# IntentFlow: Supporting Interactive and Fluid Intent Communication with Large Language Models

Yoonsu Kim
School of Computing
KAIST
Daejeon, Republic of Korea
yoonsu16@kaist.ac.kr

Brandon Chin
College of Engineering
University of California Berkeley
Berkeley, California, USA
brandoncjw@hkn.eecs.berkeley.edu

Kihoon Son
School of Computing
KAIST
Daejeon, Republic of Korea
kihoon.son@kaist.ac.kr

Seoyoung Kim
School of Computing
KAIST
Daejeon, Republic of Korea
youthskim@kaist.ac.kr

Juho Kim
School of Computing
KAIST
Daejeon, Republic of Korea
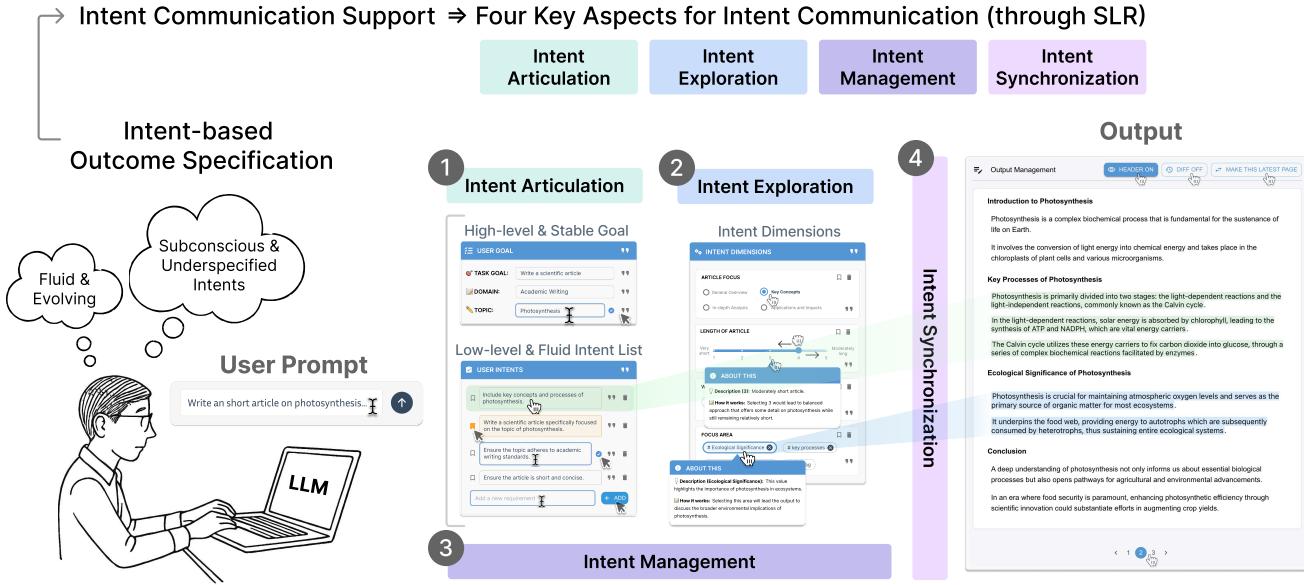juhokim@kaist.ac.kr

**Figure 1: Our work addresses the challenges of intent-based outcome specification in human–LLM interaction, where user intents are often subconscious and underspecified, and tend to be fluid and evolving. Through a systematic literature review, we identified four aspects of intent communication support: articulation, exploration, management, and synchronization. We instantiate these aspects in IntentFlow, an LLM system for writing tasks that (1) helps users articulate vague or subconscious intents, (2) supports exploration of alternative or emerging directions, (3) manages evolving intents over time, and (4) synchronizes intents with generated outputs through linking.**

## ABSTRACT

Effective collaboration with generative AI systems requires users to clearly communicate their intents (intent-based outcome specification). Yet such intents are often underspecified and evolve during interaction, dynamic support for intent communication is essential. Through a systematic literature review of 33 papers, we synthesize a structured understanding of intent communication, identifying four key aspects: **articulation**, **exploration**, **management**, and **synchronization**. Building on these findings, we derived design implications that translate them into actionable design and implemented IntentFlow, a system for LLM-based writing that realizes

these implications through adjustable UIs, intent-to-output linking, and versioned refinement. A technical evaluation (N=60) and a within-subjects study (N=12) confirm that INTENTFLOW helps users discover, elaborate, and consolidate their intents into a curated set. Interaction logs further reveal a shift from reactive error correction to proactive intent refinement. Our work demonstrates how a system effectively designed to support these four communication aspects can substantially enhance human-LLM interaction.

## CCS CONCEPTS

• **Human-centered computing → Interactive systems and tools**.

## KEYWORDS

Large Language Models, Generative AI Systems, Human-AI Interaction, Intent Communication, AI alignment, Transparency, Writing Assistant

## 1 INTRODUCTION

Advances in generative AI, such as large language models (LLMs) and vision-language models (VLMs), now allow people to easily obtain high-quality outputs, ranging from polished text [30, 36, 53, 74, 78] and functional code [40, 43] to complete visual designs [54], simply by describing what they want in a natural language prompt. This represents what Jakob Nielsen calls a new UI paradigm in computing history—*intent-based outcome specification*—where users tell the computer *what they want*, rather than *how to do it* [47].

While this prompt-based interaction appears straightforward, a closer look at this process reveals a deeper complexity in how users communicate their needs. Drawing from communication theory [21, 34, 62], we can understand prompting as conveying two key types of information: the user's **goal** and their **intent**. A **goal** represents the high-level objective (e.g., writing a cover letter), which tends to be explicit and stable throughout the process. In contrast, **intents** encompass the more low-level strategies, preferences, or constraints for achieving that goal, often emerging or shifting during the process (e.g., deciding whether to introduce one's background before the motivation, adopting a polite tone, or keeping certain sections concise). Communication scholars emphasize that intent is inherently fluid and often subconscious [21, 34, 35, 62]. People may not be fully aware of their own intents at the outset; these may evolve, surface gradually, or even fade as they reflect on outcomes and iterate on their actions. HCI researchers have similarly noted that users' intents are often not fully articulated at the beginning of interaction [5, 19, 56, 65].

However, current LLM interfaces, especially those built on linear, prompt-based interactions, offer limited support for this evolving nature of intent. Prior research in HCI has identified several challenges users face when trying to communicate their intent to these systems: the **difficulty of articulating vague or unspecified**

**intent** through natural language alone [66, 76], **the opacity of model outputs** that hinders users from building a clear mental model of how their prompts influence different parts of the output [27, 63, 79], and the **lack of a structured method** to manage and modify evolving intents, which are often scattered across the conversation history [65]. These are not isolated problems but symptoms of a broader issue: the lack of support for intent communication as a dynamic, fluid process.

