

# CoPrompt: Supporting Prompt Sharing and Referring in Collaborative Natural Language Programming

Felicia Li Feng\*  
Computational Media and Arts Thrust  
The Hong Kong University of Science  
and Technology (Guangzhou)  
Guangzhou, China

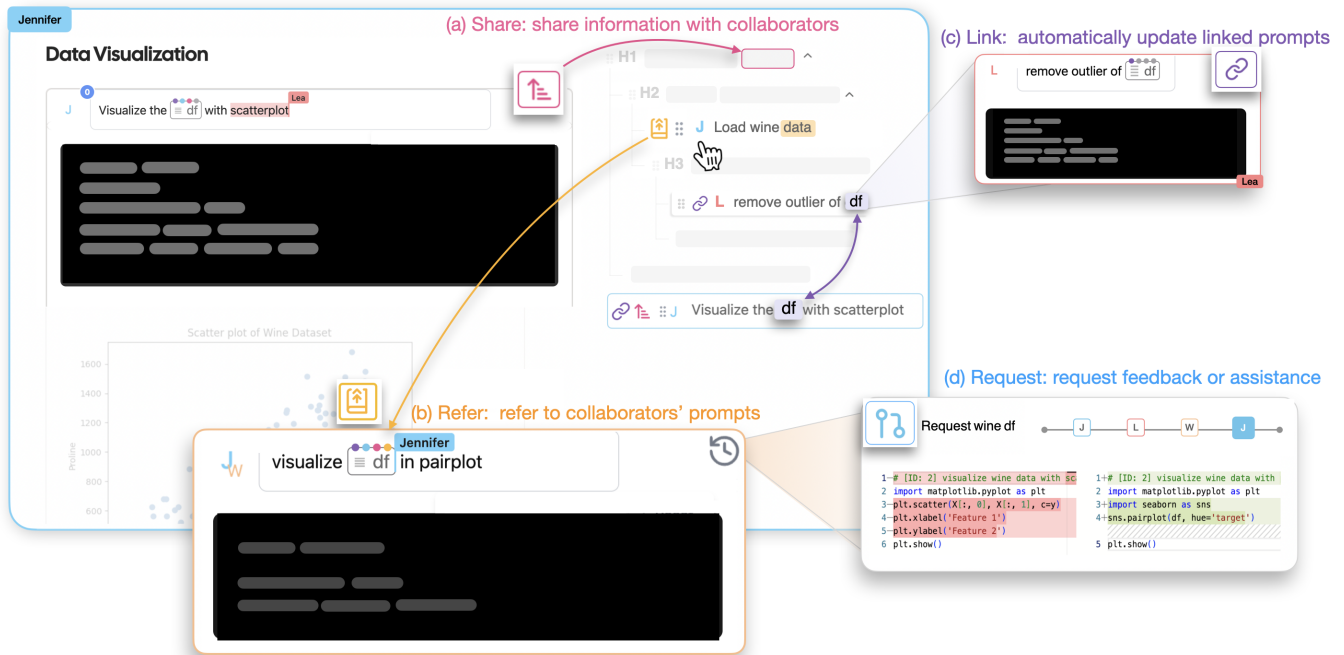
Ryan Yen\*  
School of Computer Science  
University of Waterloo  
Waterloo, Ontario, Canada

Yuzhe You  
School of Computer Science  
University of Waterloo  
Waterloo, Ontario, Canada

Mingming Fan†  
Computational Media and Arts Thrust  
The Hong Kong University of Science  
and Technology (Guangzhou)  
Guangzhou, China

Jian Zhao  
School of Computer Science  
University of Waterloo  
Waterloo, Ontario, Canada

Zhichong Lu  
Department of Computer Science  
City University of Hong Kong  
Hong Kong, China



**Figure 1:** *CoPrompt* enables programmers to conduct collaborative prompt engineering by building upon collaborators' prompts in natural language programming. It provides four mechanisms: (a) *share* mechanism enables programmers to share information with collaborators without much effort or interrupting collaborators' work. (b) *refer* mechanism assists programmers to modify prompts with reference to collaborators' prompts. (c) *link* mechanism automatically updates linked prompts. (d) *request* mechanism enables programmers to request collaborators' assistance or feedback without interrupting collaborators' progress.

## ABSTRACT

Natural language (NL) programming has become more approachable due to the powerful code-generation capability of large language models (LLMs). This shift to using NL to program enhances collaborative programming by reducing communication barriers

and context-switching among programmers from varying backgrounds. However, programmers may face challenges during prompt engineering in a collaborative setting as they need to actively keep aware of their collaborators' progress and intents. In this paper, we aim to investigate ways to assist programmers' prompt engineering in a collaborative context. We first conducted a formative study to understand the workflows and challenges of programmers when using NL for collaborative programming. Based on our findings, we implemented a prototype, *CoPrompt*, to support collaborative

\*Both authors contributed equally.

†Also with Division of Integrative Systems & Department of Computer Science and Engineering, The Hong Kong University of Science and Technology.

prompt engineering by providing referring, requesting, sharing, and linking mechanisms. Our user study indicates that *CoPrompt* assists programmers in comprehending collaborators' prompts and building on their collaborators' work, reducing repetitive updates and communication costs.

## 1 INTRODUCTION

Collaborative programming has been widely studied and supported through a range of collaborative systems [16, 17, 50, 66, 70, 84–86]. These systems assist programmers in collaboratively writing, discussing, and debugging code across various contexts, such as data science [70, 84–86] and software development [16, 17, 66]. While code has traditionally played a central role in collaborative programming, the emergence of large language models (LLMs) introduces an alternative approach: using natural language (NL) for programming and collaboration. Leveraging the context offered by NL prompts allows programmers to effectively communicate with their collaborators using a more intuitive language without getting into the intricacies of low-level code [85]. This benefit aligns with programmers' preference for understanding collaborator tasks from a high-level perspective [84].

To obtain desired code generation results, programmers often need to conduct *prompt engineering*, which involves iteratively refining prompts to guide LLMs in solving programming tasks by evaluating the generated results [18, 45, 71]. However, prompt engineering is challenging in collaborative programming. To effectively collaborate with others and ensure that their engineered prompts align with the ongoing work of their collaborators, programmers need to stay informed about their collaborators' progress. This involves regularly reviewing their collaborators' code and engaging in clear communication without causing any disruptions or difficulties for their fellow collaborators. However, programmers encounter difficulties when switching between their code and that of others [16], including issues such as a lack of contextual references [86] and limited support for interactive sharing of intermediate results [7, 47, 64]. Additionally, balancing the inclusion of contextual information in prompts can be challenging for programmers [44, 75]. Deciding on the right amount of detail to incorporate is not always straightforward, resulting in prompts that may either lack essential information (e.g., “*Web Scraping*”) or become overly detailed (e.g., “*Extract all anchor <a> tags from the parsed HTML and iterate through each*”). This variation can lead to confusion among collaborators and impact the readability and reusability of prompts. In addition, sometimes a prompt is written for better code generation, which only needs to be comprehended by LLMs instead of collaborators (e.g., “*example execution results*”). These issues increase the cognitive load of making sense of prompts and thus increase the communication cost of collaboration.

The purpose of this research is thus to explore the design of workflows that support programmers in *Prompt Co-Engineering*, which involves collaboratively refining and sensemaking prompts during NL programming. We selected data science work as a case to demonstrate the potential workflows of prompt co-engineering. Given the exploratory and explanatory nature of data science, it requires programmers to collaborate closely by sharing intermediate results [7, 47, 64] and engaging in discussions with the assistance

of contextual references [86]. However, our primary focus is to comprehend and facilitate the workflow of prompt co-engineering rather than addressing specific data science tasks.

We conducted a formative study to gain insights into potential prompt co-engineering workflows and challenges. Our findings revealed that programmers struggled to maintain a shared common ground, track collaborators' revision histories of prompts and code, and comprehensively understand the code solely from prompts due to their iterative nature. They also encountered challenges in managing procedural dependencies [48, 76] when dealing with variables represented by different NL prompts, which resulted in repetitive updates. These findings highlight the importance of supporting the prompt co-engineering workflow that encompasses comprehending collaborators' work and leveraging it to construct their own prompts.

Building on these findings, we introduce four innovative mechanisms for natural language programming aimed at reducing the effort required to build upon others' work and facilitate information sharing among collaborators: *referring*, *requesting*, *sharing*, and *linking* (Figure 1). *Referring* enables programmers to locate and access their collaborators' prompts by presenting user-defined tasks and prompts within a shared multi-level hierarchy view of prompts. *Requesting* and *sharing* enable programmers to share information with their collaborators and solicit feedback to enhance their prompts. The *linking* mechanism facilitates automatic updates for elements with procedural dependency, reducing the necessity for repetitive prompt modifications.

Incorporating these mechanisms, we designed *CoPrompt*, a prototype system that assists the workflow of *Prompt Co-Engineering*. *CoPrompt* supports programmers in making sense of collaborators' work with multi-level hierarchical interactions and contextual prompt information, as well as leveraging collaborators' work and sharing information. To evaluate the usefulness of *CoPrompt* in assisting prompt co-engineering workflow during collaborative NL programming, we conducted a 2-part user study involving 12 experienced programmers familiar with LLM-based code assistants and collaborative programming. Participants were asked to complete a real-time collaborative programming task in pairs, working together simultaneously. They were then tasked with following up on the work of other participant pairs to simulate asynchronous collaboration. This involved comprehending and modifying existing work by others and independently addressing tasks without online collaboration. The results showed that *CoPrompt* effectively supported programmers in understanding their collaborators' prompts and facilitated communication among them to build upon each other's work. In summary, this research makes the following contributions:

- A formative study that uncovered the workflow and challenges of collaborative NL programming.
- *CoPrompt*, a prototype system that supports programmers' prompt co-engineering workflow during collaborative NL programming by comprehending, referring, requesting, sharing, and linking with collaborators' prompts.
- A user study that provided insights into the usability and usefulness of *CoPrompt* and design implications for future systems assisting prompt engineering in collaborative NL programming contexts.

## 2 RELATED WORK

As our research aims to address the challenges of collaborative prompt engineering in collaborative NL programming, we review prior work on natural language programming and LLMs, as well as collaborative programming groupware.

### 2.1 Natural Language Programming and Prompt Engineering

Natural Language (NL) Programming is the process of using NL to express programming ideas for desired output [56, 57]. Prior work provided insights on how people express computer-like procedures “naturally” and on what features programming languages should include to be more “natural-like” [57]. With the development of natural language processing (NLP) [10], it has become feasible to use NL to conduct more programming tasks, as it allows more free-form NL utterances to be translated into program code [42]. This advancement increased the accessibility of programming to non-expert users [56] and end-user programmers [35] who lack training in computing.

Recently, the advances in generative AI [58, 89], especially LLMs [8], fostered the capability of generating code from NL prompts, by allowing a wider space of utterances to be transformed into satisfying code snippets. This advance significantly enhanced the performance of AI-driven code assistants [1] and thus improved the satisfaction and accessibility of programming with NL prompts [32, 75]. LLM-powered code assistants allow programmers to write at different levels of abstraction when developing code, which provides a greater degree of freedom [27, 75]. Prior work has investigated the design space of AI-powered code assistants for computational notebooks [53] and other popular code editors [81]. However, the multiple levels of abstractions [22] in the NL prompts also resulted in the abstraction matching problem when using NL for programming, where programmers find it difficult to select an utterance that will translate into the desired system action [33, 49, 67, 75]. A case study investigating the NL prompting process of prototyping also highlighted the difficulty of evaluating whether a prompt is improving [32]. This issue is rooted in the challenges of translating user instructions into executable computer tasks [30].

To mitigate the abstraction matching issue, prior work has investigated ways of prompt engineering, which is the process of engineering an NL prompt to make it more effective in generating desired results [6, 71]. Liu et al. proposed design guidelines for prompt engineering for text-to-image generative models [46]. Common practices in prompt engineering include appending information like explanations [37], demonstrations [13, 41], table schema [80], and relevant examples (few-shot prompts) [31, 45]. Although few-shot prompts have become a popular strategy, they may still behave worse than zero-shot prompts sometimes [71]. To enhance the effectiveness of the prompts, programmers can further specify tasks by constructing the signifier, memetic proxy, and specifying truth-seeking patterns [71].

Other prompt engineering methods include combining specific task information with general intentions (meta-prompts) [71], generating mutations of the prompt [40], eliciting feedback with small data [78], summarizing complicated prompts [36], defining prompt grammar [18], uncertainty highlighting [82] and introducing a new

programming language [3, 29]. Prior research into natural language interfaces suggests the benefit of managing expectations and gradually revealing the capabilities of the system through user interaction and intervention [74, 82]. There are also practices of breaking down tasks [72] and dividing complex tasks into chained series of sub-tasks [93]. By breaking down complex problems into sub-tasks, the gap in abstraction is reduced, enabling successful guidance of the model to generate code that matches the programmer’s intents [4].

However, prompt engineering in the collaborative programming context has not been investigated, which contains specific challenges regarding understanding and leveraging collaborators’ work. *CoPrompt* aims to investigate the challenges and the potential benefits the NL prompt with certain levels of abstraction can bring to the collaboration.

### 2.2 Collaborative Programming

Extensive research in HCI and CSCW has investigated challenges and system designs to assist collaborative programming. Synchronization is a challenging yet significant part of the collaboration, as programmers need to synchronize with their collaborators in various artifacts like data frames, variables, and archives [84]. It is challenging as the artifacts in programming involve procedural dependencies [48, 76]: if one part of the code changes, all related code snippets must be updated to prevent conflict and errors. In addition, programmers often encounter difficulties when switching between their code and that of others [16]. The problems of context switching and knowledge sharing are common, especially in the data science domain where programmers need to frequently share intermediate results [7, 47, 64] and discuss with the assistance of contextual references [86].

Establishing group awareness [15] can reduce communication costs and thus improve collaboration efficiency [83]. It involves understanding the activities of others, information sharing, and knowledge of group and individual contexts [23, 24, 28, 38, 59]. This is particularly important yet challenging in the domain of data science due to the diversity of artifacts and individuals involved in data science work [11]. To facilitate comprehending the complex dependencies and relationships among collaborators’ work, Albireo displays the relationships between the cells of a computational notebook using a dynamic graph structure [91]. Documentation plays an important role in maintaining shared understanding and group awareness [14, 34, 51, 73, 77]. To document the development progress, programmers write comments to make the code easier for both themselves and others to understand [63]. Comments are also essential for sharing intermediate results in data science work [65]. However, writing comments is tedious which makes many programmers not bother to write comments in time. The lack of detailed explanations and intention-revealing comments causes trouble for others understanding their work [19]. To make the commenting process easier, Wang et al. built Themisto, which leveraged AI to provide AI-assisted comments based on deep learning, query, and prompt [85]. Their user study suggested that the collaboration between data scientists and Themisto significantly reduced task completion time and resulted in satisfaction.

