool (e.g., the most basic Atomic Flow around GPT-4 would specify the prompts and the generation parameters); and (ii) it abstracts the complexity of the internal computation by exposing only a standard message-based interface for exchanging information with other Flows. Examples of Atomic Flows include wrappers around chain-of-thought prompted GPT-4 for solving math reasoning problems, few-shot prompted LLaMA for question answering, an existing chatbot, a search engine API, or an interface with a human.

Composite Flows accomplish more challenging, higher-level goals by leveraging and coordinating other Flows. Crucially, thanks to their local state and standardized interface, Composite Flows can readily invoke Atomic Flows or other Composite Flows as part of compositional, structured interactions of arbitrary complexity. Enabling research on effective patterns of interaction is one of the main goals of our work. General examples of such patterns include (i) factorizing the problem into simpler problems (i.e., divide and conquer); (ii) evaluating (sub-)solutions at inference time (i.e., feedback); and (iii) incorporating external information or a tool. Importantly, Flows can readily invoke other, potentially heavily optimized, specialized Flows to complete specific (sub-)tasks as part of an interaction, leading to complicated behavior. One example of a Composite Flow is ReAct (Yao et al., 2023b). ReAct is a sequential Flow that structures the problem-solving procedure in two steps: a Flow selects the next action out of a predefined set of actions, and another Flow executes it. The two steps are performed until an answer is obtained. Another prominent example, AutoGPT, extends the ReAct Flow with a Memory Flow and an optional Human Feedback Flow. More generally, our framework provides a unified view of prior work, which we make explicit in Appendix A.3.

Importantly, as illustrated in Fig. 1, Composite Flows can script an arbitrarily complex pattern (i) precisely specifying an interaction (e.g., generate code, execute tests, brainstorm potential reasons for failure, etc.); or (ii) defining a high-level, meta-reasoning process in which a Flow could bring about dynamic unconstrained interactions.

²The concept of a Flow is sufficient for modeling any interaction. We introduce this distinction as it improves the exposition and simplifies the implementation.

244	Key properties. The proposed framework is characterized by the following key properties:	292
245		293
246	• Flows are the compositional building blocks of computation.	294
247		295
248	• Flows encapsulate a local, isolated state.	296
249		297
250	• Flows interact only via messages.	298
251		299
252	• Flows’ behaviour depends only on their internal state and the input message.	300
253		301
254	• Flows can send messages to other Flows and create new Flows.	302
255		303
256	Connection to the Actor model. <i>Flows</i> is fundamentally a framework modeling the computation underlying interactions. As such, it shares key design principles with the <i>Actor</i> model (Hewitt et al., 1973) — a mathematical model of concurrent computation. Similarly to <i>Flows</i> , in the <i>Actor</i> model, an Actor is a concurrent computation entity that can communicate with other Actors exclusively through an asynchronous message-passing interface. By encapsulating the state and the computation within individual Actors, the model provides a high-level abstraction for effectively managing and reasoning about complex concurrent and distributed systems, completely avoiding issues associated with shared states, race conditions, and deadlocks. These benefits are similar in nature to those observed in the domain of interactions. The main distinction between the proposed framework and the <i>Actor</i> model lies in their respective communication protocols. Concretely, while the <i>Actor</i> model prescribes purely asynchronous communication, <i>Flows</i> natively supports synchronous communication, which is essential for the implementation of structured reasoning. Interestingly, a similar deviation from the “pure” <i>Actor</i> model can be identified in the implementation of Erlang, a concurrent programming language based on it (Armstrong, 2003). Overall, the shared design choices still make <i>Flows</i> inherently concurrency-friendly from the practical perspective and are sufficient for important results from the five decades of extensive studies of the <i>Actor</i> model, such as the fact that every physically possible computation can be directly implemented using Actors (Hewitt, 2010), to transfer to <i>Flows</i> .	304
257		305
258		306
259		307
260		308
261		309
262		310
263		311
264		312
265		313
266		314
267		315
268		316
269		317
270		318
271		319
272		320
273		321
274		322
275		323
276		324
277		325
278		326
279		327
280		328
281		329
282		330
283		331
284		332
285		333
286		334
287		335
288	3.2 Why Flows?	336
289	Modularity. <i>Flows</i> introduces a higher-level abstraction that isolates the state of individual Flows and specifies message-based communication as the	337
290		338
291		339
	only interface through which Flows can interact. This ensures perfect modularity by design.	
	Reduction of complexity. The framework ensures the complexity of the computation performed by a Flow is fully abstracted behind the universal message-based interface. This enables an intuitive and simple design of arbitrarily complex interactions from basic building blocks.	
	Systematicity, flexibility, and reusability. The separation of responsibility allows for modules to be developed and studied systematically in isolation or as part of different interactions. Once the correctness and the benefits of a Flow have been established, it can be readily used in developing novel Flows or as a drop-in replacement for less effective Flows leveraged in completing similar goals.	
	Concurrency. The proposed framework’s design is consistent with the Actor model, one of the most prominent models of concurrent computation. As a consequence, <i>Flows</i> can readily support any setting in which Flows run concurrently.	
	3.3 The aiFlows Library	
	Accompanying <i>Flows</i> , we release the aiFlows library, which embodies the framework. In addition to the inherent benefits that come with the framework, the library comes with the following add-ons: (i) FlowVerse: a repository (to which anyone can contribute) of Flows that can be readily used, extended, or composed into novel, more complex Flows. <i>Flows</i> allows for existing “tools” (as well as “models”, “chains”, “agents”, etc.) to be readily incorporated by wrapping them in an Atomic Flow; (ii) a detailed logging infrastructure enabling transparent debugging, analysis, and research in optimizing (i.e., learning or fine-tuning) Flows; (iii) FlowViz: a visualization toolkit to examine the Flows’ execution through an intuitive interface.	
	4 Competitive Coding Flows	
	This work investigates the potential of structured interactions for solving competitive coding (CC) problems. In CC, given a natural language description and a few input–output examples, the task is to generate code that will produce the expected output for all of the hidden input–output test cases associated with the problem. Fig. 4 provides examples.	
	We focus the analysis on three canonical dimensions of interactions: (i) problem decomposition as structured reasoning; (ii) human-AI collabora-	

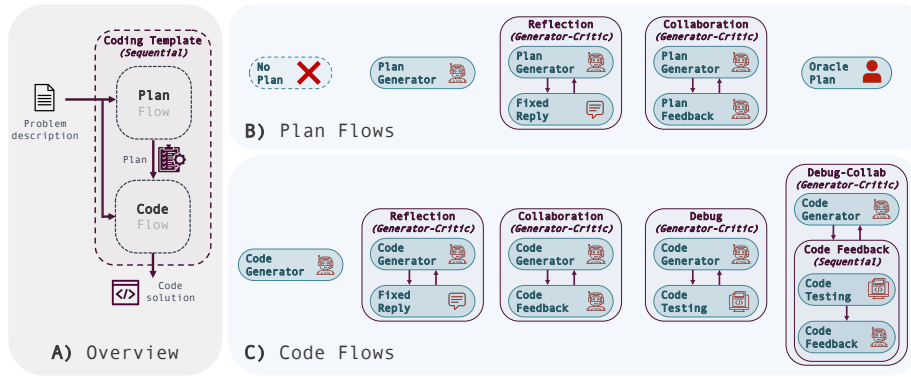


Figure 2: **Competitive coding Flows.** At the highest level, we consider planning as a specific structured reasoning pattern for problem decomposition. In particular, the Plan Flow generates a solution strategy and passes it to the Code Flow, which implements it, as depicted in A). B) and C) depict the different choices of sub-Flows used as Plan and Code Flows in the experiments. Notably, we explore the impact of human-AI collaboration at the plan level and refinement with different types of *feedback*: i) fixed reply encouraging reflection; ii) AI generated feedback; iii) code testing results as feedback; iv) AI generated feedback grounded in code testing results.

tion; and (iii) refinement with various feedback types. By providing a common language for clearly specifying interactions as well as the capability to flexibly compose, exchange, and extend them, the framework makes it possible to study the space of complex interactions in a principled fashion. In the rest of the section, we describe the specific Flows used in the experiments, depicted in Fig. 2.