To develop a comprehensive understanding of how the HCI community has addressed these challenges, we conducted a systematic literature review (SLR) of 33 papers on generative AI systems. Our analysis revealed four key, interdependent aspects of intent communication: (1) **Intent Articulation** (helping users express vague intents), (2) **Intent Exploration** (supporting users in discovering and expanding their intents), (3) **Intent Management** (providing a structured way to manage evolving intents), (4) **Intent Synchronization** (ensuring a mutual understanding between the user and LLMs regarding communicated intents). Critically, we found that these four aspects form a tightly coupled cycle: articulation sparks exploration, exploration introduces new intents that require management, and synchronization closes the loop by verifying alignment—feeding back into subsequent articulation and exploration.

While existing work has contributed valuable solutions to parts of these four aspects, for instance, by allowing direct manipulation of text or reifying prompts into manipulable objects [19, 45, 54], no single system has yet synthesized these insights to support the entire, interconnected cycle of intent communication. Our primary contribution is not the invention of each individual component, but the holistic synthesis and integration of all four aspects of intent communication into a cohesive and fluid workflow.

From these findings, we derived a set of design implications that guide how systems can support intent communication holistically across all four aspects. We instantiate these design implications in INTENTFLOW, a system designed for LLM-based writing tasks, chosen as a robust testbed given its prevalence in LLM usage and its multi-faceted, evolving, and iterative nature [16, 53]. INTENTFLOW provides comprehensive support for all four aspects of intent communication. First, it facilitates **Intent Articulation** by extracting both high-level goals and fine-grained intents from user prompts and presenting them as editable UI components. This helps users externalize vague or subconscious preferences into clear, concrete instructions. Second, it supports **Intent Exploration and Management** by making key intent dimensions explicitly manipulable. The system identifies adjustable aspects of an intent—such as tone, specificity, or emphasis—and represents them through dedicated UI components (e.g., sliders, radio buttons). This helps users actively explore new possibilities by directly manipulating the UI. Users can pin and manage their finalized intents while freely adjusting others to explore new directions. Third, to ensure **Intent Synchronization**, INTENTFLOW provides responsive feedback linking intents to output text: hovering over an intent highlights the corresponding segment in the generated draft. This transparency helps users understand how their inputs shape the model's output.

Through a technical evaluation of our core modules with 60 evaluators across diverse writing prompts, we confirmed the robustness and reliability of the system across varied writing contexts. With the

validated technical pipeline, we then conducted a within-subjects study (N=12) comparing IntentFlow with a chat-based LLM interface (baseline), resembling existing systems such as ChatGPT Canvas [1] and Claude Artifact [2]. Our findings show that Intent-Flow helped users express their intents significantly easily and clearly, discover new intents, and understand how their intents are reflected in the output. Interaction logs further revealed a shift in user behavior: from corrective prompting behavior (e.g., "This is not what I want, do it again") while engaging with the baseline system, to active adjustment behavior, more frequently refining their intent within IntentFlow. Moreover, this structured representation of IntentFlow enabled users to remove outdated intents more easily and helped them consolidate a final set of intents that reflected their evolving intents. Participants mentioned that this finalized set of intents could be reused in the future for similar writing tasks, supporting faster and more consistent outcomes.

Ultimately, our findings demonstrate that supporting all four aspects of intent communication enhances users' ability to express and manage their fluid, subconscious, and evolving intents, leading to more effective communication with generative AI. Our work makes the following contributions:

- A **conceptualization of four key aspects of intent communication**—articulation, exploration, management, and synchronization—derived from a systematic literature review of HCI research on intent communication with generative AI, and **design implications** that translate them into actionable direction for system design.
- **IntentFlow**, a system that supports all four aspects of intent communication in LLM-based writing tasks, chosen as a testbed for its iterative and evolving nature.
- **Technical evaluation results** on our system pipeline for intent communication, demonstrating its robustness across diverse writing contexts; and **findings from a user study**, showing its effectiveness in helping users express, manage, and refine their intents more easily and transparently.

## 2 RELATED WORK

### 2.1 Challenges of Aligning User Intent in Natural Language Interaction

Despite the remarkable capabilities of LLMs, several studies have consistently highlighted their limitations in understanding and aligning with user intent. Prior studies have emphasized that LLMs frequently generate responses that misalign with user expectations, sometimes leading to unintended consequences due to inherent biases or prompt sensitivity [1, 14, 25, 33, 72, 76]. One contributing factor lies in the interaction setup itself: natural language, the primary medium for prompting, provides important advantages by lowering the barrier to entry, enabling flexible self-expression without requiring specialized syntax [27, 76]. Yet the very qualities that make natural language accessible also introduce ambiguity and instability, making it difficult for users to articulate their intent with the precise phrasing for LLM and forcing them into cognitively demanding trial-and-error reformulation [19, 31, 43, 45, 63, 76]. In

other words, natural language both empowers users to communicate their intent and, paradoxically, complicates the alignment of that intent with model behavior. The opacity of the intent-to-output connection further worsens this issue: prompts may embed multiple intents, yet users are not informed which parts influenced the output, and small wording changes can yield unpredictable shifts [27, 45, 79]. Moreover, as user intent often evolves during interaction, current chat-based interfaces provide limited support for managing this fluid process. Intents expressed across multiple turns become fragmented within linear conversation histories, making it difficult for users to track, refine, and maintain alignment with model behavior [30, 75]. These challenges collectively highlight the need for new interaction mechanisms that help users articulate their intents, understand how their inputs influence the model, and manage them over time. Our work aims to address these issues by helping users more clearly articulate, refine, and adjust their intents within an LLM system.

### 2.2 Subconscious and Fluid Characteristics of User Intents

While intent alignment itself has been recognized as one of the biggest challenges for LLMs, this issue becomes even more complex when considering the characteristics of human intent during communication. Drawing from interpersonal communication research, we can distinguish user prompts into two types of information: **goals**, which are explicit and relatively stable, and **intents**, which refer to more fine-grained, sometimes subconscious strategies or actions that evolve dynamically throughout a conversation [34, 62]. This distinction is particularly relevant as users often interact with LLMs in ways that resemble human-to-human communication [76]. Research in cognitive task analysis reveals that intent is not spontaneously generated; rather, it emerges from a foundation of knowledge, cognitive processes, and goal configuration [13]. Consequently, user intents often shift during communication [34, 62], and these characteristics of intent complicate how users interact with LLMs, making them a crucial consideration in the design of such interactions. Furthermore, they can be particularly more apparent in creation tasks, such as writing or drawing, due to their iterative nature, requiring continuous reflection and refinement [5, 19, 56, 60, 67]. Recognizing intent as both subconscious and fluid thus highlights the need for interaction designs that can better accommodate the dynamic and situated nature of human communication with LLMs.