Reusing collaborators’ work is also challenging in collaborative programming [43, 75, 87]. To facilitate referring to collaborators’

work, chat.codes enabled programmers to link code with messages in the chatroom [61]. In addition, Codeon provides on-demand remote collaboration assistance by automatically capturing the relevant code context and allows remote helpers to respond with high-level descriptions, code snippets, and NL explanations.[9]. Communication is essential for maintaining shared understanding and group awareness in collaborative work [23]. While it could be time-consuming in software developing collaboration [28], prior work has investigated ways of reducing collaborators’ communication costs through documentation [19], comments [63], visualizations, and version control systems [9]. Considering the heavy dependencies among the artifacts involved in data science work, code-gathering tools highlight dependencies used to compute results to assist programmers to understand, reuse, and rewrite in cluttered notebooks [26].

However, these collaboration systems have not taken NL programming into consideration, where the challenges of comprehending and leveraging collaborators’ work may be different. *CoPrompt* aims to investigate the challenges and potential solutions for the challenges in collaborative NL programming, especially prompt co-engineering.

### 3 FORMATIVE STUDY

We conducted a formative study to understand the challenges faced by programmers and their needs in the workflow of prompt co-engineering. Specifically, we focused on how programmers comprehend and build upon their collaborators’ work to iteratively refine their own prompts for generating code that accurately matches their intents.

#### 3.1 Participants and procedure

Five pairs of experienced programmers familiar with LLM-based code assistants were recruited. Participants reported 6–24 months of experience with an AI code assistant and 3–5 years of experience using computational notebooks. All participants were 20–35 years old and had at least bachelor’s degrees in a CS-related field.

We asked participants to work remotely in pairs on an exploratory data programming task using a shared Jupyter notebook in the VS-Code editor [54] embedded with the GitHub CoPilot plugin. The real-time sharing functions are provided by the *Live Share* plugin [55], which synchronizes edits between users and allows collaborators to see each other’s cursors. The data programming task [2] is a popular data science task on the Kaggle platform which requires participants to use advanced regression techniques to conduct a prediction. It involves common data science operations such as data cleaning, encoding, feature transformation, and correlation analysis.

Participants were asked to join a Zoom meeting first to discuss their task distribution and collaboration workflow. Then, they started working on their own tasks using NL prompts, and they were allowed to communicate via audio in the meantime. To observe the natural prompt co-engineering workflow, participants were explicitly required to modify the NL prompts instead of directly tweaking the code. While the study session was conducted in real-time, the nature of the tasks did not require synchronous collaboration, i.e., participants had the flexibility to divide the tasks

into sub-tasks and work on them asynchronously. We did not explicitly require participants to complete the tasks synchronously or asynchronously. The collaborative programming session lasted about 90 minutes, after which participants were asked to attend 30-minute follow-up interviews. We instructed participants to try their best to complete the tasks in high quality and efficiency. While there was no external incentive for the task performance, all participants successfully completed the tasks. All participants received compensation according to local standards.

#### 3.2 Data Analysis

The studies were logged using VSCode extensions and the process was video-recorded and transcribed. Two co-authors conducted an inductive thematic analysis [5] involving cross-referencing timestamped data of prompt modifications from system logs with video recordings and interview transcripts, identifying the events that transpired before prompt engineering. Specifically, the two co-authors read through the transcripts first to familiarize themselves with the data and then performed the open coding process independently. Then, all co-authors discussed and updated the code book during the weekly project meeting for two weeks. Finally, we categorized and analyzed a total of 229 instances out of 392 recorded interactions between participants, which fall into 3 stages in the prompt co-engineering workflow: comprehension, pre-modification interactions, and prompt modification. We excluded 163 actions due to a lack of clear context or relevance to the specific collaborative communication events that preceded prompt modifications.

#### 3.3 General Workflows in Prompt Co-Engineering

In the following paragraphs, we present our findings of the workflows that participants adopted in prompt co-engineering and the challenges (C) that they encountered.

In the initial stages, participants convened online to gain an understanding of the data. Subsequently, high-level task distribution was discussed and noted down at the beginning of a collaborative computational notebook. We noticed that 4 pairs structured the task by hierarchically numbering lists or bullet points to externalize the task structure in their minds. Upon settling on a preliminary task distribution, participants started writing prompts independently to accomplish their own distributed tasks. Throughout the process, they maintained communication via Zoom to discuss ongoing and potential code implementations. We observed that all participants started with a high-level description of the task (e.g., data preprocessing, encoding) with the methods involved (e.g., low-pass filtering, linear regression), deliberately omitting details (e.g., variable names and parameters). Participants mentioned the reason is that they have to “*wait until her [collaborator’s] work is complete*” -P4 to continue adding more details to the prompt. Therefore, participants do not verify the generated code in detail at first because they know that “*it will eventually be changed later*” -P2. To further refine their prompts, participants primarily go through these three stages:

**Stage 1 — Comprehension:** Regularly checking in on collaborators’ work became a common practice. The goal was either to reuse collaborators’ work or to assess their progress and determine the

next steps. However, significant time was spent “scrolling up and down to identify changes” -P9. Participants also encountered difficulties in comprehending their collaborators’ code based solely on the prompt, often describing the prompts as “unorganized” and “vague”. Additionally, P7 highlighted another issue where the generated code “sometimes not aligned with the prompt,” further complicating the comprehension process. Lastly, participants faced challenges in tracking their collaborators’ revision history of prompt and code, which is essential for understanding “the reasoning behind code changes” -P1 (C1).

**Stage 2 – Pre-Modification Interactions:** Programmers often adapt their prompts based on their own experience and the current work, which can be challenging to document comprehensively. We thus focus on the explicit collaborative strategies employed prior to the start of prompt engineering. Based on the thematic analysis results, these pre-modification interactions consisted of a series of actions (Table 1).

**Syncing up with collaborators.** Many participants found that their initial task distribution was not detailed enough, which caused redundant effort and inappropriate workflow between collaborators: “We encoded the data at the same time” -P3. Participants also reported the tedious process of updating prompts due to procedural dependency [47], in which a downstream prompt only works if a particular upstream prompt works normally. Due to the ever-changing nature of data programming work, programmers often need to monitor their collaborators’ changes and update their prompts to align with their collaborators’ process, otherwise, they may “encounter error messages due to collaborators’ modifying the data frame halfway through the process” -P2 (C2).

**Requesting for collaborators’ feedback and assistance.** All participants in charge of visualizing correlation (P1, 4, 5, 8, 9) left comments asking for feedback, as it is an essential step for data analysis. We also observed that some participants (P3, 5, 8, 9) requested help from their collaborators. For instance, P3 asked his collaborator to handle a sub-task that he failed to complete. There are also cases where collaborators need to work closely and go through a trial-and-error process together: “I asked my collaborator to pay attention to the outliers every time the way of feature transformation is changed” -P4. However, most participants (N=7) indicated that they desire a non-interruptive method to send their collaborators requests, instead of speaking up, which is too interruptive for them to use frequently (C3).

**Referring to collaborators’ processes and prompts.** All participants checked their collaborators’ processes and referred to their prompts to improve their own for better generation results across the whole notebook. To leverage others’ prompts, participants first locate and read the target prompt to make sense of it. Then, they copy portions of the prompt relevant to their task and integrate them into their own prompts to provide contextual information. The redundant copy-pasting and modifying can be time-consuming (C4), as programmers may trial-and-error to determine the appropriate modifications of the prompts (P1: “I reused my collaborator’s prompt, which did not work as I imagined. After analyzing its context, I realized that I had to copy a prompt several blocks above that”).

**Proactive communication for sharing intermediate results and relevant information.** Many participants have shared intermediate results with their collaborators that they believe would be useful. They performed three types of sharing strategies: (1) leaving comments under the block that their collaborators were working on to attract their attention - P3; (2) leaving comments before the block of the shared information and pinning their collaborators using an “@” - P5, 6; and (3) ask their collaborators to check their current highlights block for reference - P1, 2, 9. The first strategy requires the comment receiver to locate the shared information, while the second strategy may influence the collaboration efficiency. Though many participants communicated directly through Zoom, it “disturbed my own progress a bit” -P10. These strategies are either “inefficient” -P9 or “disruptive” -P7 (C4).

**Stage 3 – Prompt Modification & Merge Conflicts.** The third stage centers on modifying (i.e., engineering) the prompt. In this stage, participants refine their prompts by copying and pasting utterances or code snippets from collaborators’ prompts to clarify details about the variable name, resource, methods, and detailed considerations. During this phase, participants may encounter merge conflicts or issues that need communication for resolution. Participants also expressed a desire to access previous versions of the code, as this helps them “recall who made specific changes to the prompt” -P4 and the “reasons behind those alterations” -P7. Some participants manually tweaked the code instead of modifying the NL prompts when they could not achieve their desired results within a few attempts: “find it more efficient to directly change code after several failed trials” -P1.

In summary, we identified the following user challenges in the formative study:

- C1: Effort of maintaining group awareness and shared understanding to enhance collaboration effectiveness.
- C2: Repetitive effort of syncing with collaborators’ work.
- C3: Inconvenient and disruptive ways of requesting collaborators’ feedback and assistance.
- C4: Repetitive copy-pasting effort for leveraging others’ work and disruptive information sharing.

## 4 DESIGN CONSIDERATIONS

Based on the findings from the formative study, we formulated four Design Considerations (Ds), to support prompt co-engineering in collaborative NL programming.

**D1: Supporting sense-making of collaborators’ progress and prompts.** Programmers encountered challenges locating collaborators’ work in a shared notebook that lacked a clear outline of NL prompts and code (C1). To support programmers’ locating and sense-making of collaborators’ progress, it is important to implement a clearer structure to show the notebook overview [18]. This structure should include multiple levels of hierarchy (e.g., tasks, sub-tasks, and prompts) to help programmers identify changes easily. The design should also incorporate assistive features to help programmers understand the code from prompts that might be too vague. In addition, a history view should be provided for programmers to track global activities (e.g., collaborators’ works) and local changes (e.g., prompts variations).

Action	Description	Example	<i>n</i>
Sync Up	Participants communicated to align their efforts and ensure consistency in their tasks. These interactions helped prevent redundancy and maintain cohesiveness in their prompt engineering process.	"I am encoding the whole dataset."	49
Reactive Communication (Request & Response)	These interactions aim to seek assistance, feedback, or validation regarding reactively handling specific aspects of the prompt or generated code.	"Can someone help me determine the function for encoding?"	61
Clarifying	Participants seek answers to queries about their collaborators' work, including seeking explanations, and verifying the correctness of specific code segments.	"Did you drop the column of xxx?"	57
Reference & Reuse	Participants occasionally referred to and reused (e.g., copy-paste) components from their collaborators' work, utilizing these references to inform their own prompt modifications. These actions fostered a sense of collaboration and knowledge exchange.	"I used the function you wrote in my block. Any concerns?"	33
Proactive Communication	Proactive communication involves participants sharing insights, updates, or information related to their prompts or coding tasks. These exchanges often contributed to a deeper understanding of the prompt's context.	"Here's an update on the changes I made to the prompt..."	29

**Table 1: Five types of actions of the pre-modification interactions.**

**D2: Automatic synchronization to reduce repetitive updates.** In the exploratory and iterative programming process, the prompts and codes might be updated several times throughout the whole process based on the collaborators' changes (C2). Automatic synchronization of variables and code snippets with procedural dependency should be provided to reduce programmers' cognitive load and enhance their collaboration efficiency. Programmers should be able to easily notice changes made by collaborators and the automatic updates applied to their own work.

**D3: Supporting requests for feedback and assistance.** Current ways of requesting collaborators' feedback and assistance are inconvenient and disruptive (C3). An efficient way of requesting feedback should be provided besides communicating through chat and voice. In addition, considering the situation that the collaborator has not finished the required prompt, programmers should be equipped with methods to request knowledge from others and refer to it later.

**D4: Reduce effort for sharing knowledge and incorporating others' work.** Current ways of referring to collaborators' prompts and reuse are tedious and time-consuming (C4). Programmers need to copy, paste, and modify, which takes a lot of time unnecessarily. A more effortless way of referring to collaborators' prompts should be provided. Additionally, the design should enable proactive sharing of intermediate results with collaborators, promoting sharing without concerns about disrupting their workflow.

## 5 ENVISIONED SCENARIO

Here, we present a motivating scenario that illustrates the workflow of using *CoPrompt* for prompt co-engineering in NL programming. For simplicity, we describe our scenario using two collaborators, Alice and Bob.

Alice and Bob are remotely collaborating on a data analysis task requiring them to predict house prices. To improve collaboration efficiency, they divide the tasks into smaller segments, allowing each to tackle different tasks separately. They decide that Alice would handle missing values, outliers and categorical features. Meanwhile, Bob is tasked with feature transformation and correlation analysis (Figure 2). To track each other's progress and offer/request help for specific tasks, they outline all sub-tasks in *CoPrompt*'s rich text editor (Figure 3 a), which is synchronously displayed in the multi-hierarchical wiki (Figure 3 b). The wiki's foldable task items provide them with a clear overview of the collaboration process. After listing all the required tasks, They begin to work on writing prompts independently.

Although Alice plans to finish encoding before Bob begins feature transformation, her progress is delayed due to technical issues. Due to time constraints, Bob cannot wait for Alice to finish encoding. With *CoPrompt*, Bob creates a `request` (Figure 5) from the `@` in the transformation prompt to Alice's encoding prompt, indicating that his feature transformation steps (e.g., calculating the skewness) require Alice's encoded result. Once *CoPrompt* detects (through semantic analysis) that Alice has completed the encoding, it automatically updates Bob's prompt to leverage encoded data to generate code. As a result, Bob no longer needs to manually modify his prompt whenever Alice updates her encoding, saving his time and allowing him to focus more on task completion rather than repetitive code upkeep.