Problem decomposition. Planning has been an integral intermediate step in recent work (Lu et al., 2023; Shen et al., 2023; Yao et al., 2023b). Similar decomposition is natural in the context of CC as well. In particular, we approach the task in two steps: generating a solution strategy by a Plan Flow and then generating the corresponding code by a Code Flow. This is depicted by panel A in Fig. 2.

Human-AI collaboration. When designing human-AI collaborations, it is essential to take the costs of human interaction into account (Horvitz, 1999; Amershi et al., 2019; Mozannar et al., 2023). By providing immense flexibility, *Flows* can support research in the design of interactions involving humans as computational building blocks in a way that maximizes the utility of the overall computation with a minimal human effort. In the context of CC, we hypothesize that a human can be effectively incorporated at the plan level to provide a short “oracle” plan in natural language. We operationalize this by an (Atomic) Human Flow, illustrated in Panel B of Fig. 2 as the *Oracle Plan Flow*.

Refinement with various feedback types. Iterative refinement is a general problem-solving strategy successfully deployed across various dis-

ciplines (Perrakis et al., 1999; Reid and Neubig, 2022; Schick et al., 2022; Saharia et al., 2021). The strategy revolves around the idea that a solution can be gradually improved through a mechanism for analysis, modification, and re-evaluation. The design of this “*feedback*” mechanism is critical for the effectiveness of the problem-solving strategy. The conceptual framework, paired with the accompanying library, provides the infrastructure to support the design, implementation, and principled research of effective refinement strategies and feedback mechanisms. In this work, we consider a canonical iterative refinement setup where a *generator* Flow is tasked with generating the solution, and a *critic* Flow provides feedback on the proposed solution. We consider two feedback types in the context of both the Plan and the Code Flow: (i) Reflection Flow: the feedback consists of a fixed message encouraging the model to reflect on important aspects of the proposed solution; (ii) Collaboration Flow: the feedback is provided by an AI system that “evaluates” the proposed solution. Furthermore, we explore two more code-specific feedback types: (i) Debug Flow: the feedback message corresponds to the results from executing the code and testing it against the examples provided in the problem description; (ii) Debug-Collab Flow: the feedback is provided by an AI system with access to the code testing results, effectively, grounding the feedback and allowing more systematic reasoning about the potential causes of failure.

We refer to Flows using the following convention: *CodeFlowName* when no plan is generated and *PlanFlowName-CodeFlowName* otherwise.

5 Experimental Setup

Data. We scrape publicly available problems from one of the most popular websites hosting CC contests, Codeforces (Mirzayanov, 2023), and LeetCode (LeetCode, 2023), which cover a broad spectrum of problems ranging from easy interview questions to hard CC problems (see Appendix A.1 for more details). The datasets cover problems from 2020-August-21 to 2023-March-26 for Codeforces, and from 2013-October-25 to 2023-April-09 for LeetCode. Importantly, to study the effect of structured interactions (i.e., different Flows) in a principled manner, it is crucial to account for the possibility of *data contamination*, i.e., that some of the test data has been seen during training (Magar and Schwartz, 2022). Containing problems published over an extended period up to a few months ago (at the time of writing), our datasets allow for reliable identification of the training data cutoff date that can help with addressing this issue. Prior code evaluation datasets like APPS (Hendrycks et al., 2021), HumanEval (Chen et al., 2021), and CodeContests (Li et al., 2022) lack problem release dates, and considering the lack of publicly available information about LLMs’ training data, can likely lead to confounded evaluation of models’ memorization and generalization abilities.

Code testing and solution evaluation. Just like a human participant, the Debug Flow has access only to the input–output example pairs contained in the problem description and, at inference time, uses a local code testing infrastructure to evaluate (intermediate) solution candidates. Crucially, these examples cover only a few simple cases, and generating outputs consistent with them does not imply the code corresponds to a correct solution. A solution is considered correct if it passes all the hidden test cases. To determine correctness, we leverage online evaluators that submit candidate solutions to the websites’ online judges, ensuring authoritative results. For many of the Codeforces problems, we also support local evaluation based on a comprehensive set of hidden test cases we managed to scrape. For more details, see Appendix A.2.

Models and Flows. We experiment with the competitive coding Flows described in Sec. 4, and GPT-4 (OpenAI, 2023) as the LLM tool of choice. See Appendix A.4 for the specific prompts. Also, the code to reproduce the experiments in the paper is available in the project’s GitHub repository.

Evaluation metrics. The most common evaluation metric for code generation is $\text{pass}@k$, corresponding to the probability that in a set of k sampled candidates, there will be at least one correct solution (Chen et al., 2021). To better align with practical use cases, we focus on $\text{pass}@1$, i.e. the solve rate when averaged across the problem set. We report a point estimate and a 95% confidence interval constructed from 1000 bootstrap resamples.

Compute and cost. All the experiments, including the most complex Flows, can be performed on commodity hardware relatively cheaply. For instance, the costs associated with querying the OpenAI API for generating Table 1 amount to \$1000.

6 Experimental Results

We first study the generalization ability of representative Flows and empirically identify GPT-4’s knowledge-cutoff date. Next, we perform a focused analysis along the dimensions described in Sec. 4.

6.1 Performance of Coding Flows on Pre- vs. Post-Knowledge-Cutoff-Date Data

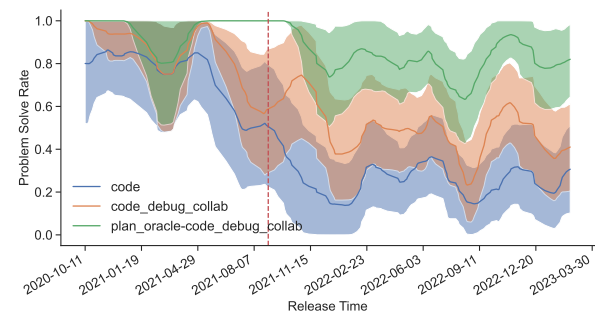


Figure 3: **Temporal analysis.** Performance is averaged over a sliding window of two months. The substantial drop in performance around the reported knowledge cutoff date for GPT-3/4 (the crimson vertical line) reveals limited generalization ability that can be alleviated through structured interactions.

In this experiment, we consider three representative Flows: (i) Code: the simplest Code Generator Flow corresponding to a single GPT-4 API call; (ii) Code_Debug_Collab: the most complex code Flow; (iii) Plan_Oracle-Code_Debug_Collab: the most complex code Flow with human guidance at the plan level. We perform the analysis by running the three Flows on Codeforces problems released from October 2020 to April 2023 and averaging the performance over a sliding window of two months. The results are reported in Fig. 3.

We observe a substantial drop in performance centered around September 2021, consistent with the knowledge cutoff date reported by OpenAI, and denote it by a vertical line on the plot. With Codeforces problems appearing in contexts outside of the contest itself (e.g., editorials), it is reasonable to assume the model has been exposed to older problems more frequently during training. This would explain why the drop spans multiple months, from May 2021 to November 2021, depending on when which data was published and crawled.

Notably, there is a stark difference in the performance of the Code Flow on problems published before and after the knowledge cutoff data, with the solve rate decreasing from around 80% to 23%. While still experiencing a substantial performance drop, the Code_Debug_Collab Flow doubles the solve rate on novel problems to around 45%. Provided with human input at the plan level, the same Flow reaches 85%. Overall, this highlights that GPT-4 performs poorly on novel complex reasoning problems, but structured interactions have the potential to enhance its generalization capabilities. As both GPT-4 (i.e., the Code Flow) and the more complex interactions (Flows) exhibit qualitatively different behavior on novel data, to draw accurate conclusions, it is critical that data contamination is taken into serious consideration when designing experiments and interpreting results.