### 2.3 Supporting Intent-Aligned Interaction with Generative AI

There are diverse attempts to support users in interacting with generative AIs in ways that better align with their intent. Prior work has explored supporting users to more effectively express their intent, such as providing prompt suggestions to scaffold ideation for text-to-image models [4], allowing easier prompting through direct manipulation [45]. Furthermore, to allow users to better attain the results aligning with their intents, there are attempts to help users understand LLMs' behaviors through a visual programming environment for hypothesis testing [2] or breaking down prompts into smaller subtasks and then aggregating the results [73].

For the cases where user intent is unclear or underspecified, there are attempts to go beyond passively responding to user prompts: proactively retrieving missing information when tasks are ambiguous [52] or asking for follow-up questions to better align with users' overarching goals [72]. For another approach to better understand user intents, IntentGPT identifies intent within user utterances, enhancing the system's ability to respond to varying goals with greater precision [55]. Moreover, to support changes in user intent, dynamic prompt middleware [15], which provides context-specific UI elements to better refine the user prompt, allows users to change their preferred options.

More recently, work has begun to reify intent as manipulable units to support more flexible and reflective workflows. IntentTagger introduces "intent tag"—atomic conceptual units that enable granular and non-linear micro-prompting, supporting intent elicitation and flexible workflows [19]. In a similar vein, AI Instruments embody prompts as reusable interface objects, reflecting multiple interpretations of ambiguous user intents and enabling iterative, non-linear exploration across creative tasks [54]. These approaches suggest new paradigms for structuring intent that go beyond linear prompting.

These efforts also exist in the context specific to writing with LLM as it involves cognitive content-generation that requires intensive and precise intent formulation and expression [20]. One line of research has proposed LLM-based writing-support systems that help users explore and prototype their writing while deciding on their writing intent [30, 64, 78]. There also exist approaches to help users better express their intent by allowing direct manipulation within the text to match their specific writing intent or style [45, 74]. Moreover, there are attempts to make AI-generated suggestions more explicit and manageable, for instance by structuring multiple variations side by side for rapid comparison or surfacing executable edits directly in the document with safeguards for accuracy [36, 53].

While prior work has contributed to improving how users express, reify, and align their intents with generative AI, it has rarely considered the subconscious and fluid nature of user intent as a basis for interaction. Building on this body of work, we conduct a systematic literature review to provide a comprehensive understanding of intent communication with generative AI and present a system that operationalizes these insights by enabling more flexible and reflective interaction with users' evolving intents in writing tasks.

# 3 SYSTEMATIC LITERATURE REVIEW: INTENT COMMUNICATION SUPPORT

In this section, we present the process and results of our systematic literature review (SLR) on existing HCI research regarding intent communication support with generative AI systems. Based on our analysis of collected papers, we identified recurring design approaches for intent communication support and the gaps in current systems. This understanding allowed us to derive a set of design goals for developing systems that holistically support intent communication with generative AI.

## 3.1 Search and Filtering Process

We conducted our SLR following the PRISMA guideline [49], which is a well-known framework for SLR. Our objective was to identify interaction support for intent communication with a Generative AI system in the HCI community.

*3.1.1 Venue.* Because our focus was on how the HCI community has addressed interaction support for intent communication, we limited our scope to the following venues: *the ACM CHI Conference on Human Factors in Computing Systems* (CHI), *the ACM Symposium on User Interface Software and Technology* (UIST), *the ACM Conference on Intelligent User Interfaces* (IUI), *the ACM SIGCHI Conference on Designing Interactive Systems* (DIS), *the ACM Conference on Computer-Supported Cooperative Work and Social Computing* (CSCW), *the ACM Conference on Creativity and Cognition* (C&C), *Proceedings of the ACM on Human-Computer Interaction* (PACMHCI), *the ACM Conference on Conversational User Interfaces* (CUI), and *the ACM Conference on User Modeling, Adaptation and Personalization* (UMAP). Since all these venues are accessible via the ACM Digital Library, we used it as our primary search source.

*3.1.2 Search Keywords and Exclusion Criteria.* We collected papers using the following keyword combination in their title or abstract: 'Generative AI (GenAI)' OR 'Large Language Model (LLM)' AND 'Intent' OR 'Intention'. To capture the period when GenAI and LLM systems became widely adopted as intent-specification interfaces, we restricted the search timeframe to December 2022 (the release of ChatGPT) through August 2025. After conducting an exhaustive search, we applied four exclusion criteria to filter papers that did not align with our SLR goal:

- EC1: We excluded papers where 'intent' was used as an adverb (e.g., 'intentionally'), as this does not relate to user intent communication.
- EC2: We focused on papers where users directly communicate their intent with GenAI or LLM systems. Studies where GenAI mediated intent communication between a user and a third party were excluded.
- EC3: We included only system or framework papers that proposed concrete interaction features, including prototypes. We excluded survey and workshop papers.
- EC4: We focused on papers involving co-creation or generative tasks (e.g., writing, coding, design) where users expressed their wants to the system and received outputs.

*3.1.3 Round 1-3.* We first conducted an extensive search on the ACM Digital Library using the keywords and search period above, which yielded 2,098 papers. We then filtered this set to only include papers from our selected venues, resulting in 275 papers (Round 1). The first author read the abstracts to apply EC1-EC4, reducing the set to 70 papers (Round 2). Subsequently, three authors independently reviewed the full papers based on EC1-EC4, resulting in 32 papers. In addition, we included one relevant study from *the Symposium on Human-Computer Interaction for Work* (CHIWORK) that closely aligned with our criteria, bringing the final set to 33 papers (Round 3). Finally, we meticulously read these 33 papers and extracted the interaction features designed to support intent communication, using these findings as the basis for our analysis. The number of papers at each round is summarized in Table 1.

| Venue | Round 1 | Round 2 | Round 3 |
|---|---|---|---|
| CHI | 144 | 32 | 19 |
| UIST | 7 | 11 | 7 |
| IUI | 28 | 6 | 2 |
| DIS | 30 | 6 | 4 |
| CSCW | 5 | 1 | 0 |
| C&C | 17 | 3 | 0 |
| PACMHCI | 12 | 3 | 0 |
| CUI | 12 | 4 | 0 |
| UMAP | 20 | 4 | 0 |
| CHIWORK | - | - | 1 |
| Sum | 275 | 70 | 33 |

**Table 1: The number of papers selected in each round of our systematic literature review across different venues**

## 3.2 Analysis

Three authors independently read the 33 finalized papers and documented interaction supports for intent communication. We then gathered this data and conducted an open coding process by grouping similar approaches and iteratively refining them into broader themes. Through discussion, our analysis converged on four key aspects for intent communication support in human-GenAI (LLM) interaction. We then returned to the papers, analyzing how each work addressed these aspects to provide a comprehensive understanding of the field's current state.