After completing the feature transformation, Bob checks the progress from the wiki (Figure 3) and notices that Alice needs the transformed data for later outlier handling. In case Alice needs to spend much time verifying the data to be used, Bob decides to proactively share it with Alice. To share this transformed data, Bob

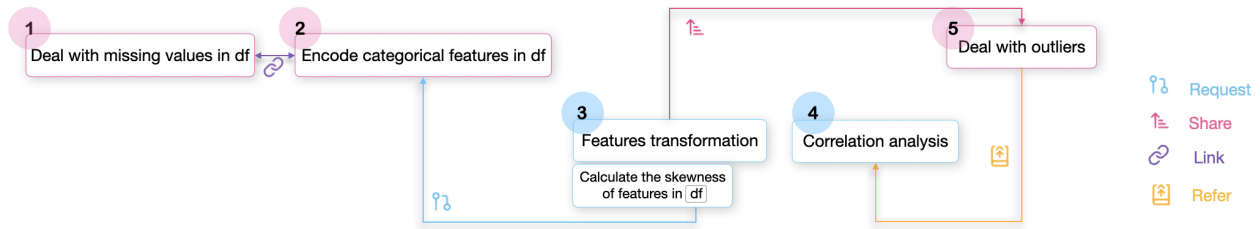


Figure 2: Envisioned Scenario of collaborative NL programming using *CoPrompt*, including tasks from 1 to 5 (pink boxes indicate Alice’s tasks, and blue boxes indicate Bob’s tasks). Four colors of arrows indicate four types of mechanisms.

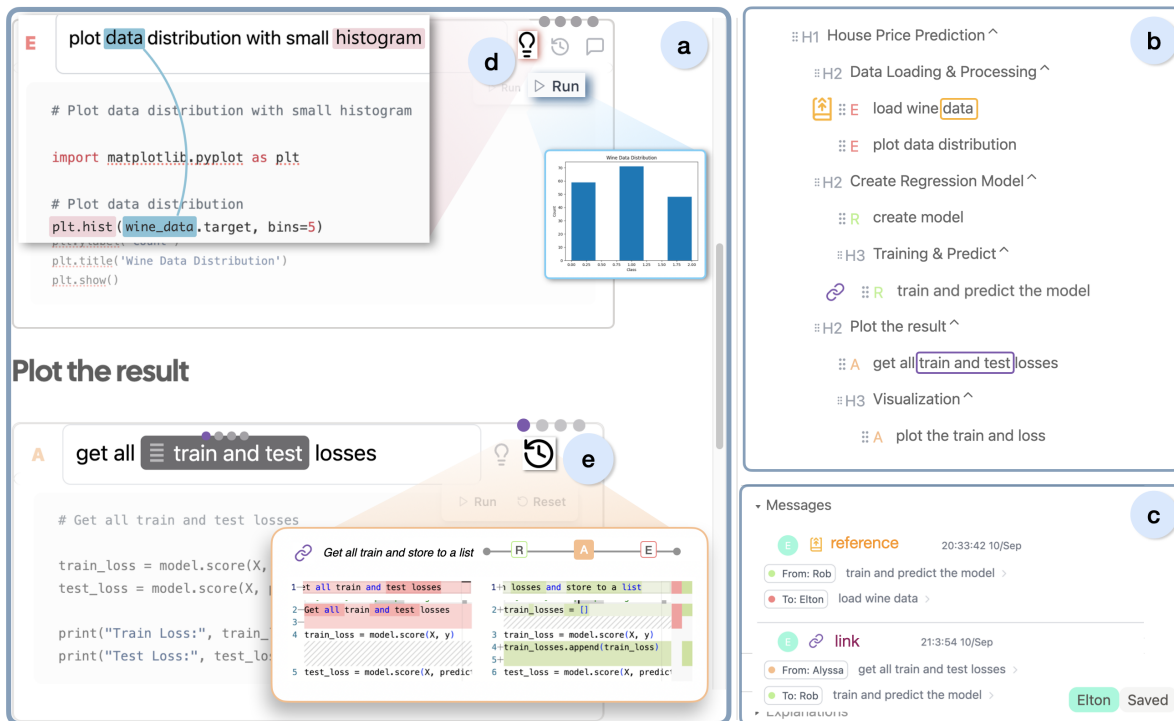


Figure 3: The *CoPrompt* user interface includes (a) a block-based rich text editor for NL inputs, which consists of prompt blocks and execution code blocks. (b) Prompt Wiki: a multi-hierarchy wiki displaying tasks and prompts, (c) a message panel providing a comprehensive log of actions, (d) an explanation view displaying an explanation for prompts and code-prompt relationship, and (e) a history view for version control.

utilizes the **share** mechanism (Figure 6) by highlighting his data frame and clicking on the share icon next to Alice’s task in the wiki. Alice receives a pop-up, allowing her to accept Bob’s data without manually locating the data frame for her subsequent workflow. Once accepted, *CoPrompt* automatically regenerates Alice’s prompts based on Bob’s shared data frame, eliminating the need for Alice to input supplementary information like variable names herself. This reduces the risk of human errors such as typos or incorrect variable references.

After handling missing values and encoding, Alice anticipates that there may be future adjustments to the missing value handling based on her past experience, which may also require updates to

her encoding methods. To avoid repetitive updating, Alice creates a **link** (Figure 7) connecting the `df` in the missing value handling prompt and the `df` in the encoding prompt. With this link, the encoding prompt gets automatically updated whenever the prompt for missing value handling changes (e.g., when `df`’s variable name or the method for handling missing values changes). This way, Alice avoids the need to repeatedly update the encoding prompt. If Alice no longer wants the two nodes to be automatically synced, she can unlink the nodes by de-highlighting the link icon associated with the nodes.

When handling outliers, Alice needs to refer to the results of Bob’s correlation analysis. However, she finds it challenging to

navigate through Bob’s prompts, which contain long execution examples and demonstrative code. To help her better understand Bob’s prompts, Alice expands the explanation view (Figure 3d) to see the highlighted prompt and annotated relationships between the NL prompts and code snippets. From the explanation view, Alice understands Bob’s considerations for correlation analysis and appropriate criteria for determining outliers. Then, Alice begins writing her prompts to address the outliers. To ensure no information (e.g., criteria for determining outliers) is overlooked, Alice wants to instruct the LLM to determine the outlier handling method based on the results of the correlation analysis. Alice employs the **refer** mechanism (Figure 4), creating a node that links the outlier handling prompt to the correlation analysis prompt. As a result, *CoPrompt* updates Alice’s prompts with appropriate methods, no longer requiring Alice to modify her prompts by copy-pasting and typing.

In asynchronous collaboration settings, the request, link and refer mechanisms work similarly as those in synchronous settings because these mechanisms do not require the collaborators to respond in real-time. However, there are some differences for the share mechanism. Specifically, if Alice is offline when Bob is sharing artifacts, the message cannot be delivered in real-time. Instead, *CoPrompt* retains the information and systematically presents each shared artifact when Alice reconnects online. Notably, any modifications made since the programmer’s last entry into the editor are highlighted to facilitate convenient inspection. Additionally, the usage details for each mechanism are readily accessible through the message panel.

## 6 DESIGNING COPROMPT

Based on our design considerations, we developed a prototype, *CoPrompt*, to support programmers in their prompt engineering workflow during collaborative NL programming. Specifically, *CoPrompt* supports sense-making of the collaboration process and prompt co-engineering. *CoPrompt* interface consists of five components (Figure 3): block-based rich text editor, prompt wiki, message panel, explanation view and history view. *CoPrompt* is also designed with a set of real-time collaborative features, such as real-time displays of collaborators’ cursor location and their text selections.

### 6.1 Sense-Making and Tracking of the Collaboration Process

*CoPrompt* introduces two custom block types: the **Prompt Block**, for creating prompts that generate code segments, and the **Execution Code Block**, which contains generated code that can be compiled and executed to check the interim results. As changes are made in the text editor, the wiki view automatically updates to display the editor’s structure with a tree-based representation. This structure provides a clear overview of the editor, allowing programmers to visualize multiple levels of hierarchy, from headings to prompts, down to individual *nodes* (i.e., phrases within the prompts). The wiki allows intuitive navigation of its content by enabling programmers to click and fold each task item, collapsing unrelated tasks and concentrating on those of interest. Furthermore, it provides a clear overview of the hierarchical relationships within

the project and keeps programmers updated about modifications made by their collaborators (**D1**).

**Message Panel.** The message panel (Figure 3 c) offers a comprehensive log of actions, allowing programmers to track their collaborative activities. These actions are tied to specific elements within the editor, such as prompts and nodes, which can be easily accessed by selecting them directly from the messages. Additionally, programmers can quickly identify essential information related to each action, including its type, creator, and timestamp. Messages are displayed in chronological order, with unprocessed actions prioritized at the top and highlighted by a small dot. Once addressed, *CoPrompt* automatically updates the panel by removing the highlights to indicate that the action was resolved.

**History View.** The history view (Figure 3 e) enhances version control and tracking of changes for prompts and their associated code. Programmers can review previous versions of selected prompts, organized chronologically. This view tracks the evolution of prompts, providing details such as actions that led to changes, individuals responsible for modifications, and alterations made to the prompts. Similar to Git version control, *CoPrompt* offers a diff-view that allows programmers to identify differences between current and previous versions of both prompts and code segments.

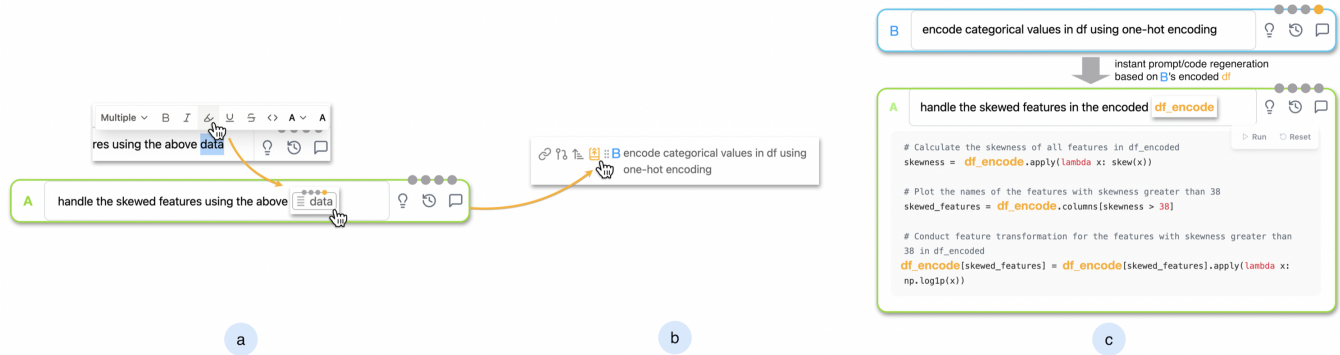
### 6.2 Supporting Programmers’ Prompt Co-engineering

To support the prompt co-engineering workflow, *CoPrompt* provides functions commonly used in collaboration interfaces (e.g., annotation and history view) to help programmers make sense of collaborators’ prompts. In addition, *CoPrompt* provides four mechanisms for programmers to (1) **refer** to collaborators’ prompts; (2) **request** intermediate results from collaborators; (3) **share** information with collaborators; and (4) **link** variables for synchronization.

**6.2.1 Making sense of collaborators’ prompts.** In order to leverage collaborators’ prompts, programmers first need to make sense of them. To check the last person who modified the prompt, programmers can check the left of each prompt where the prompt author’s icon is displayed. To assist programmers in comprehending high-level or more complex prompts from collaborators, we designed the *explanations view* (**D2**). By clicking the *Explain* button (Figure 3 d) for the prompt block, programmers can access semantic highlighting to better understand the code structure. This feature visually highlights key phrases in both the prompt and the corresponding code segments, linking them to help programmers understand the relationships between them (i.e., which phrase in the prompts led to the generation of a certain line of code). The view also provides a high-level overview of all the steps within the code, allowing programmers to quickly understand the code’s structure, logic, and functionality. To track the evolution of prompts and their previous versions, programmers can use *history view* to find historical versions of both the prompts and the corresponding code.

In the following sections, we present the design of the four core mechanisms of *CoPrompt*. All four mechanisms contain three components: (1) a *source node* in the prompt blocks that are presented in the editor; (2) a *target node* in the prompt that is listed in the *Prompt Wiki* and (3) additional messages.





**Figure 4: Workflow of refer: (a) select source node from the editor; (b) select the target node for reference from the wiki and click the refer icon; (c) A's block will instantly regenerate prompt and code based on the code and execution result of B's prompt.**

**6.2.2 Refer to Collaborators' Prompts for Precise Code Generation.** When programmers need to build upon their collaborators' work, they can use the refer mechanism (Figure 4) to specify which part of collaborators' work should be used to provide the context of their prompt (D3). The programmer needs to first select a *source node* from the editor, which could be a word (e.g., a variable name), a phrase (e.g., perform a function with a variable), or the whole prompt. Upon selection, the editor will toggle the highlighted mark on the selected node. Next, the programmer chooses a *target node* from the list of prompts displayed in the *Prompt Wiki*. This target node could also be a word, a phrase, or the entire prompt. Once selected, the *Prompt Wiki* panel will display the corresponding colored icon based on the current actions associated with the node. Additionally, a colored dot will be added before the prompt, signifying that this prompt contains nodes with associated actions. Programmers may also opt to include messages to clarify their intentions, thereby improving the accuracy of code generation.

**6.2.3 Request Collaborators' Assistance for Prompt Engineering.** When programmers require information from their collaborators' incomplete tasks as context for their own prompt, they can follow a simple process with the *Request* feature. First, they create a prompt with a placeholder indicating the expected result from their collaborators. Next, they select a *source node* from the placeholder in the editor and choose a *target node* from the collaborator's work listed in the *Prompt Wiki*. This action notifies the creator of the *target node* that a collaborator is awaiting the task to be resolved. Subsequently, the programmer can continue working on other tasks and return to check on this specific task once the collaborator has finished the requested task (D3). After creating the request, *CoPrompt* will log the actions and maintain a list of unresolved actions in the cache. When collaborators create prompt blocks that address the actions in the cache, *CoPrompt* will automatically execute the corresponding task (D2).

**6.2.4 Share Context to Collaborators.** To share intermediate results or other contexts with collaborators, programmers can employ the share mechanism (Figure 6). After choosing a *source node* in the

editor and selecting a *target node* in the *Prompt Wiki*, the collaborator will receive a pop-up notification indicating that shared context is available. Upon accepting the context, the code associated with the prompt that contains the *target node* will be updated with the context and information provided by the source node. Programmers can also select the *target node* as one of the headings that correspond to a task, especially when there are no existing suitable prompt blocks under that task. In such cases, the recipients can assign this contextual information to the prompt block they create afterward.