6.2 Comparing Competitive Coding Flows

Table 1 reports the performance of the systematically chosen set of Flows described in Sec. 4. Rows 6–10 correspond to Flows comprising planning and coding, while rows 1–5 perform the coding directly. In line with the findings of the previous section, we separately consider the performance on problems published before and after the knowledge cutoff date of September 2021.

Problem decomposition. The idea behind planning before implementing the solution is to decouple the high-level reasoning from the code implementation. To analyze the effectiveness of this pattern, we compare the Code and the Plan-Code Flow. Looking at the point estimates, in the pre-cutoff problems, introducing the plan Flow leads to decreased performance (-1.6 for Codeforces and -3.1/2.3/-9.2 for LeetCode easy/medium/hard). However, in the post-cutoff problems, incorporating a plan Flow leads to gains for Codeforces (+8) and LeetCode easy and medium (+2.3 and +3.2).

While these trends are consistent, considering the confidence intervals, we see that they are not statistically significant. Crucially, these results do not imply that this specific problem decomposition is not valuable as it creates a lot of potential in designing an effective human-AI collaboration.

Human-AI collaboration. After every contest, the Codeforces community publishes an editorial that, in addition to the code implementation, provides a short natural language description of the solution. To simulate a Flow where a human provides high-level guidance at the core of the reasoning process, we scrape the solution descriptions and pass them as human-generated plans. The results are striking: despite being only a few sentences long, human-provided plans lead to a substantial performance increase (from 26.9% to 74.5% and from 47.5% to 80.8% on novel problems, when the code is generated by Code and Code_Debug_Collab Flows, respectively). First and foremost, these results showcase the opportunities created by *Flows* for designing, implementing, and studying Human-AI collaboration as a key component of structured interactions. Second, specific to the problem of competitive coding, they validate the hypothesis that high-quality plans are important, suggesting that the design of more effective plan Flows is a promising direction to explore in the future. Last but not least, the results highlight the necessity of more systematic research, as patterns seemingly not valuable in one Flow, such as the simple plan-code structured reasoning problem decomposition, can provide immense value as part of another Flow.

Refinement with various feedback types. Among the code Flows, we find that Code_Reflection and Code_Collaboration lead to limited improvements. The two exceptions are Codeforces pre-cutoff (+9.3) for the former and Codeforces post-cutoff (+9.6) for the latter pattern. While close, these results are not statistically significant. On the other hand, the Flows providing grounded feedback, Code_Debug and Code_Debug_Collab, lead to consistent and statistically significant improvements, most notable on the novel Codeforces problems where performance increases from 26.9, without feedback, to 47.5, when the refinement is based on AI-generated feedback grounded in tests. On LeetCode these improvements are smaller in magnitude. We suspect this is a consequence of the examples provided with the problem description being more simplistic than those in Codeforces,

	Codeforces		Leetcode					
	Pre-cutoff	Post-cutoff	Easy	Pre-cutoff Medium	Hard	Easy	Post-cutoff Medium	Hard
Code	71.8 ±11.0	26.9 ±11.0	97.8 ±3.1	93.4 ±5.4	66.7 ±10.9	76.3 ±8.6	25.1 ±8.9	8.0 ±5.5
Code_Reflection	81.1 ±9.7	26.9 ±10.6	97.8 ±3.1	93.4 ±5.4	67.9 ±10.6	77.4 ±8.1	30.5 ±9.4	11.5 ±6.6
Code_Collaboration	76.6 ±10.5	36.5 ±11.8	97.8 ±3.1	91.1 ±6.0	66.6 ±10.9	73.1 ±8.7	25.1 ±8.7	9.2 ±5.9
Code_Debug	84.5 ±8.6	34.8 ±11.6	97.8 ±3.1	94.5 ±5.0	73.6 ±10.0	84.0 ±7.3	32.8 ±9.6	10.4 ±6.3
Code_Debug_Collab	84.4 ±8.9	47.5 ±12.1	97.8 ±3.1	93.4 ±5.4	72.2 ±10.4	83.8 ±7.4	34.9 ±9.7	9.2 ±6.0
Plan-Code	70.2 ±11.0	34.9 ±11.6	94.7 ±4.5	91.1 ±5.9	57.0 ±11.2	78.6 ±8.3	28.3 ±9.1	4.6 ±4.3
Plan_Reflection-Code	68.5 ±11.6	31.7 ±11.6	95.7 ±4.1	88.9 ±6.6	63.6 ±10.7	77.5 ±8.3	21.8 ±8.5	8.0 ±5.5
Plan_Collaboration-Code	67.0 ±11.5	33.2 ±11.4	96.7 ±3.7	91.1 ±6.1	59.5 ±11.2	74.3 ±8.6	25.2 ±9.0	9.2 ±5.8
Plan_Oracle-Code	82.8 ±9.4	74.5 ±10.7	-	-	-	-	-	-
Plan_Oracle-Code_ Debug_Collab	95.4 ±5.2	80.8 ±9.5	-	-	-	-	-	-

Table 1: **Main Results.** Performance of competitive coding Flows on Codeforces and LeetCode.

leading to false positives and, thereby, incorrect grounding, affecting the feedback quality. This could be addressed by generating additional tests with a Test_Case_Generator Flow, a direction we leave for future work to explore. Finally, in the plan Flows, where we consider Reflection and Collaboration (without grounding), we find that refinement does not provide statistically significant benefits.

Overall, our findings offer several important insights: (i) the direct benefit of problem decomposition hinges on the quality of the intermediate steps; (ii) involving humans at the core high-level reasoning process yields major improvements as humans can easily provide high-quality, grounded feedback; (iii) strategic problem decomposition is a powerful strategy for creating opportunities for effective Human–AI collaboration; (iv) the effectiveness from refinement patterns is not universal and depends on the quality of the starting solution and the feedback (e.g., the level of grounding), and the model’s ability to incorporate that feedback modulated through the feedback’s specificity and the model’s capabilities. The analysis paints a substantially more complex picture than what is reported by prior work for simple interactions.

7 Discussion and Conclusion

Simplicity and systematicity. Thanks to its key properties, *Flows*, together with aiFlows, provides an infrastructure that greatly simplifies the design and implementation of open-ended interactions, with a capability to flexibly isolate, compose, replace, or modify sub-Flows. The experiments demonstrate that carefully designed interactions can substantially improve generalization. However, our analysis also reveals that the effectiveness of particular interaction patterns is not universal; instead, there are many factors at play. As

researchers, we need to clearly specify the patterns we are studying, clearly communicate our hypotheses, and study them both in isolation and as sub-parts of other interactions across different datasets or/and tasks. Furthermore, it is critical that data contamination is taken into serious consideration when designing experiments and drawing conclusions, and error bars become a standard in the field.

Cost and Performance Optimization. In our experiments, we used “off-the-shelf” LLMs that have not been specifically optimized for collaboration. We posit that we can substantially improve performance and/or compute cost by fine-tuning models to collaborate more effectively, generally or toward specialized roles (e.g., controller or critic). To support research in this direction, aiFlows implements detailed logging mechanisms of Flow runs.

Meta-reasoning Flows. Cognitive science research in metacognition and meta-reasoning suggests the existence of meta-level monitoring and control processes underlying cognition (Ackerman and Thompson, 2017). Exploring the development of similar mechanisms in the context of powerful autonomous AI systems and moving beyond a single LLM call serving as a controller (Nakajima, 2023; Richards, 2023) could be a promising area of research. Flows can support such research in higher-level meta-reasoning patterns of interaction.

On the one hand, *Flows* provides a high-level abstraction enabling the design and implementation of interactions of arbitrary complexity. On the other, it offers a common framework for reasoning about interaction patterns, specifying hypotheses, and structuring research. We hope the framework will serve as a solid basis for practical and theoretical innovations, paving the way toward ever more useful AI, similar to the Actor model’s role for concurrent and distributed systems.