## 3.3 Findings

*3.3.1 Four Key Aspects of Supporting Intent Communication in Human-GenAI Interaction.* Our analysis revealed four central aspects of how prior systems supported intent communication in human–GenAI interaction:

- **Intent Articulation**: Helping users express their subconscious or underspecified intents more clearly and precisely. This is a convergent process aimed at helping users externalize their vague intent into concrete instructions.
- **Intent Exploration**: Supporting users in discovering new, emerging intents that they may not have been initially aware of. This is a divergent process that encourages users to explore and expand their initial scope.
- **Intent Management**: Providing a structured representation for users to manage the fluid, evolving nature of their intents over time. This offers a persistent and organized way to track, adjust, and revisit their intents throughout the task.
- **Intent Synchronization**: Aligning the user's communicated intents with LLM's output by making transparent how each intent is reflected in the generated output and how modifications of intents trigger corresponding updates. This process enables users to verify whether their intents are realized by the LLM as intended.

These aspects encompass the different yet interdependent ways in which systems can facilitate users in externalizing, expanding, managing, and aligning their intent during interactions with generative AI systems.

*3.3.2 Lack of Comprehensive Support in Prior Work.* Our review of 33 papers shows that existing work rarely addresses all four aspects holistically (Table 2). In fact, none of the reviewed works supported all four aspects together. Many focused primarily on • Articulation(68.75%) and • Exploration(62.50%). This trend reflects the difficulty of expressing underspecified intents through natural language alone in GenAI systems [43, 63] — leading to articulation support through prompt refinement and suggestion, direct object manipulation, or keyword-based intent representation [4, 19, 45]. Moreover, since our review centered on generation tasks, many of the examined works emphasized creative workflows, where exploration was particularly prominent, through mechanisms such as branching, remixing, or probing alternative directions [12, 54]. By contrast, less attention has been paid to • Management(37.5 %) and • Synchronization(31.25 %). Some works introduced structures for organizing intents [51, 71, 78], and others offered mechanisms for making model interpretations more transparent [9, 43, 71], but these were rarely considered together with articulation and exploration. Considering the dynamic nature of intent communication, supporting all four aspects is essential.

*3.3.3 Why All Four Aspects Are Necessary.* Intent communication is a dynamic process: users often begin with a vague intent, uncover subconscious intents along the way, discover new emerging intents, revise earlier ones, or even return to previous directions. We note that the four aspects are not independent but complementary in supporting this process. • Articulation enables users to specify vague intents that are difficult to express. • Exploration helps uncover subconscious or emerging intents that surface during interaction. • Management provides continuity, allowing users to track, revise, and revisit evolving intents over time. • Synchronization ensures that users can verify whether their intents are realized as intended, helping them refine their mental model of the system. This, in turn, strengthens subsequent articulation and exploration. Because these aspects work hand in hand, intent communication should be supported as a cycle in which articulation sparks exploration, exploration produces new intents that require management, management maintains continuity, and synchronization closes the loop. Therefore, we argue that integrating all four aspects together is essential for supporting intent communication in human–GenAI interaction.

## 4 DESIGN IMPLICATIONS

Building on our findings from SLR, we derived four design implications that operationalize the identified aspects of intent communication into actionable directions for system design..

**DI1: Distinguish and externalize goals and intents as the system's visible interpretation.** User prompts often conflate high-level goals with low-level intents [21, 34, 62]. The system should parse prompts into two layers: *stable, overarching goals* and more *fluid, actionable intents*. Extraction should include both explicit and implicit intents, each mapped to system behaviors as subtasks and presented in editable form, enabling users to revise the system's interpretation. This makes the system's interpretation visible (• Synchronization) while giving users a way to refine vague instructions into concrete components (• Articulation).

| Intent Communication Support Aspects | | | Count | Papers |
|---|---|---|---|---|
| • Articulation | • Exploration | • Management | 3 | [4, 19, 78] |
| • Articulation | • Exploration | • Synchronization | 1 | [50] |
| • Articulation | • Management | • Synchronization | 2 | [12, 80] |
| • Articulation | • Exploration | | 10 | [18, 28, 29, 36, 41, 54, 57, 64, 69, 77] |
| • Articulation | • Management | | 1 | [15] |
| • Articulation | • Synchronization | | 3 | [42, 45, 68] |
| • Exploration | • Management | | 3 | [38, 53, 61] |
| • Management | • Synchronization | | 2 | [32, 71] |
| • Articulation | | | 2 | [26, 70] |
| • Exploration | | | 3 | [8, 11, 37] |
| • Management | | | 1 | [51] |
| • Synchronization | | | 2 | [9, 43] |
| Total | | | 33 | |

**Table 2: Combinations of intent communication support aspects identified in our systematic literature review. The table reports the frequency of each combination (Count) and lists the corresponding papers.**

**DI2: Provide easily adjustable exploratory spaces for intents.** Beyond identifying intents, systems should open up exploratory spaces for them, surfacing alternative options such as different tones, structures, or emphases. Crucially, these spaces must be easy to manipulate, such as through direct manipulation [24, 59], so that users can smoothly probe variations. Such mechanisms reduce the effort of exploration and help users uncover latent or subconscious intents. (• Exploration).

**DI3: Support versioning and curation of evolving intents.** Systems should maintain intents across interactions in a structured, persistent form, allowing users to revisit and compare different versions over time. Within this structure, intents that users find effective should be marked or fixed, so that they can be selectively retained or released as needed. Through such selective management, users gradually curate a set of intents that reflect their evolving strategies, supporting both the refinement of vague intents (• Articulation) and the organization of them over time (• Management).

**DI4: Make intent–output connections transparent.** Systems should explicitly link each intent to the parts of the output they influence, making these connections clearly visible. Users should be able to see, for example, which segments of the output correspond to a chosen intent and how modifying it may propagate changes. When alternative intent options are suggested, the system should preview their potential effects on the output, helping users anticipate outcomes before committing. These mechanisms foster transparency and alignment, enabling users to refine their mental models of the system (• Synchronization).

## 5 INTENTFLOW

Based on the design implications, we developed INTENTFLOW, a system for LLM-based writing tasks that supports • Articulation, • Exploration, • Management, and • Synchronization of user intents. We chose writing tasks as our focus because of their prevalence in LLM use and their multi-faceted, evolving, and iterative nature [16, 53]. This makes it an ideal testbed for examining how a system can holistically support intent communication.

## 5.1 Interface & Features

As shown in Figure 2, the interface of INTENTFLOW consists of three main panels: **Chat Panel**, **Intent Panel**, and **Output Panel**. Internally, INTENTFLOW operates through a pipeline of LLM modules, each prompted to perform a specific function, as illustrated in Figure 3. Full system prompts for each module are included in **??** in the Appendix. Note that all prior user queries and outputs are preserved in the chat history and passed to each module at every turn, enabling incremental updates grounded in context. In this section, we present its overall architecture and describe how its key modules and interface components reflect the design implications.