**6.2.5 Link Elements for Automatic Synchronization.** In a programming context, a variable, prompt, or code segment within a project may undergo multiple changes and could have various names in blocks authored by different programmers. To reduce repetitive updates, programmers can link together variables with the same value but different names. By doing so, when the prompt linked to one variable is modified, the associated prompt will automatically update based on the changes. This mechanism also extends to variables with procedural dependencies. It not only streamlines future code generation but also helps prevent potential conflicts (D4).

## 6.3 System Implementation

*CoPrompt* is a web-based application built using Next.js and React TypeScript. The core of the block-based editor is constructed using TipTap [79], which serves as a headless wrapper for ProseMirror [25], providing the foundation for a rich text WYSIWYG editor. The overall architecture of *CoPrompt* is illustrated in Fig. 11.

*CoPrompt* enables programmers to begin by creating a prompt block with generated code beneath it. Programmers can also initiate a code block and manually input code, pausing midway to generate a prompt block from the code comments. *CoPrompt* implements collaborative editing by leveraging the power of Y.js, which is a CRDT-based approach to handling shared document editing. Changes made by users are distributed and merged using WebSocket communication, with the Hocuspocus Server serving as the backend for handling real-time synchronization of documents among users. This approach enables real-time collaboration, syncing between

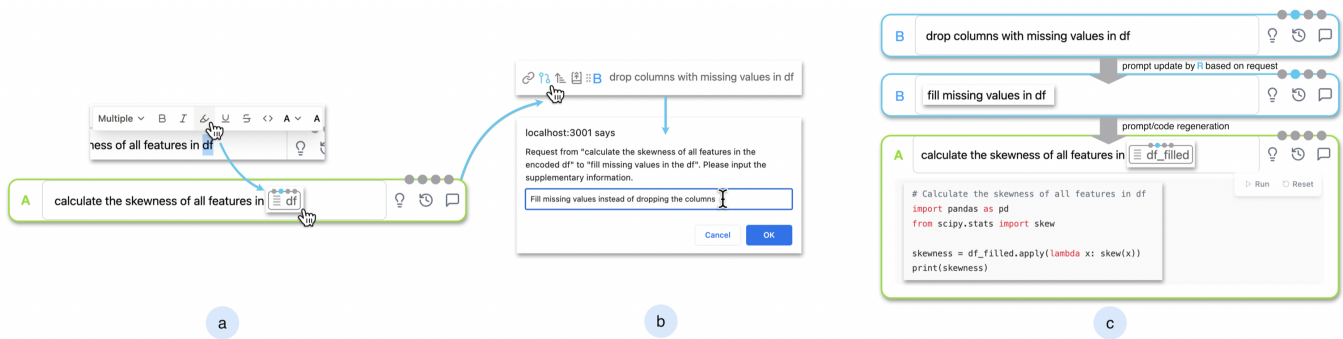


Figure 5: Workflow of request: (a) select a source node from the editor; (b) select a target node that needs the collaborator to finish and share the result from wiki and click the request icon, fill in a descriptive message for the request; (c) when B updates how they deal with missing values, A’s block will regenerate prompt and code accordingly.

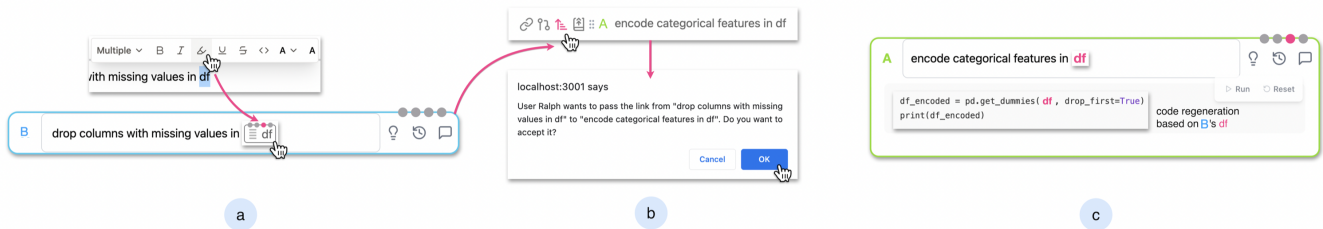


Figure 6: Workflow of share: (a) select a source node to be shared from the editor; (b) select a target node to be updated with the shared content from the wiki and click the share icon to highlight the node, then A will see a pop-up message, indicating that B would like to share some information with A; (c) when A accepts, A’s highlighted block will update its code based on B’s df.

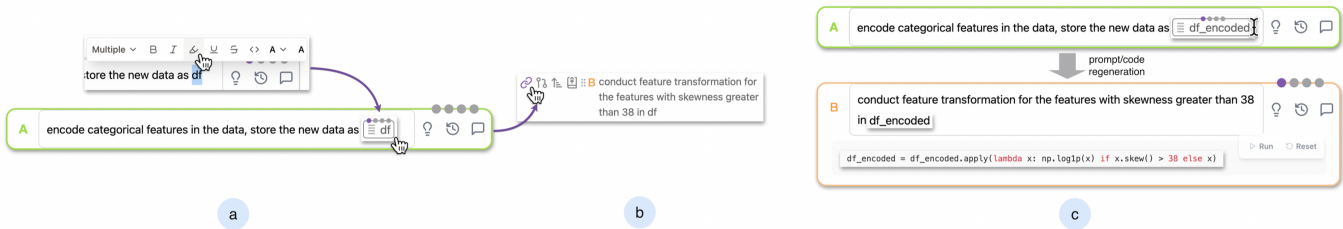


Figure 7: Workflow of link: (a) select a source node from the editor; (b) select a target node from the wiki that is to be synced with the source node and click the link icon; (c) when A updates the name of the df, B’s block will regenerate prompt/code accordingly.

devices, and the ability to work offline while maintaining consistency in the edited documents. Actions and messages are shared in real-time through the Firebase Real-Time Database [21]. The event listener updates the prompt wiki whenever actions are triggered within the four mechanisms.

To facilitate code generation, *CoPrompt* employs the OpenAI GPT-4 API [62] in combination with custom prompt templates (all prompt templates are provided in Appendix A.2). The code generation process leverages the context provided by the prompt wiki, as well as the current prompt and code blocks within block-based text editors. While *CoPrompt* allows programmers to structure their

prompts in any format using any techniques, it is not specifically designed for these techniques (e.g., prompt decomposition or few-shot learning). The primary intention is to enable them to write simple prompts and allow *CoPrompt* to generate the desired code. In the case of link, *CoPrompt* incorporates automatic self-check mechanisms with a specific prompt template to validate whether changes are necessary when one side is updated. Regarding the request, all programmers’ requests are queued and cross-verified against the expected results specified by requestors. These two mechanisms make use of prompting techniques that draw upon

few-shot learning [6] and Chain-of-Thought techniques [90] to enhance accuracy in the code generation process.

The Python code execution in the web app is made possible through Pyodide [69]. Pyodide represents a port of CPython to WebAssembly, enabling the installation and execution of Python packages directly within the browser using micropip. For the execution of Python code in a separate thread, a communication channel is established between the main thread and the Pyodide worker, incorporating a defined communication protocol. It is worth noting that, despite sharing the same context within the collaborative editor, *CoPrompt* ensures that the OpenAI and Python kernels do not overlap, preserving stability and functionality.

## 7 USER STUDY

We conducted a user study to evaluate the effectiveness of *CoPrompt* in assisting programmers' prompt co-engineering during NL programming by answering the following research questions:

- **RQ1:** How does the system support programmers to understand collaborators' progress and prompts?
- **RQ2:** How does the system support programmers to build on top of collaborators' work?
- **RQ3:** How does the system minimize the redundant updates of prompts or code?

### 7.1 Participants and Tasks

We recruited 12 participants (7 female, 5 male, aged 20-31) from a local university. All participants have more than two months of experience in AI-based code assistants and more than three years of experience in using computational notebooks. Participants were compensated \$50 for the 120-minute study. The data science task was modified from a Kaggle competition [12] that predicts survival from a disaster. To scope the task within the study duration, we asked participants to only perform exploratory data analysis. We also divided the high-level task into several sub-tasks with the assistance of two expert data scientists and provided a basic task division plan for participants' reference.

### 7.2 Procedure

Participants were first informed of the aim of this study and gave their consent. Then, they were asked to participate in a two-part study including (1) real-time collaboration and (2) following up with the others' work.

**7.2.1 Part 1: Real-time collaboration (90 minutes).** Participants were asked to perform a data analysis task in pairs using NL prompts collaboratively. Each pair of participants was asked to complete two sessions of collaborative programming: one session to use our *CoPrompt* prototype and another to use the baseline system (VS-Code live share with CoPilot plugin). The sequence of the sessions was counter-balanced.

For each session, the experimenter first introduced features of the system and gave each participant a 10-minute training session on the system, with example tasks to complete. Then, each pair of participants was asked to work on a data science task collaboratively by prompting for around 30 minutes. After that, participants were asked to rate their experience of *CoPrompt* and the baseline system

on a 7-point Likert scale. The experimenters then conducted a semi-structured interview based on the results and observed use patterns to learn participants' perspectives.

**7.2.2 Part 2: Following up with the Collaboration Process (30 minutes).** In stage 2, we evaluated how a new collaborator followed up with an ongoing collaborative project using *CoPrompt*. We asked participants to review the work of another pair of participants collected from part 1 of the study using *CoPrompt* (e.g., P1 and P2 reviewed the work of the second pair - P3 and P4). After reviewing the work for 10 minutes, we asked participants to complete additional tasks to modify the existing work, such as changing the way of dealing with outliers. The experimenters then conducted a semi-structured interview regarding the user experience of reviewing and revising others' work using *CoPrompt*.

### 7.3 Data Analysis

All study sessions were recorded and transcribed. Data collection included server-side logs, screen recordings, and interviews. Additionally, we made observational notes during the study. Our analytical approach involved the articulation of codes and themes, employing a combination of inductive and deductive thematic analysis. Two authors independently coded the transcripts and identified themes to gain insights into how participants utilized and evaluated *CoPrompt* and the baseline condition. The themes generated encompass both parts one and two of the studies, which address the three RQs. We explicitly specified if the results pertain exclusively to part two in the results section.

For all survey data, we opted for non-parametric statistical methods given the ordinal nature of Likert-scale responses and the small sample size. Specifically, we employed the Wilcoxon signed-rank test to compare responses between the two conditions. Additionally, prompts and code collected from participants' logs underwent open coding based on the reasons for manual modifications. Two researchers independently open-coded 30% of the data to establish a codebook, identifying major reasons for prompt and code modifications. These codes were then applied to the remaining data, resulting in a 74% agreement, which was subsequently refined iteratively until reaching 100%. Results from self-defined Likert scale data will be highlighted with question numbers (Fig. 9).

## 8 RESULTS

Here we report the findings from analyzing participants' survey responses, interview transcripts, think-aloud feedback, and system usage logs to understand 1) how *CoPrompt* support prompt co-engineering and 2) how participants perceive the utility of the four mechanisms for leveraging collaborators' work and sharing information among collaborators.

### 8.1 Overall Collaboration Behaviors and User Perceptions

**8.1.1 Completion Time.** All participants successfully completed the programming tasks and utilized all four types of mechanisms in the study. However, participants took significantly more time to complete the tasks in the baseline condition ( $M_{CoPrompt}=23.21 < M_{baseline}=26.82$  minutes,  $p=.002$ ,  $r=.57$ ). This could be attributed to

the increased need for synchronous communication in the baseline condition, whereas in *CoPrompt*, participants were able to utilize the four mechanisms, which served as “a way for quicker communication” -P3.

**8.1.2 Overall Collaborative Workflow.** In both conditions, all participants initiated their work by crafting prompts at a higher level of abstraction, such as defining the task as “data visualization,” and then iteratively refined these prompts to reach a lower level of detail, like specifying “pair plot df.” This process was facilitated by the wiki provided by *CoPrompt*, which allowed participants to easily locate the target prompt block that was available for further iteration. The four mechanisms further supported participants in collaborative efforts without requiring context switching to external communication tools. P4 explained, “Using *CoPrompt*, I do not need to wait for my collaborator to finish encoding, as I can write prompts for transformation first and then refer to my collaborator’s prompt.” In general, *CoPrompt* facilitated parallel work on programming tasks without being hindered by collaborators’ workflows compared to the baseline.

**8.1.3 System Usability & Cognitive Load.** To measure the usability of *CoPrompt*, we computed the SUS scores based on the UMUX-LITE [39]. The average SUS scores were significantly greater ( $p = 0.02$ ) for *CoPrompt* (Mdn = 90.61), compared to baseline (Mdn = 68.94). We also used NASA-TLX to measure participants’ perceptions of the cognitive workload of using the systems. Compared to baseline, *CoPrompt* had lower mental (Mdn = 3.5 < 5.5,  $p = 0.040$ ), physical (Mdn = 1.0 < 3.0,  $p = 0.0179$ ), and temporal (Mdn = 3.0 < 5.0,  $p = 0.0082$ ) demand, required less effort (Mdn = 3.0 < 5.5,  $p = 0.033$ ), and led to better performance (Mdn = 4.5 > 3.0,  $p = 0.0532$ ) and less frustration (Mdn = 1.5 < 3,  $p = 0.0187$ ). The overall perceived workload, obtained by averaging all six raw NASA-TLX scores (with the “Performance” measure inverted), was also lower for *CoPrompt* than baseline (Mdn = 2.5 < 4,  $p = 0.0532$ ).

**8.1.4 Code Edit & Prompt Edit.** We compared the overall code and prompt edit counts between the baseline and *CoPrompt* conditions (Figure 8 Left). We observed that while there is a less significant difference in adding prompts ( $Mdn_{CoPrompt}=35.0 < Mdn_{baseline}=46.5$ ,  $p=.003$ ), there are much larger significant differences in code editing ( $Mdn_{CoPrompt} = 45.5 < Mdn_{baseline}=111.5$ ,  $p=3.42 \times 10^{-8}$ ) and prompt editing ( $Mdn_{CoPrompt}=67.0 < Mdn_{baseline}=99.0$ ,  $p=6.94 \times 10^{-4}$ ). These substantial differences can be collectively attributed to the four mechanisms and are explained in Sec 8.4.