665 Limitations

666 **Cost and latency.** aiFlows is fundamentally a
667 framework modeling the computation that under-
668 lies structured reasoning and collaboration, which
669 inherently involves multiple calls. Naturally, this
670 will result in higher latency, which impacts the user
671 experience, and cost in comparison to a single call.

672 **Evaluation limitations..** This work provides the
673 infrastructure to support a systematic study of struc-
674 tured interactions, and demonstrates its utility by
675 providing a thorough evaluation using a single
676 model and a specific subset of interactions on the
677 task of competitive coding. However, as discussed
678 in Sec. 7, many factors determine the effectiveness
679 of structured interactions, and future work should
680 continue exploring the vast space of models and
681 conceivable interactions across the many complex
682 tasks that can be addressed in this setup.

683 **Risk and biases associated with tools.** Flows rely
684 on the computation performed by the tools (e.g.,
685 LLMs, search engines, etc.) and, therefore, will
686 be exposed to the risks and biases associated with
687 their usage.

688 **Cost and performance optimization.** As dis-
689 cussed in Sec. 7, the "off-the-shelf" LLM used in
690 the experiments has not been specifically optimized
691 for effectiveness in structured interactions. Albeit
692 for the better, fine-tuning with aiFlows in mind
693 would substantially affect cost and performance.

694 References

695 Rakefet Ackerman and Valerie A. Thompson. 2017.
696 [Meta-reasoning: Monitoring and control of think-
697 ing and reasoning.](#) *Trends in Cognitive Sciences*,
698 21(8):607–617.

699 Saleema Amershi, Daniel S. Weld, Mihaela Vorvoreanu,
700 Adam Fourney, Besmira Nushi, Penny Collisson,
701 Jina Suh, Shamsi T. Iqbal, Paul N. Bennett,
702 Kori Inkpen, Jaime Teevan, Ruth Kikin-Gil, and Eric
703 Horvitz. 2019. [Guidelines for human-ai interaction.](#)
704 In *Proceedings of the 2019 CHI Conference on Human
705 Factors in Computing Systems, CHI 2019, Glasgow,
706 Scotland, UK, May 04-09, 2019*, page 3. ACM.

707 Joe Armstrong. 2003. [Making reliable distributed sys-
708 tems in the presence of software errors.](#) Ph.D. thesis,
709 Royal Institute of Technology, Stockholm, Sweden.

710 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
711 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
712 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
713 Askell, Sandhini Agarwal, Ariel Herbert-Voss,

Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners.](#) In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc. 714 715 716 717 718 719 720 721 722 723

Harrison Chase. 2022. Langchain. <https://github.com/hwchase17/langchain>. 724 725

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde, Jared Kaplan, Harrison Edwards, Yura Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, David W. Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William H. Guss, Alex Nichol, Igor Babuschkin, S. Arun Balaji, Shantanu Jain, Andrew Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew M. Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. [Evaluating large language models trained on code.](#) *ArXiv*, abs/2107.03374. 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744

Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2022. [Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks.](#) *ArXiv*, abs/2211.12588. 745 746 747 748 749

Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2023. [Teaching large language models to self-debug.](#) *ArXiv*, abs/2304.05128. 750 751 752

Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. 2021. [Measuring coding challenge competence with apps.](#) *NeurIPS*. 753 754 755 756 757

Carl E. Hewitt. 2010. [Actor model of computation: Scalable robust information systems.](#) *arXiv: Programming Languages*. 758 759 760

Carl E. Hewitt, Peter Boehler Bishop, and Richard Steiger. 1973. [A universal modular actor formalism for artificial intelligence.](#) In *International Joint Conference on Artificial Intelligence*. 761 762 763 764

Eric Horvitz. 1999. [Principles of mixed-initiative user interfaces.](#) In *Proceeding of the CHI '99 Conference on Human Factors in Computing Systems: The CHI is the Limit, Pittsburgh, PA, USA, May 15-20, 1999*, pages 159–166. ACM. 765 766 767 768 769

770	Geunwoo Kim, Pierre Baldi, and Stephen McAleer.	821
771	2023. Language models can solve computer tasks.	822
772	<i>ArXiv</i> , abs/2303.17491.	823
773	Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yu-	824
774	taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-	825
775	guage models are zero-shot reasoners . In <i>Advances in</i>	826
776	<i>Neural Information Processing Systems</i> , volume 35,	827
777	pages 22199–22213. Curran Associates, Inc.	828
778	LeetCode. 2023. Leetcode.com .	829
779	Guohao Li, Hasan Abed Al Kader Hammoud, Hani	
780	Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023.	
781	Camel: Communicative agents for "mind" explo-	
782	ration of large scale language model society. <i>arXiv</i>	
783	<i>preprint arXiv:2303.17760</i> .	
784	Yujia Li, David Choi, Junyoung Chung, Nate Kushman,	
785	Julian Schrittwieser, Rémi Leblond, Tom Eccles,	
786	James Keeling, Felix Gimeno, Agustin Dal Lago,	
787	et al. 2022. Competition-level code generation with	
788	alphacode. <i>Science</i> , 378(6624):1092–1097.	
789	Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-	
790	Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jian-	
791	feng Gao. 2023. Chameleon: Plug-and-play compo-	
792	sitional reasoning with large language models. <i>ArXiv</i> ,	
793	abs/2304.09842.	
794	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	
795	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	
796	Nouha Dziri, Shrimai Prabhunoye, Yiming Yang,	
797	et al. 2023. Self-refine: Iterative refinement with	
798	self-feedback. <i>arXiv preprint arXiv:2303.17651</i> .	
799	Inbal Magar and Roy Schwartz. 2022. Data contamina-	
800	tion: From memorization to exploitation . In <i>Proceed-</i>	
801	<i>ings of the 60th Annual Meeting of the Association for</i>	
802	<i>Computational Linguistics (Volume 2: Short Papers)</i> ,	
803	<i>ACL 2022, Dublin, Ireland, May 22-27, 2022</i> , pages	
804	157–165. Association for Computational Linguistics.	
805	Mike Mirzayanov. 2023. Codeforces.com .	
806	Hussein Mozannar, Gagan Bansal, Adam Fourney, and	
807	Eric Horvitz. 2023. When to show a suggestion? inte-	
808	grating human feedback in ai-assisted programming .	
809	<i>CoRR</i> , abs/2306.04930.	
810	Yohei Nakajima. 2023. Babyagi. https://github.	
811	com/yoheinakajima/babyagi .	
812	Maxwell I. Nye, Anders Johan Andreassen, Guy Gur-	
813	Ari, Henryk Michalewski, Jacob Austin, David	
814	Bieber, David Dohan, Aitor Lewkowycz, Maarten	
815	Bosma, David Luan, Charles Sutton, and Augustus	
816	Odena. 2021. Show your work: Scratchpads for inter-	
817	mediate computation with language models . <i>CoRR</i> ,	
818	abs/2112.00114.	
819	OpenAI. 2023. Gpt-4 technical report. <i>ArXiv</i> ,	
820	abs/2303.08774.	
	Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beat-	
	riz Borges, Antoine Bosselut, Robert West, and	
	Boi Faltings. 2023. Refiner: Reasoning feedback	
	on intermediate representations . <i>arXiv preprint</i>	
	<i>arXiv:2304.01904</i> .	
	Anastassis Perrakis, Richard J. Morris, and Victor S.	
	Lamzin. 1999. Automated protein model building	
	combined with iterative structure refinement . <i>Nature</i>	
	<i>Structural Biology</i> , 6:458–463.	
	Machel Reid and Graham Neubig. 2022. Learning to	
	model editing processes . In <i>Conference on Empirical</i>	
	<i>Methods in Natural Language Processing</i> .	
	Toran Bruce Richards. 2023. Autogpt. https://	
	github.com/Significant-Gravitas/Auto-GPT .	
	Robin Rombach, Andreas Blattmann, Dominik Lorenz,	
	Patrick Esser, and Björn Ommer. 2021. High-	
	resolution image synthesis with latent diffusion mod-	
	els . <i>CoRR</i> , abs/2112.10752.	
	Chitwan Saharia, Jonathan Ho, William Chan, Tim	
	Salimans, David J. Fleet, and Mohammad Norouzi.	
	2021. Image super-resolution via iterative refinement .	
	<i>IEEE Transactions on Pattern Analysis and Machine</i>	
	<i>Intelligence</i> , 45:4713–4726.	
	Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta	
	Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola	
	Cancedda, and Thomas Scialom. 2023. Toolformer:	
	Language models can teach themselves to use tools.	
	<i>ArXiv</i> , abs/2302.04761.	
	Timo Schick, Jane Dwivedi-Yu, Zhengbao Jiang, Fabio	
	Petroni, Patrick Lewis, Gautier Izacard, Qingfei You,	
	Christoforos Nalmpantis, Edouard Grave, and Sebas-	
	tian Riedel. 2022. Peer: A collaborative language	
	model . <i>ArXiv</i> , abs/2208.11663.	
	Yongliang Shen, Kaitao Song, Xu Tan, Dong Sheng Li,	
	Weiming Lu, and Yue Ting Zhuang. 2023. Hugging-	
	gpt: Solving ai tasks with chatgpt and its friends in	
	huggingface. <i>ArXiv</i> , abs/2303.17580.	
	Noah Shinn, Federico Cassano, Beck Labash, Ash-	
	win Gopinath, Karthik Narasimhan, and Shunyu	
	Yao. 2023. Reflexion: Language agents with	
	verbal reinforcement learning. <i>arXiv preprint</i>	
	<i>arXiv:2303.11366</i> .	
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	
	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	
	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	
	Azhar, Aurélien Rodriguez, Armand Joulin, Edouard	
	Grave, and Guillaume Lample. 2023a. Llama: Open	
	and efficient foundation language models. <i>ArXiv</i> ,	
	abs/2302.13971.	
	Hugo Touvron, Louis Martin, Kevin R. Stone, Peter	
	Albert, Amjad Almahairi, Yasmine Babaei, Niko-	
	lay Bashlykov, Soumya Batra, Prajjwal Bhargava,	
	Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cris-	
	tian Canton Ferrer, Moya Chen, Guillem Cucurull,	
	David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin	