*5.1.1* **Chat Panel (Left, Figure 4a)**. The **Chat Panel** serves as a conversational space where users can enter free-form natural language prompts. The *Entrypoint Chat Module*, which works in this panel not only generates a direct reply in the chat (Figure 4a-e) but also triggers updates in the **Intent Panel**, accompanied by a status update message that indicates which parts are being updated (Figure 4a-d). Depending on the content of the prompt, the module may update goals, intents, or intent dimensions accordingly. This panel allows users to interact freely while simultaneously shaping a structured representation of their intents.

*5.1.2* **Intent Panel (Middle)**. The **Intent Panel** provides a structured representation of the user's goals and intents (**DI1**). It comprises three sections, each dynamically constructed and updated through dedicated LLM modules based on user input, and also editable for user revision.

**Goal Section**: This section captures stable elements of the task—*task goal*, *domain*, and *topic*—which serve as anchors for the writing process [16]. Displaying them explicitly allows users to verify whether the system has correctly understood their overall task context. Any field can be directly revised, and confirmed edits update the goal representation and trigger subsequent modules in the pipeline Figure 3.

**Intent List Section (Figure 4b)**: This section displays a list of discrete intents. *Intent Module* analyzes the prompt and goal to construct this list, capturing not only explicit but also implicit intents, essential or logically required ones that are necessary to
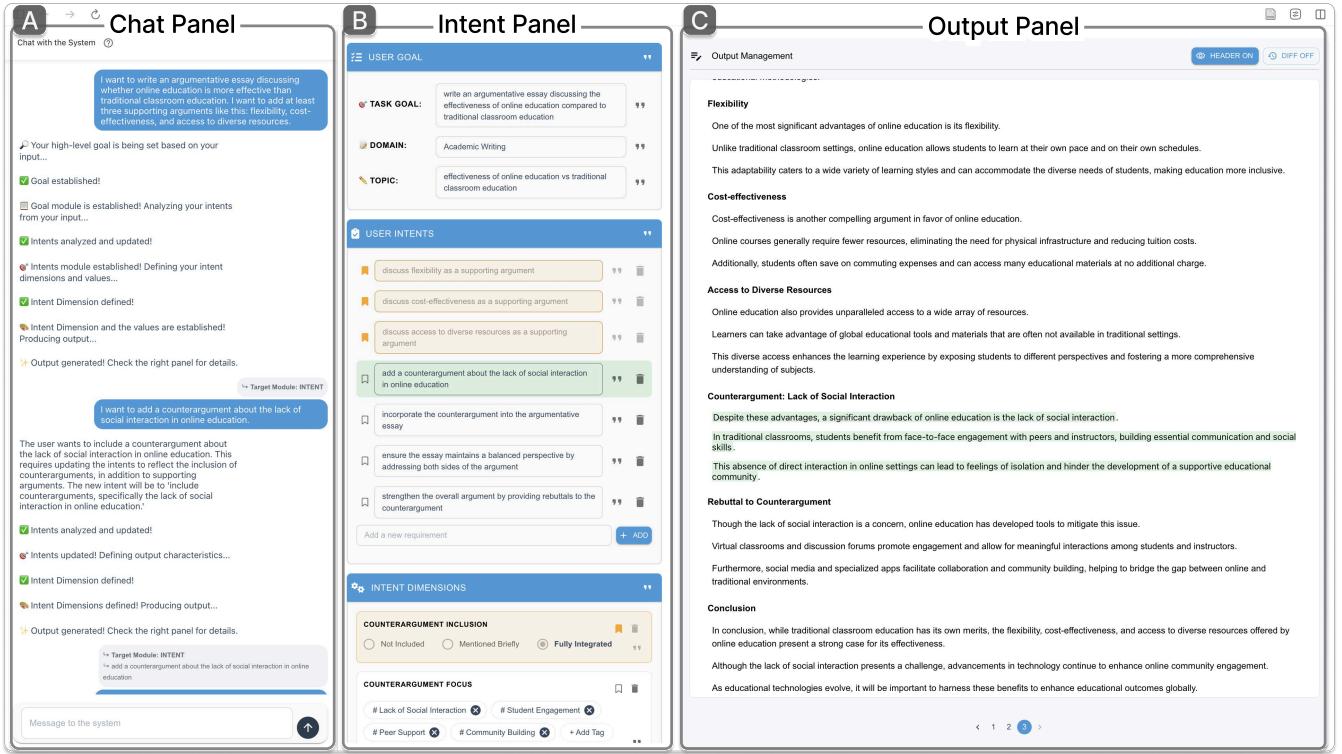
Figure 2: Overall interface of IntentFlow. It is split up into three sections: (A) Chat Panel, (B) Intent Panel, and (C) Output Panel.



Figure 3: System pipeline of IntentFlow. The *EntryPoint Chat Module* first interprets the prompt and coordinates three intent-processing modules: the *Goal Module, Intent Module,* and *Dimension Module.* These modules extract and structure the user's goal and intents, which are then used by the *Output Module* to generate text. The *Linking Module* associates each intent with corresponding segments in the generated output. Each module box shows its LLM model name, with a gray-text description below or beside it.

carry out the task. The list is also designed to reflect how the LLM decomposes the prompt into fine-grained subtasks, making the system's interpretation visible (**DI1**). Each item is editable, allowing

users to revise through direct text editing, delete (Figure 4b-b), or add new ones (Figure 4b-c). Users can click the 🔖 button to keep the intents they are satisfied with and wish to retain (Figure 4b-a)

(**DI3**). The kept intents are persistently reflected in subsequent output generations. Clicking the button again cancels the keep.

**Intent Dimension Section (Figure 4c)**: Each intent is further decomposed into adjustable dimensions. *Dimension Module* analyzes the user's prompt, goal, and intent list to determine which dimensions are relevant, what UI format is most appropriate (e.g., sliders (Figure 4c-b), radio buttons (Figure 4c-a), or tags (Figure 4c-c)), and what initial values should be assigned. In addition, explanations of each dimension value, including its meaning and how selecting it would affect the output, are generated by *Preview Module* and shown as a popover when the user hovers over the value (Figure 4c-d) (**DI4**). Through this, users can easily adjust and explore with alternatives during the task (**DI2**).

*5.1.3* **Output Panel (Right, Figure 5)**. The **Output Panel** displays the writing output, which is generated by the *Output Module* based on the user's prompt, along with the full state of the **Intent Panel**. After each generation, *Linking Module* is immediately invoked to identify which parts of the output correspond to specific intents and dimension values. As a result, users can hover over items in the **Intent Panel** to see the relevant parts of the output highlighted (Figure 5-a), enabling more transparent traceability between expressed intent and generated text (**DI4**). The output panel also maintains a version history: each output is saved with its associated goals, intents, and dimensions. Users can browse past versions, inspect changes using `DIFF ON` toggle, and rollback to any prior version by clicking `MAKE THIS LATEST PAGE` button. This action restores the selected output and its associated intent structure as the current state, allowing users to continue their work from a preferred point and supporting the curation of an evolving set of intents over time (**DI3**). To enhance readability, the output is presented in a structured format with optional section headers, which users can toggle on or off (Figure 5-d) according to their preferences.