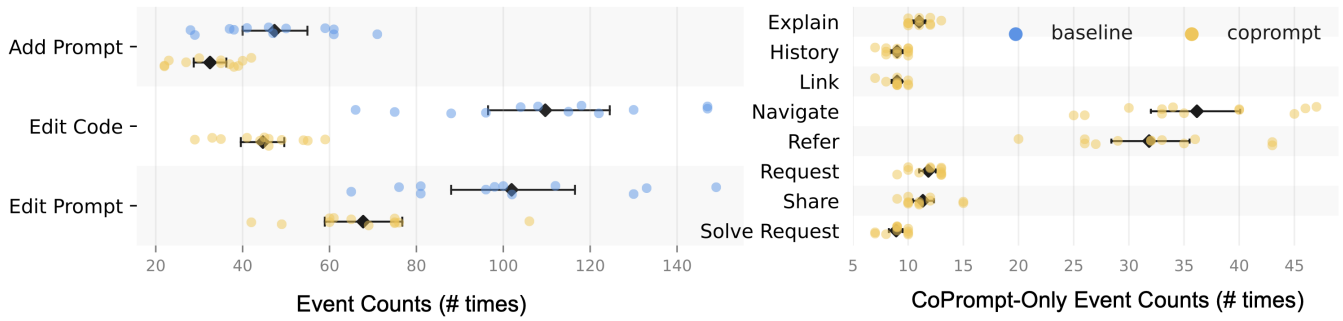
## 8.2 RQ1: How does the system support programmers to understand collaborators’ progress and prompts?

*CoPrompt* supported programmers’ sense-making of collaborators’ work by providing the hierarchical overview, generating explanations associated with the prompts, and displaying historical views. In the following sections, we report the detailed usage and perceived utility of these features.

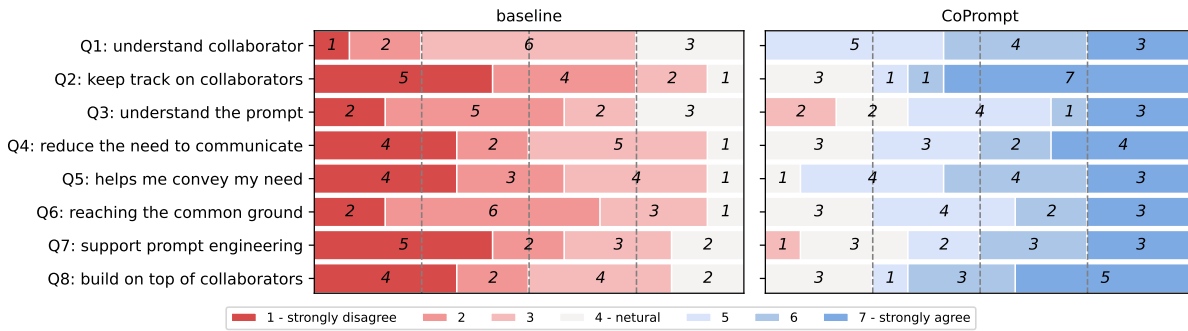
**8.2.1 The holistic overview with multi-levels of details facilitates users tracking and locating collaborators’ processes (D1).** All participants found that the prompt wiki reduces the need for programmers to keep track of collaborators’ progress (Q2: Mdn = 6.5 > 2,  $p = 0.0034$ , Figure 9) and helps them understand what the collaborators are doing (Q1: Mdn = 5.5 > 3,  $p = 0.0034$ ). The foldable table-of-content-like view allowed participants to view all the sections of the notebook at a higher level, which facilitated participants in locating their target regions and “understand the overall structure easily.” -P4 With the navigation function, participants could navigate to their target regions by clicking on the items in the prompt wiki, which is straightforward as “there is no need to scroll back and forth.” -P1 The event count from the log data (Figure 8 Right) shows that participants frequently utilize this navigation feature ( $M_{CoPrompt}=36.17$ ,  $SD=7.47$ ). P2 expressed their preference for the annotation in the wiki as it indicated the owner of the prompt: “With the icon before the prompt, the ownership of the prompt is clear.” We also noticed that most participants (N=9) collapsed their own sections while leaving the collaborator’s sections expanded to “track on collaborators’ work” -P4 and “easily identify changes not made by me [participant].” -P11 In part two of the study, participants utilize the prompt wiki to gain a quick understanding of the asynchronous work accomplished by their collaborators. They subsequently used the navigation feature to look into lower-level code details when they encountered “tasks that were unclear.” -P5

**8.2.2 Generating and associating explanations assisted participants’ comprehension of prompts (D2).** All participants found that the system helps them understand the prompt written by collaborators (Q3: Mdn = 5 > 2.5,  $p = 0.0304$ ). They mentioned that the high-level, step-by-step explanations of the generated code and the semantic connections between prompts and code segments are “especially useful when the prompt or the code is way too long.” -P1 Participants typically required these explanations before utilizing the **refer** mechanism, as sometimes they struggled to fully grasp their collaborators’ intentions through the prompts alone. The prompt explanations played a vital role in scaffolding participants’ comprehension by “providing more details about the code.” -P1 Consequently, participants could reuse or refer to their collaborators’ work without the need for direct communication, as highlighted by one participant, “I do not need to keep bothering my partner to ask what the code is about.” -P9 Furthermore, some participants perceived these explanations as a means to reduce conflicts and errors during collaboration. As one participant noted, “I became less likely to misunderstand my collaborators’ intention and modify their code, which always causes conflicts.” -P1

**8.2.3 Providing a history view allows users to have a clear view of the changes (D1).** The history view presents a comprehensive record of all historical versions of a specific prompt, offering users insight into its evolving process, as noted by P5, “It is good to view the changes.” This view not only captures the modifications made but also the “interactions that happened” -P1 and the “generated result for each iteration” -P2, thus facilitating a holistic understanding of the prompt’s evolution. Participants frequently utilized this feature when changes were enacted after any mechanism had been activated (e.g., the requested message had been resolved). In such cases, the history view allowed them to identify “who altered my code.” -P2



**Figure 8: Distribution of event counts per participant. The left image compares the event counts between the baseline and CoPrompt, and the right image shows the count of CoPrompt-only event types.**



**Figure 9: The results from the self-defined Likert scale questionnaire comparing between condition baseline and CoPrompt**

Most participants (N=8) also leveraged the diff view of the code and prompt to resolve the conflict and restore the version. Overall, participants reported that the system supports programmers in reaching common ground with collaborators (Mdn = 5 > 2,  $p = 0.0034$ ). Interestingly, participants also employed the history view to check if their collaborators had “started work on my requests.” -P12 Moreover, the historical view also helped participants refine their prompts. For instance, P7 mentioned, “it helped me recall what or why I made certain changes.”

### 8.3 RQ2: How does the system support programmers to leverage collaborators’ work?

All participants agreed that the four types of mechanisms are **(1) easy to use**: refer (Mdn=6.5, SD=1.76), request (Mdn=5.5, SD=1.35), share (Mdn=6, SD=1.45) and link (Mdn=6, SD=1.27); **(2) intuitive to learn**: refer (Mdn=6, SD=1.15), request (Mdn=5, SD=1.68), share (Mdn=5.5, SD=1.31) and link (Mdn=6, SD=1.30); **(3) could fulfill their requirements**: refer (Mdn=6, SD=1.68), request (Mdn=5.5, SD=1.53), share (Mdn=6, SD=1.61) and link (Mdn=6, SD=1.36); and **(4) easy to control**: refer (Mdn=6, SD=1.37), request (Mdn=5, SD=1.03), share (Mdn=5, SD=1.61) and link (Mdn=5.5, SD=1.51). Overall, CoPrompt reduces participants’ need to communicate with

collaborators (Q4: Mdn = 5 > 3,  $p = 0.0122$ ) and facilitates participants to modify their prompt efficiently (Q7: Mdn = 5.5 > 2.5,  $p = 0.0065$ ).

**8.3.1 Request and auto-updates reduce mental load (D3).** CoPrompt helps participants to convey their needs clearly to the collaborator (Q5: Mdn = 5.5 > 2.5,  $p = 0.0122$ ) and reduce the need to “remember what tasks have not yet been done.” -P9 In collaborative programming, since the task distribution may involve procedural dependencies, there are often cases where participants need to wait for their collaborators to complete certain tasks. With the request-detect-update mechanism, participants just need to send request with brief descriptions. Participants reported that although it “took more time to get familiar with [request]” -P10, it ultimately saved them a significant amount of time. They could request tasks at higher-level headings if a suitable prompt had not yet been created by another collaborator.

Compared to communicating with collaborators using messages or audio in baseline, the request takes less cognitive effort since “it provides contextual information through the node” -P3 and does not request programmers to “write too much text to describe the changes.” -P2 Some participants (N=4) highlighted that this feature reduces cognitive load of ensuring a polite tone, as they no longer need to carefully craft messages to their collaborators, “I do not need to care about the manner.” -P3 In addition, the auto-update of the prompts sending requests after target input is detected from

collaborators offloads the tedious review and update process. P9 and P10 expressed their satisfaction with these features: “*Detection and auto-update saved me a lot of effort.*” -P10 P9 added, “*I only need to know what the task is about, without thinking about where and when to address it.*” However, the automatic update feature makes P12 feel a slight loss of control, and sometimes they would like to decide whether to handle updates manually or automatically.

**8.3.2 Sharing knowledge both proactively and reactively (D4).** Participants utilized the share mechanism for both proactive and reactive sharing of output or prompt segments. Proactively, some participants shared their results with collaborators, “*I know the encoded result would be used for following step.*” -P1 They also shared insights from the result through share, “*I would share some text that may not necessary is the prompt to share insights I got.*” -P7 Reactively, participants responded to their collaborators’ requests for specific data processing results through direct communication by sharing requested results. Participants also made use of the hierarchical structure to share elements by placing them under relevant headings. This occurred when collaborators had not yet delved into those specific details.

For the receiver of the shared content, most participants indicated that the content was useful for clarifying their own prompts and the pop-up message was clear enough for quick comprehension, “*very convenient as I do not need to find and refer to my collaborators’ work.*” -P5 Participants also revealed their trial-and-error strategies for dealing with the pop-up message that is too brief to understand or too long to read, “*Just like using ChatGPT, I just accept the answer and view its execution result to evaluate its effect.*” -P3 This strategy was adopted by many participants (N=9), as the trial-and-error cost is low with the view of the history version and the “*ability to trace back to previous versions.*” -P4

**8.3.3 Referencing reduces the effort of careful reading and copy-pasting (D4).** *CoPrompt* helps participants build on top of collaborators’ work easily (Q8:  $Mdn = 6 > 3, p = 0.0049$ ). All participants used the refer feature significantly more than other links when using *CoPrompt* to perform prompt co-engineering (Figure 8 Right). With the refer mechanism, participants no longer need to read the whole prompt and select utterances for copy-pasting to modify their own prompts. Instead, they handed off the comprehension work to *CoPrompt* by guiding it with the refer mechanism. The multi-level hierarchy display of prompts supported programmers in pinpointing the prompts they wished to refer to and enabled them to select the most appropriate level for reference. All participants agreed that the interaction process of refer “*reduced the cognitive switching between communication and code.*” -P5

Analyzing the open-coded results regarding the reasons for manually modifying prompts and code (Fig. 10), we observed that participants using *CoPrompt* made significantly fewer modifications due to “*Missing Requirements*” in both code and prompt edits (Code Edit:  $Mdn_{CoPrompt}=2 < Mdn_{baseline}=36, p < 5.10 \times 10^{-8}$ ; Prompt Edit:  $Mdn_{CoPrompt}=4.5 < Mdn_{baseline}=25.5, p = 2.57 \times 10^{-9}$ ). Similarly, we found that participants in the baseline condition needed to make significantly more modifications to the prompt and code due to “*Wrong Variable*”, which was caused by outdated or incorrect variables generated by the AI model (Code Edit:  $Mdn_{CoPrompt}=2$

$< Mdn_{baseline}=9.5, p = 1.93 \times 10^{-5}$ ; Prompt Edit:  $Mdn_{CoPrompt}=2 < Mdn_{baseline}=15, p = 1.39 \times 10^{-5}$ ). While these results can be attributed to all four mechanisms that collectively reduce the overall need for manual modifications, the substantial reduction in the need to modify due to missing requirements and wrong variables is more likely a result of the mechanism refer.

#### 8.4 RQ3: How does the system reduce the repetitive updates of prompts or code?

The link mechanism effectively reduced the frequency of repetitive updates by offering automatic synchronization (D2). Participants made significantly fewer modifications to code and prompts due to “*Sync*” when using *CoPrompt* ( $Mdn_{CoPrompt}=1 < Mdn_{baseline}=28, p < .001$ ), indicating a reduced need to update in response to changes made by others (Fig. 10). *CoPrompt* decreased the need for participants to repeatedly and iteratively modify prompts and significantly facilitated their ability to establish a shared understanding with collaborators (Q6:  $Mdn = 5 > 2, p = 0.0068$ ).

The link mechanism was deemed the most intuitive by all participants, and they unanimously agreed that it significantly reduced the workload associated with repetitive updates due to procedural dependencies. P4 expressed that it “*saved a lot of time on keep going back and forth between users*” and “*helped to offload some mental model*” -P5 without the need to keep track of collaborators’ changes on a certain task. Additionally, three participants employed links between headings to convey the synchronization of an entire subsection.

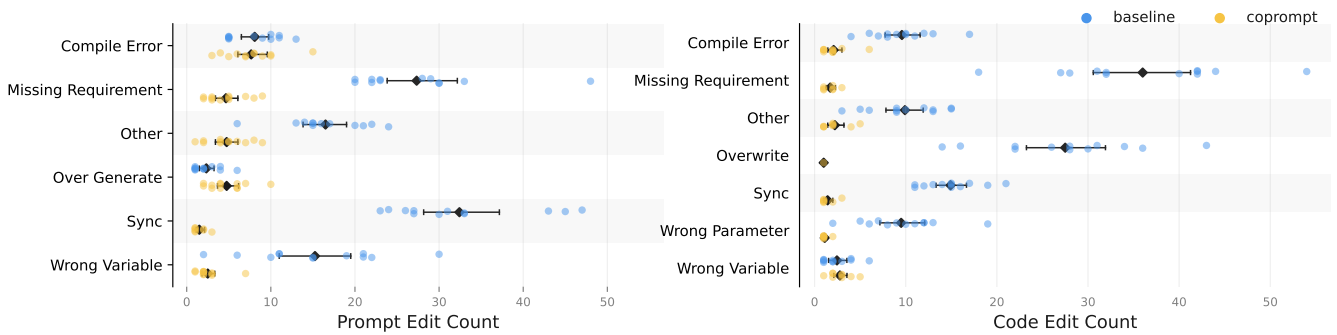
All participants mentioned that the refer and share mechanism simplifies the interaction of updating prompts. When there is a need to update existing prompts, many participants (N=9) leveraged refer so that they just needed to indicate the part of the prompt that requires modification and the target prompt reference. The share mechanism assisted programmers’ prompt engineering by allowing collaborators to pass knowledge to others and it only requires receivers to accept the shared prompt or code snippets and coarsely navigate to the target prompt, instead of “*carefully locating and manual copy-pasting.*” -P7 However, it also increases the likelihood of the models generating excessive content (as shown in Figure 10 left), requiring participants to manually refine the prompts ( $Mdn_{CoPrompt}=3 > M_{baseline}=2, p = .63$ ).