876 Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami,
877 Naman Goyal, Anthony S. Hartshorn, Saghar Hos-
878 seini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor
879 Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V.
880 Korenev, Punit Singh Koura, Marie-Anne Lachaux,
881 Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai
882 Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov,
883 Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew
884 Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan
885 Saladi, Alan Schelten, Ruan Silva, Eric Michael
886 Smith, R. Subramanian, Xiaoqing Ellen Tan, Binh
887 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan,
888 Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang,
889 Angela Fan, Melanie Kambadur, Sharan Narang, Au-
890 relien Rodriguez, Robert Stojnic, Sergey Edunov, and
891 Thomas Scialom. 2023b. Llama 2: Open foundation
892 and fine-tuned chat models. *arXiv*.

893 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
894 Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,
895 et al. 2022. Chain-of-thought prompting elicits rea-
896 soning in large language models. *Advances in Neural
897 Information Processing Systems*, 35:24824–24837.

898 Sean Welleck, Ximing Lu, Peter West, Faeze Brah-
899 man, Tianxiao Shen, Daniel Khashabi, and Yejin
900 Choi. 2023. [Generating sequences by learning to
901 self-correct](#). In *The Eleventh International Confer-
902 ence on Learning Representations*.

903 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,
904 Thomas L. Griffiths, Yuan Cao, and Karthik
905 Narasimhan. 2023a. Tree of thoughts: Deliberate
906 problem solving with large language models. *ArXiv*,
907 abs/2305.10601.

908 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak
909 Shafran, Karthik R Narasimhan, and Yuan Cao.
910 2023b. [React: Synergizing reasoning and acting
911 in language models](#). In *The Eleventh International
912 Conference on Learning Representations*.

913 Ori Yoran, Tomer Wolfson, Ben Bogin, Uri Katz, Daniel
914 Deutch, and Jonathan Berant. 2023. Answering
915 questions by meta-reasoning over multiple chains
916 of thought. *ArXiv*, abs/2304.13007.

917 Eric Zelikman, Qian Huang, Gabriel Poesia, Noah D
918 Goodman, and Nick Haber. 2022. [Parsel: A \(de-
919 \)compositional framework for algorithmic reasoning
920 with language models](#).

921 Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao,
922 George Karypis, and Alexander J. Smola. 2023. Mul-
923 timodal chain-of-thought reasoning in language mod-
924 els. *ArXiv*, abs/2302.00923.

925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973

A Appendix

A.1 Data

Example Codeforces and LeetCode problems are provided in Fig. 4.

In the first experiment, the temporal analysis, we use 239 Codeforces problems ranging from October 2020 to April 2023. In the second experiment, we have 136 problems for Codeforces (some problems are dropped in order to keep the pre-cutoff and post-cutoff buckets equal to 68) and 558 problems for LeetCode (93 for each of the six buckets). Additionally, to support research in the area, we set up an AI competitive coding challenge based on a dataset of Codeforces problems of various difficulties published after the knowledge cutoff date. More details about the CC competition are available in Appendix A.5.

A.2 Code Testing and Solution Evaluation

The solution evaluation requires a set of input–output pairs, hidden from the user, that comprehensively test the behavior of the program. To compute the final results, we have implemented an online evaluation infrastructure that submits the candidate solutions to the websites’ online judges and automatically scrapes the judgment. This mechanism ensures authoritative results.

For many of the Codeforces problems, we managed to scrape (sometimes a subset) of the hidden tests, allowing us to use a faster, local infrastructure for evaluating candidate solutions. On the other hand, LeetCode does not expose any of the hidden tests publicly.

For code testing at inference time, just like a human would, we rely on tests constructed from the (public) input–output example pairs contained in the problem description.

A.3 Concurrent and Previous Works as Specific Instances of Flows

The introduction of LLMs such as BARD, GPT-3, ChatGPT, and its latest version, GPT-4, has led to a breakthrough in AI. This has enabled many exciting developments like CoT, HuggingGPT, AutoGPT, AgentGPT, and BabyAGI. In this section, we demonstrate how *Flows* provides a unified view encompassing concurrent and previous work as specific Flow instances. The details are provided in Figure 5 and Table. 2.

1. **Few shot Prompting (FS)** (Brown et al., 2020) consists in providing a few input-output

examples within the prompt, acting as demonstrations to enable the LLM to perform a specific task. This technique relies on the LLM’s emergent in-context learning ability to extrapolate from these limited examples and infer how to solve the task in general.

2. **Chain of Thoughts (CoT)** (Wei et al., 2022) is a prompting method (atomic Flow) that allows LLMs to generate a series of intermediate natural language reasoning steps that lead to the final output.
3. **Tree of Thoughts (ToT)** (Yao et al., 2023a) is a framework that enables (*orchestration*) exploration over coherent units of text (thoughts) that serve as intermediate steps toward problem-solving. ToT allows LLMs to perform deliberate decision-making by considering multiple different reasoning paths and self-evaluating choices to decide the next course of action, as well as looking ahead or backtracking when necessary to make global choices.
4. **Program of Thoughts (PoT)** (Chen et al., 2022) is a prompting method that allows language models (mainly Codex) to express the reasoning process as a program. The computation is relegated to an external program, which executes the generated programs to derive the answer.
5. **Multimodal CoT (M-CoT)** (Zhang et al., 2023) is a method that incorporates language (text) and vision (images) modalities into a two-stage framework that separates rationale generation and answer inference. To facilitate the interaction between modalities in M-CoT, smaller language models (LMs) are fine-tuned by fusing multimodal features.
6. **ToolFormer** (Schick et al., 2023) is a model that is trained to decide which APIs to call, when to call them, what arguments to pass, and how to incorporate the results into future tokens prediction.
7. **ReAct** (Yao et al., 2023b) is a framework that uses LLMs to generate reasoning traces and task-specific actions sequentially. The framework allows for greater synergy between the two: reasoning traces help the model induce, track, and update action plans and han-

974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021



You have received data from a Bubble bot. You know your task is to make factory facilities, but before you even start, you need to know how big the factory is and how many rooms it has. When you look at the data you see that you have the dimensions of the construction, which is in rectangle shape: $N \times M$.