## 5.2 Interaction Flow

To illustrate how INTENTFLOW supports iterative intent communication during the writing process, we present an example user scenario. Emma, a graduate student in environmental science, is writing an article on photosynthesis for a science magazine. Unsure where to begin or how to frame her ideas, she turns to INTENTFLOW to help surface, organize, and refine her writing intents.

*5.2.1* **From vague prompts to structured intents and linked draft.** Emma begins with a high-level prompt: *"Write a scientific and concise article on photosynthesis"*. *Entrypoint Chat Module* interprets this as the starting input and initiates the pipeline. *Goal Module* extracts stable task elements—task goal, domain, and topic—providing a high-level reference. *Intent Module* then infers both explicit and implicit intents needed to achieve Emma's goal, which are shown as editable items in the Intent Panel (Figure 4b). Next, the *Dimension Module* generates adjustable dimensions for intents (e.g., length, focus, tone) with suitable UI controls and initial values (Figure 4c). Based on this structured representation, *Output Module* generates a draft, and *Linking Module* connects each intent and dimension value to specific parts of the text, allowing Emma to hover over them and immediately see their influence on the output

(Figure 5-a). Through this, Emma can see how her vague prompt is interpreted into concrete intents and dimensions, and how they shape output.

*5.2.2* **Exploring Intents through Direct Manipulation.** After seeing how her vague prompt is decomposed into intents and dimensions, Emma starts to explore alternatives by directly manipulating them. She notices that the article's current focus on *'Key concepts'* may be too generic for her purpose. Using the *'Article Focus'* control, she hovers over each item to see how it would affect the output and switches to *'In-depth Analysis'*, which expands the explanations of processes and mechanisms (Figure 4c). Wanting to bring in her disciplinary perspective, Emma adds a new intent in the **Intent List**: *"Highlight the ecological importance of photosynthesis in sustaining ecosystems"*. This triggers *Dimension Module* to surface a new dimension, *Ecological Emphasis* with a hashtag-style options, *'#Ecosystem Sustainability'*, *'#Climate Impact'*. Seeing the result, Emma realizes that this better captures the direction she wanted, helping her uncover an intent she had not explicitly considered at the outset.

*5.2.3* **Refining intents through dual prompting.** While direct manipulation helps Emma explore alternatives, some needs go beyond manipulating dimensions or editing text fields. Sometimes she wants to suggest broader changes or articulate ideas that are not yet captured. To address this, INTENTFLOW supports two complementary prompting mechanisms.

With **Chat-based Prompting**, Emma can type free-form prompts at any point. For example, after reviewing the output, Emma realizes that the article is too technical for a general audience, and writes: *"I want to make the article easier for readers without a science background to understand, while keeping the academic tone"*. *Entrypoint Chat Module* interprets this as keeping the same goal but adding a new intent. It therefore calls the *Intent Module* and it adds a new intent, *'Use simpler terminology to explain the scientific concepts of photosynthesis'*, and *Dimension Module* adds corresponding dimensions, *'Terminology Complexity'* and *'Target Audience'*. A new draft is generated with these changes, while the previous version is kept in history.

With **Intent-based Prompting**, Emma can refine a specific intent without rewriting the whole prompt, using the 🗩 button next to each element in the **Intent Panel**. For example, Emma notices that the introduction of the article feels too brief to convey the topic's significance. Instead of rewriting the entire prompt, she clicks 🗩 button next to the intent *"Introduce the topic concisely"* and enters a targeted instruction: *"Add a bit more background about why photosynthesis is important for the environment."*. *Intent Module* updates the selected intent based on this instruction and slightly expands the introduction section in the output. Meanwhile, the rest of the article remains unchanged. She hovers over the updated intent to check which part of the output has changed, and clicks `DIFF ON` button to easily see the exact modifications.

*5.2.4* **Curating and Managing Intents over Time.** As Emma continues experimenting and refining, INTENTFLOW stores each draft together with its associated intent state. When she encounters an intent she finds satisfying, she can mark it by clicking the 🔖 button next to it, which keeps the intent persistent across subsequent
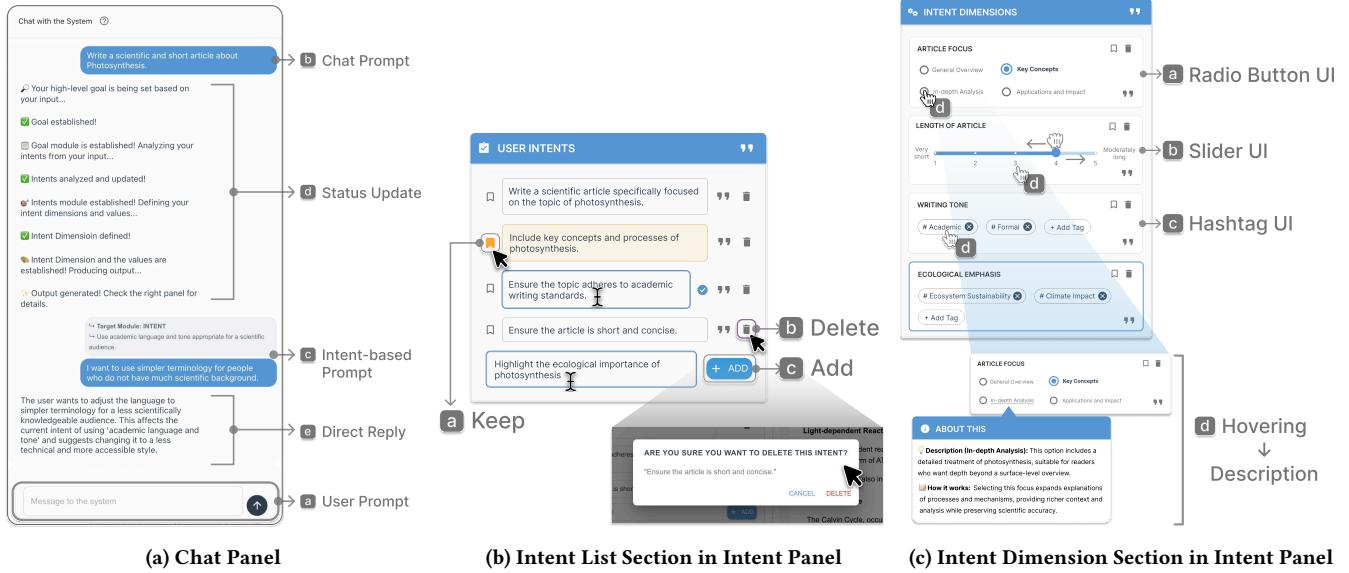
**(a) Chat Panel**          **(b) Intent List Section in Intent Panel**          **(c) Intent Dimension Section in Intent Panel**