## 9 DISCUSSION

In this section, we discuss various topics that emerged during our study, offering insights and design implications. Additionally, we outline the primary limitations of our study and highlight areas for future research.

### 9.1 Using NL-predominant Expression in Collaborative Programming

Compared to code, NL-predominant expressions are easier to understand [57] and therefore helpful for maintaining shared understanding among collaborators [63, 65]. However, programmers still have difficulty making sense of the NL expressions used by their collaborators when they are too vague or at a high level. *CoPrompt* provided the hierarchical wiki, explanation view, and history view



**Figure 10: The left image displays edit prompt counts categorized by reasons of change for both the baseline and *CoPrompt* conditions; The right image illustrates code edit counts categorized by reasons of change for the same two conditions.**

to assist programmers’ comprehension with different levels of details. *CoPrompt* also allows programmers to share knowledge without losing context through four mechanisms, reducing their need for repetitive updates, copy-pasting, and synchronization in collaboration. Our findings reveal the possible benefits of prompt sharing and referring, including reducing task completion time, the need for communication and cognitive load. Participants found it useful to share intermediate results with context like the method details. However, it also increased the likelihood of the models generating excessive content, requiring participants to manually refine it.

## 9.2 Sense of Control in Code Generation Process

In *CoPrompt*, the four mechanisms involve a high level of automation [52] as LLMs are in charge of comprehending. This design significantly reduced the programmers’ cognitive load, while it raised concerns about the generated content: How to make the generation more satisfying and how to efficiently deal with unsatisfying results. Prior work in the domain of human-AI collaboration emphasized the importance of preserving user control when collaborating with LLMs [95] at different granularities [93]. To maintain user control in code generation process, *CoPrompt* allows programmers to manually tweak the code whenever the result is unsatisfying, and provides convenient history views for each prompt for checking and rolling back. Our design serves as a first step towards interaction design in prompt co-engineering, using cases of simple prompting templates to demonstrate the effectiveness of the workflow and core mechanisms. Future work could extend the design by considering more complex prompting strategies like few-shot prompts. The level of generation automation and controllability provided to the programmers could also be further investigated.

## 9.3 Synchronous and Asynchronous Collaboration

Overall, *CoPrompt* can be used in both synchronous and asynchronous collaboration settings. While the overall design of *CoPrompt* is catered towards synchronous programming, some features could assist asynchronous tasks. For instance, the message panel stores and displays each programmer’s usage of the mechanisms, which facilitates programmers who are initially offline to catch up with their collaborators’ work. Although the formative study was conducted

in real-time settings, participants overall followed the scatter-gather interaction [88], where they did not synchronously work on the task. Many interactions and communications were not carried out in real-time, where participants worked on their distributed tasks independently and then gathered the results without close and back-and-forth communication. All participants spent more than 15 minutes working asynchronously. We also evaluated the usage of *CoPrompt* in both synchronous and asynchronous settings using the 2-part user study: part 1 as real-time collaboration and part 2 as asynchronous collaboration. This demonstrates that *CoPrompt* is able to support both sync and async collaborative programming. However, there are some specific challenges for asynchronous collaboration that the current *CoPrompt* could not fully tackle, such as preserving essential elements when automatically updating off-line collaborators’ work. Future work could investigate ways of balancing user control and level of automation.

## 9.4 Supporting Mixed Collaborative Programming Styles

Although *CoPrompt* aims to facilitate prompt co-engineering, it does not depend on the “NL-first” style of collaboration. It also supports code-level collaboration, which allows programmers to manually write and modify code in the generation block. If programmers would like to have multiple prompts and code in one block, they can write multiple prompts inside the block, each starting a new line. The code will then be generated exactly below the prompts.

Participants edit both prompts and code when using *CoPrompt* (Figure 8), but they tend to add prompt blocks more frequently than code blocks. We observed that three participants opted to *add* new code blocks directly and used the four mechanisms on the code itself. They explained that this approach provided them with greater “controllability” -P9 over their program. Additionally, we observed several differences between NL-first and code-level collaboration when participants interacted with *CoPrompt*. 1) NL prompts served as high-level summaries, while code sharing often involved larger code chunks that required more navigation; 2) NL prompts were generally declarative and conveyed meaning directly, though exceptions existed for participants who added code blocks when code was in one line and “declarative” -P5 enough; 3) Participants’ familiarity with programming tasks also influenced their

preference. Those less familiar with data science preferred prompts that mostly explain the intentions behind them. Future research could further investigate the differences in using *CoPrompt* with NL and code. We also envision the possibility of a mixed-methods approach that combines elements of code and NL for more efficient collaboration.

Further, all four mechanisms in *CoPrompt* can be generalized to code-level manipulation by selecting code snippets as nodes. For instance, when utilizing the link to link two variables in code, the code will automatically update based on the context whenever the user-specified condition is met. This feature is particularly valuable in programming languages plagued by issues related to procedural dependencies [48]. However, the effectiveness of the share mechanism might be limited at the code level. When a programmer wants to share a piece of code with his collaborator, merely sharing the code without prompt information might result in misinterpretation for both the LLM and the collaborator. Future research could further distinguish between code and prompt sharing in collaborative NL programming.

## 9.5 Limitations

Our study of the *CoPrompt* presents important findings in the domain of collaborative NL programming. However, certain limitations should be recognized. First, *CoPrompt* was exclusively tested with programmers experienced in using LLM-driven code assistants. While NL programming is beneficial to a wide range of expertise levels, our pilot study indicated that programmers with limited experience often found it difficult to refine prompts effectively and were less active in the collaborative process. Future studies should be conducted including programmers unfamiliar with AI-driven code assistants to understand how varying expertise levels influence the collaborative workflows.

Secondly, we chose data science work as our case study due to its involvement of diverse participants [11] and the existence of procedural dependencies among various artifacts [48, 76]. The exploratory and explanatory aspects of data science necessitate close collaboration, frequent information exchange [7, 47, 64], and discussions [86]. These characteristics made data science an ideal initial case to explore the concept of prompt co-engineering. Beyond data science, challenges such as repetitive updating and synchronization persist [94] in general collaborative programming [92]. Since the workflow and mechanisms of *CoPrompt* are designed for prompt co-engineering rather than specific data science tasks, they remain applicable to reduce the need for repetitive updating and the effort for synchronization. Nonetheless, in certain cases, there might be some challenges that the current *CoPrompt* could not fully tackle, such as dealing with compilation error [20] and organizing spaghetti code [60]. For instance, *CoPrompt* could automatically update some linked artifacts in the spaghetti code, but could not ensure that there are no conflicts due to the entangled code structure. To better cope with the specific challenges in the wider domain of collaborative programming, future work should incorporate more customized designs.

While *CoPrompt* leverages a rich text editor with block-formatting features, it also works for less structured files or projects spanning multiple files, which are common in general programming tasks

besides data science. The coding blocks used in *CoPrompt* are the objectification of code snippets and prompt phrases, leading to the prompt co-engineering to a node-based workflow. By selecting any part of the project as a node, all mechanisms could be applied to the selected node. Meanwhile, the four mechanisms enabled participants to distribute the tasks into smaller sub-tasks. Prior work mainly investigated the collaborative styles of single authoring, divide & conquer and competitive authoring [68, 84]. Our proposed mechanisms could lower the cost of merging task results and leveraging others' output. This could enable the task distribution to be more nuanced so that every collaborator could program simultaneously, resulting in a more parallel collaboration and thus higher efficiency. This may also prevent the formation of spaghetti code generally.

## 10 CONCLUSION

In this work, we investigated the potential workflow of using NL prompts to conduct collaborative programming, especially prompt co-engineering. Our formative study revealed the workflow of prompt co-engineering and identified four challenges related to comprehension, synchronization, feedback, and reference of collaborators' prompts. To tackle these challenges, we introduced *CoPrompt*, a prototype to support prompt co-engineering through four novel mechanisms: share, refer, request, and link. Our user study indicated that *CoPrompt* effectively supported programmers' workflow of prompt co-engineering from comprehending collaborators' work to leveraging and sharing work.

## REFERENCES

- [1] 2023. Copilot: Your AI pair programmer. <https://github.com/features/copilot>
- [2] DataCanary Anna Montoya. 2016. House Prices - Advanced Regression Techniques. <https://kaggle.com/competitions/house-prices-advanced-regression-techniques>
- [3] Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. 2023. Prompting Is Programming: A Query Language for Large Language Models. *Proc. ACM Program. Lang.* 7, PLDI, Article 186 (Jun 2023), 24 pages. <https://doi.org/10.1145/3591300>
- [4] Susanne Bodker. 2015. Third-wave HCI, 10 years later—participation and sharing. *interactions* 22, 5 (2015), 24–31.
- [5] Virginia Braun and Victoria Clarke. 2012. *Thematic analysis*. American Psychological Association.
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [7] Souti Chattopadhyay, Ishita Prasad, Austin Z. Henley, Anita Sarma, and Titus Barik. 2020. What's Wrong with Computational Notebooks? Pain Points, Needs, and Design Opportunities. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '20). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3313831.3376729>
- [8] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri 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, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebguss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew 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:2107.03374 [cs.LG]
- [9] Yan Chen, Sang Won Lee, Yin Xie, YiWei Yang, Walter S. Lasecki, and Steve Oney. 2017. Codeon: On-Demand Software Development Assistance. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (Denver,