Then in the next N lines you have M numbers. These numbers represent factory tiles and they can go from 0 to 15. Each of these numbers should be looked in its binary form. Because from each number you know on which side the tile has walls. For example number 10 in its binary form is 1010, which means that it has a wall from the North side, it doesn't have a wall from the East, it has a wall on the South side and it doesn't have a wall on the West side. So it goes North, East, South, West. It is guaranteed that the construction always has walls on its edges. The input will be correct. Your task is to print the size of the rooms from biggest to smallest.

Example 1:
Input: 4 5
 9 14 11 12 13
 5 15 11 6 7
 5 9 14 9 14
 3 2 14 3 14
Output: 9 4 4 2 1



Given an input string (s) and a pattern (p), implement wildcard pattern matching with support for '?' and '*' where:
 '?' Matches any single character.
 '*' Matches any sequence of characters (including the empty sequence).
 The matching should cover the entire input string (not partial).

Example 1:
Input: $s = "aa", p = "a"$
Output: false
Explanation: "a" does not match the entire string "aa".

Figure 4: Examples of competitive coding problems from Codeforces and LeetCode.

1022	dle exceptions, while actions allow it to interface with external sources, such as knowledge bases or environments, to gather additional information.	GPT), one for performing the task and another for criticizing.	1053
1023			1054
1024			
1025			
1026	8. Parsel (Zelikman et al., 2022) is a framework that enables the automatic implementation and validation of complex algorithms with code LLMs. The framework first synthesizes an intermediate representation based on the Parsel language and can then apply a variety of post-processing tools. Code is generated in a next step.	12. Self-Correct (Welleck et al., 2023) is a framework that decouples a flawed base generator (an LLM) from a separate corrector that learns to iteratively correct imperfect generations. The imperfect base generator can be an off-the-self LLM or a supervised model, and the corrector model is trained.	1055
1027			1056
1028			1057
1029			1058
1030			1059
1031			1060
1032			1061
1033			
1034	9. REFINER (Paul et al., 2023) is a framework for LMs to explicitly generate intermediate reasoning steps while interacting with a critic model that provides automated feedback on the reasoning.	13. Self-Debug (Chen et al., 2023) is a framework that relies on external tools (SQL application or Python interpreter) to help large language models revise and debug SQL commands or Python code with bugs.	1062
1035			1063
1036			1064
1037			1065
1038			1066
1039	10. Self-Refine (Madaan et al., 2023) is a framework for LLMs to generate coherent outputs. The main idea is that an LLM will initially generate an output while the same LLM provides feedback for its output and uses it to refine itself iteratively.	14. Reflexion (Shinn et al., 2023) is a framework that provides a free-form reflection on whether a step was executed by LLM correctly or not and potential improvements. Unlike self-refine and self-debug, Reflexion builds a persisting memory of self-reflective experiences, which enables an agent to identify its own errors and self-suggest lessons to learn from its mistakes over time.	1067
1040			1068
1041			1069
1042			1070
1043			1071
1044			1072
1045	11. Recursively Criticize and Improve (RCI) (Kim et al., 2023) showed that a pre-trained large language model (LLM) agent could execute computer tasks guided by natural language using a simple prompting scheme where the agent Recursively Criticizes and Improves its output (RCI). Unlike Self-refine, this method uses two separate LLMs (Chat-	15. Meta-Reasoner (Yoran et al., 2023) is an approach which prompts large language models to meta-reason over multiple chains of thought rather than aggregating their answers. This approach included two steps: (i) ask LLM to generate multiple reasoning chains, (ii) ask another LLM (meta-reasoner) to reason over the multiple reasoning chains to arrive at the correct answer.	1073
1046			1074
1047			1075
1048			
1049			1076
1050			1077
1051			1078
1052			1079
			1080
			1081
			1082
			1083
			1084

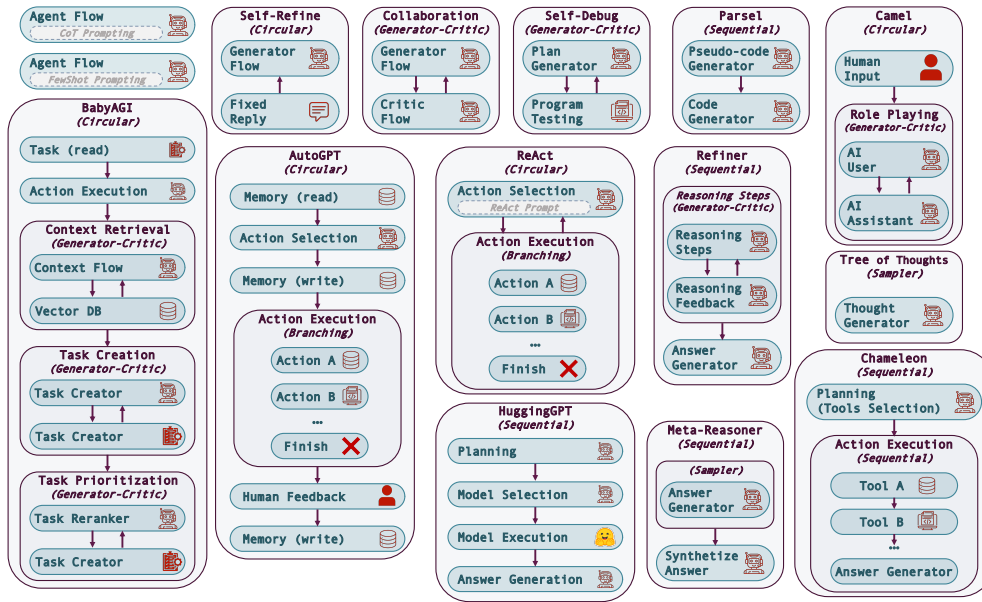


Figure 5: **Previous works are specific Flows.** We depict a selected subset of previous works incorporating structured reasoning and/or interactions between AI agents, tools, and humans, through the lens of the Flows framework. This demonstrates that Flows is a powerful language for describing, conceptualizing, and disseminating structured interaction patterns.

- 1085 16. **HuggingGPT** (Shen et al., 2023) is a framework that leverages LLMs (e.g., ChatGPT) to connect various AI models in machine learning communities (e.g., Hugging Face) to solve numerous sophisticated AI tasks in different modalities (such as language, vision, speech) and domains.
- 1086
- 1087
- 1088
- 1089
- 1090
- 1091
- 1092 17. **Camel** (Li et al., 2023) is a communicative agent framework involving inception prompting to guide chat agents toward task completion while maintaining consistency with human intentions.
- 1093
- 1094
- 1095
- 1096
- 1097 18. **Chameleon** (Lu et al., 2023) is a plug-and-play compositional reasoning framework that augments external tools with LLMs in a plug-and-play manner. The core idea is that an LLM-based planner assembles a sequence of tools to execute to generate the final response. The assumption is that this will be less error-prone, easily expandable to new modules, and user-friendly.
- 1098
- 1099
- 1100
- 1101
- 1102
- 1103
- 1104
- 1105
- 1106 19. **AutoGPT** (Richards, 2023) is an experimental open-source application that leverages the capabilities of large language models (LLMs) and Chatbots such as OpenAI’s GPT-4 and Chat-GPT to create fully autonomous and customizable AI agents. It has internet access,

1120 long-term and short-term memory management.