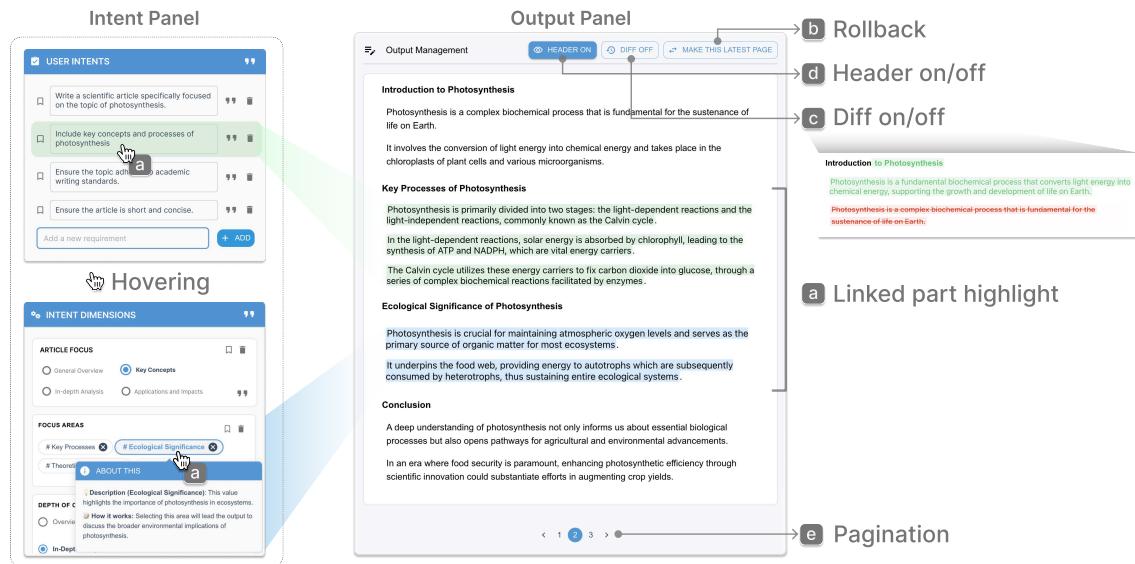**Figure 4: Chat Panel and Intent Panel**



**Figure 5: The Output Panel displays generated output and links between the output and the Intent Panel. (a) Hovering highlights linked parts in green (user intents) or blue (intent dimensions). (b) Users can roll back to any prior version, which brings the selected output and its associated intents to the latest position in the workflow. (c) Diff view compares old and new outputs. (d) Section headers can be toggled. (e) All outputs are paginated, allowing users to browse previous outputs.**

interactions. She can also unkeep it by clicking the button again. Emma can browse earlier drafts and their associated intents, compare them using ⟲ DIFF ON toggle, and restore a preferred version with ⇄ MAKE THIS LATEST PAGE button. Through this, she can continue building on a version that aligns more closely with her evolving intents. This allows her to curate a set of intents that increasingly reflects the directions she is most satisfied with.

### 5.3 Implementation

INTENTFLOW was implemented as a web-based application with React for the frontend and Flask for the backend. The backend leveraged the OpenAI API to handle language model functionalities, with each module responsible for a specific role in the pipeline, as illustrated in Figure 3. To improve responsiveness in real-time interactions, we tested our pipeline using a smaller model and observed

comparable performance in the *Goal Module*, *Intent Module*, and *Dimension Module*. Therefore, we used `gpt-4o-mini-2024-07-18` for those modules. For all other modules, we used `gpt-4o-2024-08-06`. For the *Linking Module*, which processes all intents and dimension values against the output, we employed multiprocessing to reduce loading time. Full prompts for each module are provided **??** in the Appendix.

## 5.4 Technical Evaluation

We conducted a technical evaluation to assess the quality of outputs generated by our pipeline (Figure 3), which extracts task goals (*Goal Module*), a set of user intents (*Intent Module*), intent dimensions with UI components and initial values (*Dimension Module*), and intent-to-output links(*Linking Module*). We focused our evaluation specifically on these four modules, as they were newly introduced to support the intent communication process. Our goal was to verify whether each module operated accurately and appropriately across diverse writing contexts. We focused on one-step generation outputs, as iterative refinement can be examined in the user study.

*5.4.1 Evaluation Setup.* We selected 12 representative prompts—two from each of the six writing contexts defined by Lee et al. [39]: academic, creative, journalistic, personal, professional, and technical. For each prompt, we generated a full pipeline output and created a corresponding survey form for evaluation. To reduce ambiguity and individual variance, the survey used binary (Yes/No) questions, complemented by an optional free-form section. We recruited a total of 60 evaluators (5 per prompt) via Prolific [3]. To ensure evaluation quality, we excluded evaluators who had never used LLMs, lacked experience with the given writing task, or reported unfamiliarity with the given writing topic. Each evaluator was compensated with £5 for completing a 30-minute evaluation task. Detailed information about the prompts, writing tasks, and topics used in the evaluation is provided in the Appendix A.2.2.

*5.4.2 Evaluation Criteria.* We designed evaluation questions as follows. For clarity, we refer to each question by its title below; full phrasing is included in Appendix A.2.1.

[*Goal Module* − **Q1. Goal Alignment**]: Whether the extracted goal, domain, and topic appropriately reflect the user's overall objective.

[*Intent Module*]

- [**Set of Intents** − **Q2. Completeness**]: Whether the intents cover all key aspects of the prompt.
- [**Set of Intents** − **Q3. Distinctiveness**]: Whether the intents are distinct without redundancy.
- [**Individual Intents** − **Q4. Relevance**]: Whether each intent is relevant to the prompt.

[*Dimension Module*]

- [**Q5. Relevance**]: Whether each intent dimension is relevant to the prompt.
- [**Q6. UI Appropriateness**]: Whether the UI component fits the nature of the dimension.
- [**Q7. Value Appropriateness**]: Whether the given value is appropriate for the prompt.

---
[3]https://prolific.com

| Evaluation Criteria | "Agree"(%) |
|---|---|
| Q1. Goal Alignment | 95.00 |
| Q2. Set of Intents – Completeness | 95.00 |
| Q3. Set of Intents – Distinctiveness | 86.67 |
| Q4. Individual Intents – Relevance | 94.08 |
| Q5. Intent Dimension – Relevance | 86.56 |
| Q6. Intent Dimension – UI Appropriateness | 86.78 |
| Q7. Intent Dimension – Value Appropriateness | 86.14 |
| Q8. Intent-to-Output Linking – Link Accuracy | 94.04 |

**Table 3: Evaluation results for each question in the technical evaluation. Values indicate the percentage of "Agree" responses aggregated across all prompts and participants.**

[*Linking Module* − **Q8. Link Accuracy**]: Whether the highlighted output reflects the corresponding intent.