- Colorado, USA) (*CHI '17*). Association for Computing Machinery, New York, NY, USA, 6220–6231. <https://doi.org/10.1145/3025453.3025972>
- [10] KR1442 Chowdhary and KR Chowdhary. 2020. Natural language processing. *Fundamentals of artificial intelligence* (2020), 603–649.
- [11] Anamaria Crisan, Brittany Fiore-Gartland, and Melanie Tory. 2021. Passing the Data Baton : A Retrospective Analysis on Data Science Work and Workers. *IEEE Transactions on Visualization and Computer Graphics* 27, 2 (2021), 1860–1870. <https://doi.org/10.1109/TVCG.2020.3030340>
- [12] Will Cukierski. 2012. Titanic - Machine Learning from Disaster. <https://kaggle.com/competitions/titanic>
- [13] Allen Cypher and Daniel Conrad Halbert. 1993. *Watch what I do: programming by demonstration*. MIT press.
- [14] Sergio Cozzetti B. de Souza, Nicolas Anquetil, and Káthia M. de Oliveira. 2005. A Study of the Documentation Essential to Software Maintenance. In *Proceedings of the 23rd Annual International Conference on Design of Communication: Documenting & Designing for Pervasive Information* (Coventry, United Kingdom) (*SIGDOC '05*). Association for Computing Machinery, New York, NY, USA, 68–75. <https://doi.org/10.1145/1085313.1085331>
- [15] Paul Dourish and Victoria Bellotti. 1992. Awareness and coordination in shared workspaces. In *Proceedings of the 1992 ACM conference on Computer-supported cooperative work*. 107–114.
- [16] Hongfei Fan, Jiayao Gao, Hongming Zhu, Qin Liu, Yang Shi, and Chengzheng Sun. 2017. Balancing Conflict Prevention and Concurrent Work in Real-Time Collaborative Programming. In *Proceedings of the 12th Chinese Conference on Computer Supported Cooperative Work and Social Computing* (Chongqing, China) (*ChineseCSCW '17*). Association for Computing Machinery, New York, NY, USA, 217–220. <https://doi.org/10.1145/3127404.3127447>
- [17] Hongfei Fan, Chengzheng Sun, and Haifeng Shen. 2012. ATCoPE: Any-Time Collaborative Programming Environment for Seamless Integration of Real-Time and Non-Real-Time Teamwork in Software Development. In *Proceedings of the 2012 ACM International Conference on Supporting Group Work* (Sanibel Island, Florida, USA) (*GROUPE '12*). Association for Computing Machinery, New York, NY, USA, 107–116. <https://doi.org/10.1145/2389176.2389194>
- [18] Alexander J. Fiannaca, Chinmay Kulkarni, Carrie J Cai, and Michael Terry. 2023. Programming without a Programming Language: Challenges and Opportunities for Designing Developer Tools for Prompt Programming. In *Extended Abstracts of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (*CHI EA '23*). Association for Computing Machinery, New York, NY, USA, Article 235, 7 pages. <https://doi.org/10.1145/3544549.3585737>
- [19] R Stuart Geiger, Nelle Varoquaux, Charlotte Mazel-Cabasse, and Chris Holdgraf. 2018. The types, roles, and practices of documentation in data analytics open source software libraries: a collaborative ethnography of documentation work. *Computer Supported Cooperative Work (CSCW)* 27, 3-6 (2018), 767–802.
- [20] Max Goldman, Greg Little, and Robert C. Miller. 2011. Real-Time Collaborative Coding in a Web IDE. In *Proceedings of the 24th Annual ACM Symposium on User Interface Software and Technology* (Santa Barbara, California, USA) (*UIST '11*). Association for Computing Machinery, New York, NY, USA, 155–164. <https://doi.org/10.1145/2047196.2047215>
- [21] Google. 2023. Firebase is an app development platform that helps you build and grow apps and games users love. Backed by Google and trusted by millions of businesses around the world. <https://firebase.google.com/>
- [22] Thomas R. G. Green and Marian Petre. 1996. Usability analysis of visual programming environments: a ‘cognitive dimensions’ framework. *Journal of Visual Languages & Computing* 7, 2 (1996), 131–174.
- [23] Carl Gutwin and Saul Greenberg. 1998. Effects of awareness support on groupware usability. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 511–518.
- [24] Carl Gutwin and Saul Greenberg. 2002. A Descriptive Framework of Workspace Awareness for Real-Time Groupware. *Computer Supported Cooperative Work (CSCW)* 11, 3-4 (Sept. 2002), 411–446. <https://doi.org/10.1023/a:1021271517844>
- [25] Marijn Haverbeke. 2023. A toolkit for building rich-text editors on the web. <https://prosemirror.net/>
- [26] Andrew Head, Fred Hohman, Titus Barik, Steven M. Drucker, and Robert DeLine. 2019. Managing Messes in Computational Notebooks. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland UK) (*CHI '19*). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3290605.3300500>
- [27] Austin Z. Henley and Scott D. Fleming. 2014. The Patchworks Code Editor: Toward Faster Navigation with Less Code Arranging and Fewer Navigation Mistakes. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Toronto, Ontario, Canada) (*CHI '14*). Association for Computing Machinery, New York, NY, USA, 2511–2520. <https://doi.org/10.1145/2556288.2557073>
- [28] James D Herbsleb, Helen Klein, Gary M Olson, Hans Brunner, Judith S Olson, and Joe Harding. 1995. Object-oriented analysis and design in software project teams. *Human-Computer Interaction* 10, 2-3 (1995), 249–292.
- [29] Di Huang, Ziyuan Nan, Xing Hu, Pengwei Jin, Shaohui Peng, Yuanbo Wen, Rui Zhang, Zidong Du, Qi Guo, Yewen Pu, and Yunji Chen. 2023. ANPL: Compiling Natural Programs with Interactive Decomposition. [arXiv:2305.18498](https://arxiv.org/abs/2305.18498) [cs.LG]
- [30] Edwin L Hutchins, James D Hollan, and Donald A Norman. 1985. Direct manipulation interfaces. *Human-computer interaction* 1, 4 (1985), 311–338.
- [31] Naman Jain, Skanda Vaidyanath, Arun Iyer, Nagarajan Natarajan, Suresh Parthasarathy, Sriram Rajamani, and Rahul Sharma. 2022. Jigsaw: Large Language Models Meet Program Synthesis. In *Proceedings of the 44th International Conference on Software Engineering* (Pittsburgh, Pennsylvania) (*ICSE '22*). Association for Computing Machinery, New York, NY, USA, 1219–1231. <https://doi.org/10.1145/3510003.3510203>
- [32] Ellen Jiang, Kristen Olson, Edwin Toh, Alejandra Molina, Aaron Donsbach, Michael Terry, and Carrie J Cai. 2022. PromptMaker: Prompt-Based Prototyping with Large Language Models. In *Extended Abstracts of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (*CHI EA '22*). Association for Computing Machinery, New York, NY, USA, Article 35, 8 pages. <https://doi.org/10.1145/3491101.3503564>
- [33] Ellen Jiang, Edwin Toh, Alejandra Molina, Kristen Olson, Claire Kayacik, Aaron Donsbach, Carrie J Cai, and Michael Terry. 2022. Discovering the Syntax and Strategies of Natural Language Programming with Generative Language Models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (*CHI '22*). Association for Computing Machinery, New York, NY, USA, Article 386, 19 pages. <https://doi.org/10.1145/3491102.3501870>
- [34] Mira Kajko-Mattsson. 2005. A survey of documentation practice within corrective maintenance. *Empirical Software Engineering* 10 (2005), 31–55.
- [35] Amy J. Ko, Robin Abraham, Laura Beckwith, Alan Blackwell, Margaret Burnett, Martin Erwig, Chris Scaffidi, Joseph Lawrence, Henry Lieberman, Brad Myers, Mary Beth Rosson, Gregg Rothermel, Mary Shaw, and Susan Wiedenbeck. 2011. The State of the Art in End-User Software Engineering. *ACM Comput. Surv.* 43, 3, Article 21 (apr 2011), 44 pages. <https://doi.org/10.1145/1922649.1922658>
- [36] Kirby Kuznia, Swaroop Mishra, Mihir Parmar, and Chitta Baral. 2022. Less is More: Summary of Long Instructions is Better for Program Synthesis. [arXiv:2203.08597](https://arxiv.org/abs/2203.08597) [cs.CL]
- [37] Andrew K. Lampinen, Ishita Dasgupta, Stephanie C. Y. Chan, Kory Matthewson, Michael Henry Tessler, Antonia Creswell, James L. McClelland, Jane X. Wang, and Felix Hill. 2022. Can language models learn from explanations in context? [arXiv:2204.02329](https://arxiv.org/abs/2204.02329) [cs.CL]
- [38] J Chris Lauwers and Keith A Lantz. 1990. Collaboration awareness in support of collaboration transparency: Requirements for the next generation of shared window systems. In *Proceedings of the SIGCHI conference on Human factors in computing systems*. 303–311.
- [39] James R. Lewis, Brian S. Utesch, and Deborah E. Maher. 2013. UMUX-LITE: When There’s No Time for the SUS. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Paris, France) (*CHI '13*). Association for Computing Machinery, New York, NY, USA, 2099–2102. <https://doi.org/10.1145/2470654.2481287>
- [40] Zongjie Li, Chaozheng Wang, Zhibo Liu, Haoxuan Wang, Dong Chen, Shuai Wang, and Cuiyun Gao. 2023. CCTEST: Testing and Repairing Code Completion Systems. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. 1238–1250. <https://doi.org/10.1109/ICSE48619.2023.00110>
- [41] Henry Lieberman. 2001. *Your wish is my command: Programming by example*. Morgan Kaufmann.
- [42] Hugo Liu and Henry Lieberman. 2005. Programmatic Semantics for Natural Language Interfaces. In *CHI '05 Extended Abstracts on Human Factors in Computing Systems* (Portland, OR, USA) (*CHI EA '05*). Association for Computing Machinery, New York, NY, USA, 1597–1600. <https://doi.org/10.1145/1056808.1056975>
- [43] Michael Xieyang Liu, Aniket Kittur, and Brad A. Myers. 2021. To Reuse or Not To Reuse? A Framework and System for Evaluating Summarized Knowledge. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW1, Article 166 (apr 2021), 35 pages. <https://doi.org/10.1145/3449240>
- [44] Michael Xieyang Liu, Advait Sarkar, Carina Negreanu, Benjamin Zorn, Jack Williams, Neil Toronto, and Andrew D. Gordon. 2023. “What It Wants Me To Say”: Bridging the Abstraction Gap Between End-User Programmers and Code-Generating Large Language Models. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (*CHI '23*). Association for Computing Machinery, New York, NY, USA, Article 598, 31 pages. <https://doi.org/10.1145/3544548.3580817>
- [45] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-Train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing. *ACM Comput. Surv.* 55, 9, Article 195 (jan 2023), 35 pages. <https://doi.org/10.1145/3560815>
- [46] Vivian Liu and Lydia B Chilton. 2022. Design Guidelines for Prompt Engineering Text-to-Image Generative Models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (*CHI '22*). Association for Computing Machinery, New York, NY, USA, Article 384, 23 pages. <https://doi.org/10.1145/3491102.3501825>
- [47] Yang Liu, Tim Althoff, and Jeffrey Heer. 2020. Paths Explored, Paths Omitted, Paths Obscured: Decision Points & Selective Reporting in End-to-End Data Analysis. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (*CHI '20*). Association for Computing Machinery, New York, NY, USA, 1–14. <https://doi.org/10.1145/3313831.3376533>

- [48] Yang Liu, Alex Kale, Tim Althoff, and Jeffrey Heer. 2021. Boba: Authoring and Visualizing Multiverse Analyses. *IEEE Transactions on Visualization and Computer Graphics* 27, 2 (2021), 1753–1763. <https://doi.org/10.1109/TVCG.2020.3028985>
- [49] Ewa Luger and Abigail Sellen. 2016. “Like Having a Really Bad PA”: The Gulf between User Expectation and Experience of Conversational Agents. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (San Jose, California, USA) (CHI ’16). Association for Computing Machinery, New York, NY, USA, 5286–5297. <https://doi.org/10.1145/2858036.2858288>
- [50] Yifan Ma, Batu Qi, Wenhua Xu, Mingjie Wang, Bowen Du, and Hongfei Fan. 2022. Integrating Real-Time and Non-Real-Time Collaborative Programming: Workflow, Techniques, and Prototypes. *Proc. ACM Hum.-Comput. Interact.* 7, GROUP, Article 13 (dec 2022), 19 pages. <https://doi.org/10.1145/3567563>
- [51] Walid Maalej and Martin P. Robillard. 2013. Patterns of Knowledge in API Reference Documentation. *IEEE Transactions on Software Engineering* 39, 9 (2013), 1264–1282. <https://doi.org/10.1109/TSE.2013.12>
- [52] Maximilian Mackeprang, Claudia Müller-Birn, and Maximilian Timo Stauss. 2019. Discovering the Sweet Spot of Human-Computer Configurations: A Case Study in Information Extraction. 3, CSCW, Article 195 (nov 2019), 30 pages. <https://doi.org/10.1145/3359297>
- [53] Andrew M McNutt, Chenglong Wang, Robert A Deline, and Steven M. Drucker. 2023. On the Design of AI-Powered Code Assistants for Notebooks. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (CHI ’23). Association for Computing Machinery, New York, NY, USA, Article 434, 16 pages. <https://doi.org/10.1145/3544548.3580940>
- [54] Microsoft. 2021. Visual studio code for the web. <https://code.visualstudio.com/docs/editor/vscode-web>
- [55] code Microsoft. 2021. Use Microsoft Live share to collaborate with Visual Studio Code. <https://code.visualstudio.com/learn/collaboration/live-share>
- [56] Rada Mihalcea, Hugo Liu, and Henry Lieberman. 2006. NLP (natural language processing) for NLP (natural language programming). In *Computational Linguistics and Intelligent Text Processing: 7th International Conference, CICLing 2006, Mexico City, Mexico, February 19-25, 2006. Proceedings* 7. Springer, 319–330.
- [57] L. A. Miller. 1981. Natural language programming: Styles, strategies, and contrasts. *IBM Systems Journal* 20, 2 (1981), 184–215. <https://doi.org/10.1147/sj.202.0184>
- [58] Michael Muller, Lydia B Chilton, Anna Kantosalo, Q. Vera Liao, Mary Lou Maher, Charles Patrick Martin, and Greg Walsh. 2023. GenAICHI 2023: Generative AI and HCI at CHI 2023. In *Extended Abstracts of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (CHI EA ’23). Association for Computing Machinery, New York, NY, USA, Article 350, 7 pages. <https://doi.org/10.1145/3544549.3573794>
- [59] Michael Muller, Ingrid Lange, Dakuo Wang, David Piorkowski, Jason Tsay, Q. Vera Liao, Casey Dugan, and Thomas Erickson. 2019. How Data Science Workers Work with Data: Discovery, Capture, Curation, Design, Creation. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk) (CHI ’19). Association for Computing Machinery, New York, NY, USA, 1–15. <https://doi.org/10.1145/3290605.3300356>
- [60] Brad A. Myers. 1991. Separating Application Code from Toolkits: Eliminating the Spaghetti of Call-Backs. In *Proceedings of the 4th Annual ACM Symposium on User Interface Software and Technology* (Hilton Head, South Carolina, USA) (UIST ’91). Association for Computing Machinery, New York, NY, USA, 211–220. <https://doi.org/10.1145/120782.120805>
- [61] Steve Oney, Christopher Brooks, and Paul Resnick. 2018. Creating Guided Code Explanations with ChatCodes. *Proc. ACM Hum.-Comput. Interact.* 2, CSCW, Article 131 (nov 2018), 20 pages. <https://doi.org/10.1145/3274400>
- [62] OpenAI. 2023. GPT-4 Technical Report. arXiv:2303.08774 [cs.CL]
- [63] Yoann Padioleau, Lin Tan, and Yuanyuan Zhou. 2009. Listening to programmers – Taxonomies and characteristics of comments in operating system code. In *2009 IEEE 31st International Conference on Software Engineering*. 331–341. <https://doi.org/10.1109/ICSE.2009.5070533>
- [64] Rock Yuren Pang, Ruotong Wang, Joely Nelson, and Leilani Battle. 2022. How Do Data Science Workers Communicate Intermediate Results?. In *2022 IEEE Visualization in Data Science (VDS)*. IEEE, 46–54.
- [65] Rock Yuren Pang, Ruotong Wang, Joely Nelson, and Leilani Battle. 2022. How Do Data Science Workers Communicate Intermediate Results?. In *2022 IEEE Visualization in Data Science (VDS)*. 46–54. <https://doi.org/10.1109/VDS57266.2022.00010>
- [66] Soya Park, Amy X. Zhang, and David R. Karger. 2018. Post-Literate Programming: Linking Discussion and Code in Software Development Teams. In *Adjunct Proceedings of the 31st Annual ACM Symposium on User Interface Software and Technology* (Berlin, Germany) (UIST ’18 Adjunct). Association for Computing Machinery, New York, NY, USA, 51–53. <https://doi.org/10.1145/3266037.3266098>
- [67] Hammond Pearce, Benjamin Tan, Baleegh Ahmad, Ramesh Karri, and Brendan Dolan-Gavitt. 2023. Examining Zero-Shot Vulnerability Repair with Large Language Models. In *2023 IEEE Symposium on Security and Privacy (SP)*. 2339–2356. <https://doi.org/10.1109/SP46215.2023.10179324>
- [68] I.R. Posner and R.M. Baecker. 1992. How people write together (groupware). In *Proceedings of the Twenty-Fifth Hawaii International Conference on System Sciences*, Vol. iv. 127–138 vol.4. <https://doi.org/10.1109/HICSS.1992.183420>
- [69] Pyodide. 2023. Pyodide is a Python distribution for the browser and Node.js based on WebAssembly. <https://github.com/pyodide/pyodide>
- [70] Luigi Quaranta, Fabio Calefato, and Filippo Lanubile. 2022. Eliciting Best Practices for Collaboration with Computational Notebooks. *Proc. ACM Hum.-Comput. Interact.* 6, CSCW1, Article 87 (apr 2022), 41 pages. <https://doi.org/10.1145/3512934>
- [71] Laria Reynolds and Kyle McDonell. 2021. Prompt Programming for Large Language Models: Beyond the Few-Shot Paradigm. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (CHI EA ’21). Association for Computing Machinery, New York, NY, USA, Article 314, 7 pages. <https://doi.org/10.1145/3411763.3451760>
- [72] Nico Ritschel, Felipe Fronchetti, Reid Holmes, Ronald Garcia, and David C. Shepherd. 2022. Can Guided Decomposition Help End-Users Write Larger Block-Based Programs? A Mobile Robot Experiment. *Proc. ACM Program. Lang.* 6, OOPSLA2, Article 133 (oct 2022), 26 pages. <https://doi.org/10.1145/3563296>
- [73] Tobias Roehm, Rebecca Tiarks, Rainer Koschke, and Walid Maalej. 2012. How do professional developers comprehend software?. In *2012 34th International Conference on Software Engineering (ICSE)*. 255–265. <https://doi.org/10.1109/ICSE.2012.6227188>
- [74] Steven I Ross, Fernando Martinez, Stephanie Houde, Michael Muller, and Justin D Weisz. 2023. The programmer’s assistant: Conversational interaction with a large language model for software development. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*. 491–514.
- [75] Advait Sarkar, Andrew D. Gordon, Carina Negreanu, Christian Poelitz, Sruti Srinivasa Ragavan, and Ben Zorn. 2022. What is it like to program with artificial intelligence? <https://doi.org/10.48550/arXiv.2208.06213> [cs.HC]
- [76] Abhraneel Sarma, Alex Kale, Michael Jongho Moon, Nathan Taback, Fanny Chevalier, Jessica Hullman, and Matthew Kay. 2023. Multiverse: Multiplexing Alternative Data Analyses in R Notebooks. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (CHI ’23). Association for Computing Machinery, New York, NY, USA, Article 148, 15 pages. <https://doi.org/10.1145/3544548.3580726>
- [77] Lin Shi, Hao Zhong, Tao Xie, and Mingshu Li. 2011. An empirical study on evolution of API documentation. In *Fundamental Approaches to Software Engineering: 14th International Conference, FASE 2011, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2011, Saarbrücken, Germany, March 26–April 3, 2011. Proceedings* 14. Springer, 416–431.
- [78] Hendrik Strobelt, Albert Webson, Victor Sanh, Benjamin Hoover, Johanna Beyer, Hanspeter Pfister, and Alexander M. Rush. 2023. Interactive and Visual Prompt Engineering for Ad-hoc Task Adaptation with Large Language Models. *IEEE Transactions on Visualization and Computer Graphics* 29, 1 (2023), 1146–1156. <https://doi.org/10.1109/TVCG.2022.3209479>
- [79] Tiptap. 2023. Tiptap is a suite of open source content editing and real-time collaboration tools for developers building apps like Notion or Google Docs. <https://tiptap.dev/>
- [80] Immanuel Trummer. 2022. CodexDB: Generating Code for Processing SQL Queries using GPT-3 Codex. arXiv:2204.08941 [cs.DB]
- [81] Priyan Vaithilingam, Elena L. Glassman, Peter Groenwegen, Sumit Gulwani, Austin Z. Henley, Rohan Malpani, David Pugh, Arjun Radhakrishna, Gustavo Soares, Joey Wang, and Aaron Yim. 2023. Towards More Effective AI-Assisted Programming: A Systematic Design Exploration to Improve Visual Studio IntelliCode’s User Experience. In *2023 IEEE/ACM 45th International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. 185–195. <https://doi.org/10.1109/ICSE-SEIP58684.2023.00022>
- [82] Helena Vasconcelos, Gagan Bansal, Adam Fournay, Q Vera Liao, and Jennifer Wortman Vaughan. 2023. Generation probabilities are not enough: Exploring the effectiveness of uncertainty highlighting in AI-powered code completions. *arXiv preprint arXiv:2302.07248* (2023).
- [83] April Yi Wang. 2022. Improving Real-Time Collaborative Data Science Through Context-Aware Mechanisms. In *2022 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)*. 1–3. <https://doi.org/10.1109/VL/HCC53370.2022.9833140>
- [84] April Yi Wang, Anant Mittal, Christopher Brooks, and Steve Oney. 2019. How Data Scientists Use Computational Notebooks for Real-Time Collaboration. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 39 (nov 2019), 30 pages. <https://doi.org/10.1145/3359141>
- [85] April Yi Wang, Dakuo Wang, Jaimie Drozdal, Michael Muller, Soya Park, Justin D. Weisz, Xuye Liu, Lingfei Wu, and Casey Dugan. 2022. Documentation Matters: Human-Centered AI System to Assist Data Science Code Documentation in Computational Notebooks. *ACM Trans. Comput.-Hum. Interact.* 29, 2, Article 17 (jan 2022), 33 pages. <https://doi.org/10.1145/3489465>
- [86] April Yi Wang, Zihan Wu, Christopher Brooks, and Steve Oney. 2020. Callisto: Capturing the “Why” by Connecting Conversations with Computational Narratives. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI ’20). Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/3313831.3376740>