1121

- 1122 20. **BabyAGI** (Nakajima, 2023) is an intelligent agent capable of generating and attempting to execute tasks based on a given objective. BabyAGI operates based on three LLM flows: Task creation flow, Task prioritization flow, and Execution flow.
- 1123
- 1124
- 1125
- 1126
- 1127
- 1128
- 1129
- 1130
- 1131
- 1132
- 1133
- 1134
- 1135
- 1136
- 1137
- 1138

A.4 Prompting

1120 We provide the prompts used to obtain the results in Section 6. Our evaluation is made possible thanks to the modular and compositional nature of *Flows*. Some of the experimental setups are deeply nested, and in cases where Flows build on each other, we avoid repetition. Note that the project’s GitHub repository provides the code and data to reproduce all of the experiments in the paper.

1121 Direct prompting for a solution is shown in Listing 1. To add reflection, we use a Generator-Critic Flow to combine the code generation with a fixed reply, as shown in Listing 2. In the collaboration setting, we use Listing 3 as the generator and Listing 4 as the critic.

1122 Debugging is incorporated via a testing Flow that adds formatting to the output of a code executor. The formatting templates are shown in Listing 6. To respond to the debug output, we rely on an

Flows	Flow Type	Interactions				Reasoning Patterns		Feedback	Learning
		Self	Multi-Ag.	Human	Tools	Struct.	Plan		
FS (Brown et al., 2020)	Atomic	X	X	X	X	X	X	X	X
CoT (Wei et al., 2022)	Atomic	X	X	X	X	✓	X	X	X
ToT (Yao et al., 2023a)	Circular	✓	X	X	✓	✓	X	X	X
PoT (Chen et al., 2022)	Seq.	X	X	X	✓	✓	X	X	X
M-CoT (Zhang et al., 2023)	Seq.	X	X	X	X	✓	X	X	✓
ToolFormer (Wei et al., 2022)	Seq.	X	X	X	✓	✓	X	X	✓
ReAct (Yao et al., 2023b)	Circular	X	X	X	✓	✓	X	X	X
Parsel (Zelikman et al., 2022)	Seq.	X	✓	X	✓	✓	✓	X	X
REFINER (Paul et al., 2023)	Gen-Crit	X	✓	✓	X	✓	X	✓	✓
Self-Refine (Madaan et al., 2023)	Gen-Crit	✓	X	X	X	✓	X	✓	X
RCI (Kim et al., 2023)	Gen-Crit	✓	X	X	✓	✓	X	✓	X
Self-Correct (Welleck et al., 2023)	Gen-Crit	✓	X	X	✓	✓	X	✓	X
Self-Debug (Chen et al., 2023)	Gen-Crit	✓	X	X	✓	✓	X	✓	X
Reflexion (Shinn et al., 2023)	Gen-Crit	✓	X	X	✓	X	X	✓	X
Meta-Reasoner (Yoran et al., 2023)	Seq.	✓	✓	X	X	✓	X	X	X
HuggingGPT (Shen et al., 2023)	Seq.	X	✓	X	✓	✓	✓	X	X
Camel (Li et al., 2023)	Circular	X	✓	✓	X	✓	X	✓	X
Chameleon (Lu et al., 2023)	Seq.	X	✓	X	✓	✓	✓	X	X
AutoGPT (Richards, 2023)	Circular	✓	✓	X	✓	✓	✓	✓	X
BabyAGI (Nakajima, 2023)	Circular	X	✓	X	✓	✓	✓	X	X

Table 2: **Previous work.** We compare previous work across relevant dimensions.

adjusted coding Flow 5. Adding collaboration in the debugging setting is done by introducing a critic that provides feedback grounded in the test results. This Flow is detailed in Listing 3.

The scenarios explained above also support the addition of a planning Flow. An example of plan generation is shown in Listing 8.

Listing 1: Prompts for Code Flow (Codeforces)

```
"prompt templates":
"system_message": |-
    Your goal is to provide
    executable Python code
    that solves a competitive
    programming problem. The
    code should correctly
    handle all corner cases in
    order to pass the hidden
    test cases, which are used
    to evaluate the
    correctness of the
    solution.

    The user will specify the
    problem by providing you
    with:
    - the problem statement
    - input description
    - output description
    - example test cases
    - (optional) explanation of
    the test cases
```

```
The user will provide you
with a task and an output
format that you will
strictly follow.
"query_message": |-
# Problem statement
{{problem_description}}

# Input description
{{input_description}}

# Output description
{{output_description}}

{{io_examples_and_explanation
}}

The input should be read from
the standard input and
the output should be
passed to the standard
output.
Return Python code that
solves the problem. Reply
in the following format:
```python
{{code_placeholder}}
```
"human_message": |-
{{query}}
```

1139
1140
1141
1142
1143
1144
1145
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200

Listing 2: Prompts for Fixed-Reply Flow

```

1201 "prompt templates":
1202   "fixed_reply": |-
1203     Consider the problem
1204       statement and the last
1205       proposed solution. Are you
1206       sure that the solution is
1207       provided in the requested
1208       format, and crucially,
1209       solves the problem?
1210   If that is not the case,
1211     provide the corrected
1212     version of the code in the
1213     following format:
1214     ```python
1215     {{python_code}}
1216     ```
1217   otherwise, reply:
1218   "Final answer."

```

```

{{problem_description}} 1249
# Input description 1251
{{input_description}} 1252
# Output description 1254
{{output_description}} 1255
{{io_examples_and_explanation 1256
  }} 1257
The input should be read from 1261
the standard input and 1262
the output should be 1263
passed to the standard 1264
output. 1265
Return Python code that 1266
solves the problem. Reply 1267
in the following format: 1268
```python 1269
{{code_placeholder}} 1270
``` 1271

```

Listing 3: Prompts for Code-Collab Flow (Codeforces)

```

1219 "prompt templates":
1220   "system_message": |-
1221     Your goal is to provide
1222     executable Python code
1223     that solves a competitive
1224     programming problem. The
1225     code should correctly
1226     handle all corner cases in
1227     order to pass the hidden
1228     test cases, which are used
1229     to evaluate the
1230     correctness of the
1231     solution.
1232
1233   The user will specify the
1234   problem by providing you
1235   with:
1236   - the problem statement
1237   - input description
1238   - output description
1239   - example test cases
1240   - (optional) explanation of
1241     the test cases
1242
1243   The user will provide you
1244   with a task and an output
1245   format that you will
1246   strictly follow.
1247   "query_message": |-
1248     # Problem statement

```

```

"human_message": |- 1272
# Feedback on the last 1273
proposed solution 1274
{{code_feedback}} 1275
Consider the original problem 1278
statement, the last 1279
proposed solution and the 1280
provided feedback. Does 1281
the solution need to be 1282
updated? If so, provide 1283
the corrected version of 1284
the code in the following 1285
format: 1286
```python 1287
{{code_placeholder}} 1288
``` 1289
otherwise, reply: 1290
"Final answer." 1291

```

Listing 4: Prompts for Code-Collab-Critic Flow (Codeforces)

```

"prompt templates": 1292
"system_message": |- 1293
  Your goal is to identify 1294
  potential issues with a 1295
  competitive programming 1296
  solution attempt. 1297

```

```

1298
1299 The user will specify the
1300     problem by providing you
1301     with:
1302     - the problem statement
1303     - input description
1304     - output description
1305     - example test cases
1306     - (optional) explanation of
1307       the test cases
1308     - a Python solution attempt
1309
1310 Crucially, your goal is to
1311     correctly identify
1312     potential issues with the
1313     solution attempt, and not
1314     to provide the code
1315     implementation yourself.
1316 The user will provide you
1317     with a task and an output
1318     format that you will
1319     strictly follow.
1320 "query_message": |-
1321     # Problem statement
1322     {{problem_description}}
1323
1324     # Input description
1325     {{input_description}}
1326
1327     # Output description
1328     {{output_description}}
1329
1330     {{io_examples_and_explanation
1331       }}
1332
1333     # Python solution attempt:
1334     ```python
1335     {{code}}
1336     ```
1337
1338
1339 Consider the problem
1340     statement and the solution
1341     attempt. Are there any
1342     issues with the proposed
1343     solution or it is correct?
1344     Explain your reasoning
1345     very concisely, and do not
1346     provide code.
1347 "human_message": |-
1348     {{query}}