[87] Dakuo Wang, Elizabeth Churchill, Pattie Maes, Xiangmin Fan, Ben Shneiderman, Yuanchun Shi, and Qianying Wang. 2020. From Human-Human Collaboration to Human-AI Collaboration: Designing AI Systems That Can Work Together with People. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (*CHI EA '20*). Association for Computing Machinery, New York, NY, USA, 1–6. <https://doi.org/10.1145/3334480.3381069>

[88] Dakuo Wang, Justin D. Weisz, Michael Muller, Parikshit Ram, Werner Geyer, Casey Dugan, Yla Tausczik, Horst Samulowitz, and Alexander Gray. 2019. Human-AI Collaboration in Data Science: Exploring Data Scientists’ Perceptions of Automated AI. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 211 (nov 2019), 24 pages. <https://doi.org/10.1145/3359313>

[89] Sitong Wang, Samia Menon, Tao Long, Keren Henderson, Dingzeyu Li, Kevin Crowston, Mark Hansen, Jeffrey V Nickerson, and Lydia B Chilton. 2023. Reel-Framer: Co-creating News Reels on Social Media with Generative AI. *arXiv preprint arXiv:2304.09653* (2023).

[90] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. *arXiv:2201.11903* [cs.CL]

[91] John Wenskovich, Jian Zhao, Scott Carter, Matthew Cooper, and Chris North. 2019. Albireo: An Interactive Tool for Visually Summarizing Computational Notebook Structure. In *2019 IEEE Visualization in Data Science (VDS)*. 1–10. <https://doi.org/10.1109/VDS48975.2019.8973385>

[92] L. Williams, R.R. Kessler, W. Cunningham, and R. Jeffries. 2000. Strengthening the case for pair programming. *IEEE Software* 17, 4 (2000), 19–25. <https://doi.org/10.1109/52.854064>

[93] Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022. AI Chains: Transparent and Controllable Human-AI Interaction by Chaining Large Language Model Prompts. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (*CHI '22*). Association for Computing Machinery, New York, NY, USA, Article 385, 22 pages. <https://doi.org/10.1145/3491102.3517582>

[94] Kimberly Michelle Ying and Kristy Elizabeth Boyer. 2020. Understanding Students’ Needs for Better Collaborative Coding Tools. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems* (<conf-loc>, <city>Honolulu</city>, <state>HI</state>, <country>USA</country>, </conf-loc>) (*CHI EA '20*). Association for Computing Machinery, New York, NY, USA, 1–8. <https://doi.org/10.1145/3334480.3383068>

[95] Zheng Zhang, Jie Gao, Ranjodh Singh Dhaliwal, and Toby Jia-Jun Li. 2023. VISAR: A Human-AI Argumentative Writing Assistant with Visual Programming and Rapid Draft Prototyping. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology* (San Francisco, CA, USA) (*UIST '23*). Association for Computing Machinery, New York, NY, USA, Article 5, 30 pages. <https://doi.org/10.1145/3586183.3606800>

## A SYSTEM DESIGN

### A.1 System Architecture

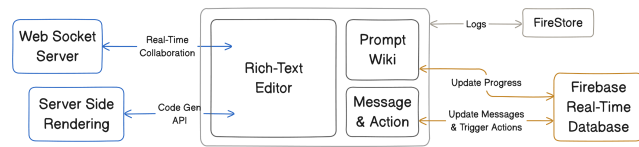


Figure 11: The system architecture of *CoPrompt* consists of a web socket server and a server-side rendering route for invoking the generative AI API. All messages and triggered actions are managed and updated using the Firebase Real-Time Database, while logs are stored in Firestore.

### A.2 Prompt Templates

[ADD] As an AI Python code-writing assistant, your task is to generate specific Python code segments based on the user's input in [SPROMPT]. When a request for code continuation is made, refer to context from [SREFER]. In such cases, concatenate the prompt with the previously generated code from [SREFER] to form a cohesive extension.

Requirements:

- Accurately generate Python code that aligns with the directives in [SPROMPT].
- In case of code continuation, seamlessly integrate new code with the existing code from [SREFER].
- Focus on the latter parts of the prompt if it contains multiple instructions, but ensure that the entire prompt is considered.
- Provide concise responses, strictly limited to the required code without additional explanations, contextual information, or comments.

Figure 12: Prompt template for Adding prompt blocks.

As an AI Python code editor, your role is to modify existing Python code based on user instructions in [SPROMPT]. Use the code provided in [SCONTEXT] as the base for your modifications. The user will specify the exact changes needed, which may include bug fixes, adding new functionalities, or altering existing ones.

Guidelines:

- Clearly understand and follow the modification instructions in [SPROMPT].
- Ensure that the edited code remains consistent with the overall logic and structure of the original code in [SCONTEXT].
- If the prompt includes multiple edit requests, prioritize them based on their order in the prompt.
- The response should be limited to the modified code only, without additional explanations or comments.
- Maintain the integrity of the original code, ensuring that changes do not introduce new errors or disrupt existing functionalities.

Figure 13: Prompt template for editing existing prompt blocks.

You are an AI model specializing in analyzing Python code dependencies. Your task is to assess pairs of code segments provided by the user, identify procedural dependencies between them, and determine if any changes in one segment necessitate updates in the other.

**Input:** A list of pairs of code segments, each with a unique identifier (id) and index. Instructions:

- **Analyze Dependencies:** For each pair of code segments, identify the dependencies and the impact of changes in one segment on the other.
- **Decide on Updates:** Based on the dependency analysis, determine if updating one segment requires changes to the other. Use logical reasoning to justify your decision.
- **Provide Updated Code:** If an update is necessary, provide the updated code for the affected segment.

**Output Format:**

- id: [Unique identifier of the code pair]
- index: [Index number of the code pair]
- updated: [Boolean indicating if an update is necessary]
- updated\_code: [String containing the updated code, if applicable]

**Example Input/Output:**

Example N:

Input:

- id: 'pair1',
- index: 1,
- Code\_Segment A:

```
def calculate_sum(a, b):
    return a + b
```

- Code Segment B:

```
result = calculate_sum(2, 3)
print("The sum is:", result)
```

Output:

- id: 'pair1',
- index: 1,
- updated: True,
- updated\_code:

```
def calculate_sum(a, b, c): # Updated to include a third parameter
    return a + b + c

result = calculate_sum(2, 3, 5) # Updated to pass a third argument
print("The sum is:", result)
```

- Reasoning: The update is required because Code Segment A has been modified to include an additional parameter in the function calculate\_sum. Consequently, Code Segment B needs to be updated to pass the additional argument when calling the function.

**Figure 14: Prompt template for link mechanism.**

You are an AI model tasked with evaluating whether the provided code fulfills specific user requests. Each request comes with a set of details including an ID, an index number, a prompt, the generated code, and supplementary information.

**Input:** A list of user requests, each containing:

- id: [string] - The ID of the request.
- index: [number] - The index of the request.
- prompt: [string] - The original prompt used for the code request.
- code: [string] - The generated code based on the prompt.
- supplementary\_info: [string] - Additional information or specific requirements related to the request.

**Task:** Evaluate if the 'code' from each request successfully meets the criteria or needs expressed in the 'supplementary\_info'. Please logically reason through each request, considering how the 'code' corresponds to the 'supplementary\_info'.

**Output Format:** Return a list of objects, each including:

- id: [string] - The same ID as in the input.
- index: [number] - The same index as in the input.
- fulfilled: [boolean] - Indicates whether the code fulfills the supplementary

**Example Input:**

```
completion: "iris = load_iris()"
requests: [{
  id: "qodwjqoiojJ@OIJoihy2ge1i12jhqiuo",
  index: 1,
  prompt: "Visualize iris data in pairplot",
  code: "import pandas as pd\nimport seaborn as sns\nsns.pairplot(iris, hue='species')",
  supplementary_info: "Plot data"
},
{
  id: "qwedojoio3289901hANUII-jio",
  index: 6,
  prompt: "Visualize wine data",
  code: "import pandas as pd\nimport seaborn as sns\nsns.pairplot(wine, hue='species')",
  supplementary_info: "Provide me wine data"
}]
```

**Example Output:**

```
The UserB is requesting code to "plot data" from the prompt "visualize iris data" UserA used to
The UserB is requesting code to "provide me wine data" from the prompt "visualize wine data" Use

[
  {
    id: "qodwjqoiojJ@OIJoihy2ge1i12jhqiuo",
    index: 1,
    fulfilled: true
  },
  {
    id: "qwedojoio3289901hANUII-jio",
    index: 6,
    fulfilled: false
  }
]
```

**Figure 15: Prompt template for request mechanism.**

## **B SURVEY QUESTIONS**

### **B.1 System Usability Likert Scale Questions**

- (1) The system helps me understand what my collaborator is doing.
- (2) The system reduces the need for me to keep track of collaborators' progress.
- (3) The system helps me understand the prompt written by collaborators.
- (4) The system reduces the need to communicate with collaborators.
- (5) The system helps me convey my needs to the collaborator.
- (6) The system supports me in reaching common ground with my collaborators.
- (7) The system supports me on how to engineer my prompt.
- (8) The system helps me build on top of collaborators' work easily.

### **B.2 Mechanism Usability Likert Scale Questions**

- (1) I think this mechanism is easy to learn.
- (2) I think this mechanism is easy to control.
- (3) I think this mechanism is easy to use.
- (4) I think this mechanism can help me achieve what I want.

### **B.3 UMUX-LITE**

- (1) This system is easy to use.
- (2) This system's capabilities meet my requirements.

### **B.4 NASA-TLX**

- (1) How mentally demanding was the task?
- (2) How physically demanding was the task?
- (3) How hurried or rushed was the pace of the task?
- (4) How successful were you in accomplishing what you were asked to do?
- (5) How hard did you have to work to accomplish your level of performance?
- (6) How insecure, discouraged, irritated, stressed, and annoyed were you?