```

Listing 5: Prompts for Code-Debug Flow (Codeforces)

```

"prompt templates":
"system_message": |-
    Your goal is to provide
    executable Python code
    that solves a competitive
    programming problem. The
    code should correctly
    handle all corner cases in
    order to pass the hidden
    test cases, which are used
    to evaluate the
    correctness of the
    solution.
The user will specify the
    problem by providing you
    with:
    - the problem statement
    - input description
    - output description
    - example test cases
    - (optional) explanation of
      the test cases
The user will provide you
    with a task and an output
    format that you will
    strictly follow.
"query_message": |-
    # Problem statement
    {{problem_description}}
    # Input description
    {{input_description}}
    # Output description
    {{output_description}}
    {{io_examples_and_explanation
      }}
    # Input description
    {{input_description}}
    # Output description
    {{output_description}}
    {{io_examples_and_explanation
      }}
The input should be read from
    the standard input and
    the output should be
    passed to the standard
    output.
Return Python code that
    solves the problem. Reply
    in the following format:

```



```

1399     ``python                                     ## [Failed test] Input      1448
1400     {{code_placeholder}}                          ```````````````````` 1449
1401     ````````````````````                          {{test_input}}        1450
1402 "human_message": |-                               ```````````````````` 1451
1403     {{testing_results_summary}}                   ## [Failed test] Runtime  1452
1404     error message                                  1453
1405     {{error_message}}                              1454
1406     Consider the problem                           "single test error": |- 1455
1407     statement, the last                            ${.issue_title}        1456
1408     proposed solution, and its                     The Python code does not 1457
1409     issue. Provide a                               solve the problem in the 1458
1410     corrected version of the                       problem description due to 1459
1411     code that solves the                           logical errors. It fails 1460
1412     original problem and                           the following test:      1461
1413     resolves the issue,                            ## [Failed test] Input  1462
1414     without any explanation,                       ```````````````````` 1463
1415     in the following format:                        {{test_input}}        1464
1416     ``python                                       ```````````````````` 1465
1417     {{code_placeholder}}                           ## [Failed test] Expected 1466
1418     ````````````````````                          output                 1467
1419                                                                 ```````````````````` 1468
1420 Listing 6: Formatting templates for Code-Testing Flow  {{expected_output}}   1469
1421 (Codeforces)                                       ```````````````````` 1470
1422 "formatting templates":                            ## [Failed test] Generated 1471
1423 "no error template": |-                            output                 1472
1424     ${.issue_title}                                ```````````````````` 1473
1425     All of the executed tests                       {{generated_output}}   1474
1426     passed.                                         ```````````````````` 1475
1427 "all tests header": |-                            "test error": |-      1476
1428     ${.issue_title}                                ## [Failed test] {{idx}} 1477
1429     The Python code does not                       ### [Failed test] {{idx}} 1478
1430     solve the problem in the                       Input                 1479
1431     problem description due to                     ```````````````````` 1480
1432     logical errors. It fails                        {{test_input}}        1481
1433     on the following tests.                         ```````````````````` 1482
1434 "compilation error template":                     ### [Failed test] {{idx}} 1483
1435     |-                                             Expected output       1484
1436     ${.issue_title}                                ```````````````````` 1485
1437     The execution resulted in a                     {{expected_output}}   1486
1438     compilation error.                              ```````````````````` 1487
1439     ## Compilation error message:                 ### [Failed test] {{idx}} 1488
1440     {{error_message}}                               Generated output      1489
1441     "timeout error template": |-                   ```````````````````` 1490
1442     ${.issue_title}                                {{generated_output}}   1491
1443     The execution timed out, the                     ```````````````````` 1492
1444     solution is not efficient
1445     enough.
1446 "runtime error template": |-
1447     ${.issue_title}
1448     The execution resulted in a
1449     runtime error on the
1450     following test.
1451     "prompt templates":
1452     "system_message": |-
1453     Your goal is to identify the
1454     issues with an incorrect

```

1497	competitive programming	```	1549
1498	solution attempt.		1550
1499		{{testing_results_summary}}	1551
1500	The user will specify the		1552
1501	problem by providing you		1553
1502	with:	Consider the problem	1554
1503	- the problem statement	statement, the solution	1555
1504	- input description	attempt and the issue. Why	1556
1505	- output description	is the solution attempt	1557
1506	- example test cases	incorrect? How should it	1558
1507	- (optional) explanation of	be fixed? Explain your	1559
1508	the test cases	reasoning very concisely,	1560
1509	- an incorrect Python	and do not provide code.	1561
1510	solution attempt and a	"human_message": -	1562
1511	description of its issue	{{query}}	1563
1512			
1513	Crucially, your goal is to	Listing 8: Prompts for Plan Flow (Codeforces)	
1514	consider all aspects of	"prompt templates":	1564
1515	the problem and pinpoint	"system_message": -	1565
1516	the issues with the	Your goal is to provide a	1566
1517	solution attempt, and not	high-level conceptual	1567
1518	to provide the code	solution that, if	1568
1519	implementation yourself.	implemented, will solve a	1569
1520	Some aspects to consider: Is	given competitive	1570
1521	the input correctly parsed	programming problem.	1571
1522	? Is the output correctly		1572
1523	formatted? Are the corner	The user will specify the	1573
1524	cases correctly handled?	problem by providing you	1574
1525	Is there a logical mistake	with:	1575
1526	with the algorithm itself	- the problem statement	1576
1527	?	- input description	1577
1528	Use the code execution	- output description	1578
1529	results provided in the	- example test cases	1579
1530	issue description to guide	- (optional) explanation of	1580
1531	your reasoning/debugging.	the test cases	1581
1532	"query_message": -		1582
1533	# Problem statement	The proposed algorithm should	1583
1534	{{problem_description}}	be computationally	1584
1535		efficient, logically	1585
1536	# Input description	correct and handle all	1586
1537	{{input_description}}	corner cases.	1587
1538			1588
1539	# Output description	The user will provide you	1589
1540	{{output_description}}	with a task and an output	1590
1541		format that you will	1591
1542	{{io_examples_and_explanation	strictly follow.	1592
1543	}}	"query_message": -	1593
1544		# Problem statement	1594
1545	# Solution attempt to be	{{problem_description}}	1595
1546	fixed		1596
1547	```python	# Input description	1597
1548	{{code}}	{{input_description}}	1598

```

1599
1600     # Output description
1601     {{ output_description }}
1602
1603     {{ io_examples_and_explanation
1604         }}
1605
1606
1607     Return a high-level
1608         conceptual solution that
1609         would solve the problem.
1610         Be very concise, and do
1611         not provide code.
1612     Reply in the following format
1613         :
1614     # Conceptual solution
1615     {{ plan_placeholder }}
1616     "human_message": |-
1617     {{ query }}

```

We will curate a leaderboard of best-performing Flows that will be publicly available on FlowVerse and provide the predictions that reproduce the reported scores using the provided infrastructure.

1647
1648
1649
1650

A.5 The CC-Flows-competition: a new form of competitive coding

Solving competitive coding challenges is an eminently hard problem. The solve rate of only 27% by directly attempting the problem and 47% by the best-performing code Flow, paired with a reliable automatic evaluation metric, make competitive programming an ideal benchmark for AI systems. Motivated by this, we propose a competition where instead of people, proposed Flows solve competitive programming problems.

The competition will leverage the comprehensive dataset of publicly available Codeforces problems and the open-source infrastructure for inference and testing used in the experiments, available at [anonymous](#). The competition will only include problems published after the knowledge-cutoff date of GPT-4. Furthermore, not to overload the Codeforces online evaluation infrastructure, we further filter this dataset to problems for which public and private tests are available, and the output format is compatible with our local code testing infrastructure. Codeforces ranks the difficulty of each problem from 800 to 2100. At the time of publishing, we have the following number of problems per difficulty (total of 416):

- difficulty 800: 149
- difficulty 900 to 1500 (inclusive): 185
- difficulty 1600 to 220 (inclusive): 82