Each question was presented with corresponding outputs or UI components, and participants were asked to evaluate their appropriateness with respect to the given prompt.

*5.4.3 Technical Evaluation Results.* We computed the proportion of "Agree" responses for each question across prompts and participants. As shown in Table 3, most components received high scores, with over 85% positive responses in all categories. We also analyzed qualitative feedback from evaluators. Several participants noted that when intent dimensions were presented without clear descriptions of value meanings, it was difficult to judge their appropriateness. For example, when the "Formality Level" dimension was presented as a slider with an initial value of 4, participants reported that it was hard to judge how formal the value 4 actually was. This likely contributed to the lower agreement in the *Intent Dimension* category since our technical evaluation did not provide a hover-based explanation for each dimension value. This underscores the importance of providing descriptive explanations for dimension values (Figure 4c-d).

## 6 USER STUDY

To understand how users engage in intent communication during writing tasks and how our system supports this process, we conducted a within-subjects user study comparing INTENTFLOW with a baseline LLM interface. The order of system and task was counterbalanced across participants. We aimed to examine both the behavioral patterns of intent communication and users' subjective experiences across the two conditions. We framed our inquiry around the following research questions:

- **RQ1**: How do users interact with INTENTFLOW in intent communication processes?
- **RQ2**: How effective is INTENTFLOW in supporting intent communication?

To examine the effects of intent communication support, we implemented a baseline system that closely resembled INTENTFLOW in overall structure, but excluded the **Intent Panel** and the *Linking Module* that support intent communication. All other functionalities—such as prompt-based generation, pagination, header on/off view, diff view, and version control using ⟲ MAKE THIS LATEST PAGE button—were preserved in the baseline to ensure comparability. The baseline used the same generation model as INTENTFLOW, `gpt-4o-2024-08-06`.

The baseline interface followed the design of recent chat-based LLM tools such as ChatGPT and Claude Artifacts, where users interact through free-form chat prompts, and the system responds conversationally, generating writing outputs in a separate panel. The screenshot of the baseline interface can be found in Figure 11 in the Appendix.

## 6.1 Participants

We recruited 12 participants (8 male, 4 female; age $M = 25.50$, age $SD = 2.75$) through online recruitment posting at our university community platforms. All participants had prior experience using large language models (LLMs) and were fluent in English reading and writing, as the study tasks were conducted entirely in English. During the recruitment, we administered a pre-survey to collect information about participants' LLM experience, familiarity with the target writing topics (e.g., Doppler Effect), and general writing background. We used this information for pre-screening, excluding those who lacked sufficient topic knowledge or writing experience, as these factors could limit meaningful engagement and evaluation quality. Participants received 30,000 KRW (around 22 USD) for the 90-minute session. Additional details, including participants' LLM experience and writing experiences, are provided in the Appendix A.3.3.

## 6.2 Tasks

Each participant completed two writing tasks, which were randomly paired with the counterbalanced study conditions. We selected writing tasks commonly used in prior HCI studies [17, 30, 53] and added short scenarios to each task to introduce more nuanced and multifaceted considerations. The first task involved **writing a social media post** explaining a scientific concept (e.g., the Doppler Effect) for a general audience with little scientific background. The second involved **writing a professional email** applying for a personal secretary position for a well-known individual outside the participant's domain. Through multiple rounds of pilot testing, we refined the scenarios to ensure that the two tasks were similar in complexity and involved a comparable level of multifaceted consideration. To guide the output scope, we set a flexible length guideline of around half an A4 page. We intentionally did not specify an exact word count to prevent users from fixating on meeting a numerical target, which could distract from focusing on the given task scenario and goal. Importantly, as our study centered on intent communication, participants were encouraged to explore and refine intents rather than produce polished drafts. To maintain this focus, we deliberately excluded manual editing of outputs, since pilot sessions showed it often shifted attention to surface-level fixes rather than intent communication. Task descriptions were always accessible during the study, though participants were instructed not to copy and paste their content into the system. Full task descriptions are available in Appendix A.3.1.

## 6.3 Procedure

The study was conducted either in person or online via Zoom [4], depending on participant availability. Each study session lasted approximately 90 minutes and followed a fixed protocol Figure 6. After

a brief introduction and consent process (5 minutes), participants were given a tutorial for the first system (8 minutes). They then completed the first writing task using that system (20 minutes), followed by a post-survey and an annotation activity (10 minutes). To gather each participant's intent communication actions, we brought them back to the system after the survey to review their interactions and annotate the purpose of each input. These annotations were used in our analysis of how users engaged in intent communication, as described in more detail in subsection 6.4. This process was then repeated for the second system and task. At the end of the sessions, participants completed a semi-structured interview (10 minutes) reflecting on their experience with both systems, including their preferences, their strategies for expressing and adjusting intent, and the usefulness and challenges of each system.

## 6.4 Measures

To evaluate how participants engaged in intent communication and how effectively each system supported this process, we collected and analyzed data from three sources: self-report ratings, interaction logs with annotated user actions, and post interviews.

*6.4.1 **Self-Report Ratings**.* Participants completed a series of 7-point Likert scale questions (1: Strongly Disagree / 7: Strongly Agree) after each writing task with each system. The questions were designed to evaluate how effectively each system supported the intent communication process. To guide the question design, we referred to self-rating measures used in prior studies on LLM-based writing support systems [30, 53, 73] and adapted them to capture multiple aspects of intent communication support. The full set of questions is as follows.

(1) **M1. Intent Expression — Ease**: *"I could easily express my intent to the system."*
(2) **M2. Intent Expression — Clarity**: *"The system helped me express my intent clearly."*
(3) **M3. Intent Discovery**: *"The system helped me recognize or discover additional intents that I had not explicitly considered at the start."*
(4) **M4. Transparency**: *"The system helped me see how each of my intents influenced the output generation."*
(5) **M5. Understanding**: *"I understood how each of my intents was reflected in the output."*
(6) **M6. Think-Through**: *"The system helped me think what kinds of intents I would want to complete the task goal, and how to complete the task."*
(7) **M7. Intent Adjustment — Ease**: *"I was able to adjust my intent to achieve the output aligned with my task goal."*
(8) **M8. Intent Elaboration**: *"The intents I created/kept were specific, detailed, and well-articulated."*
(9) **M9. Intent Match**: *"The system helped me obtain intents and an output that better matched what I wanted."*
(10) **M10. Draft Quality**: *"I felt like communicating intents using this workflow will help me have a better final draft."*
(11) **M11. Intent Reusability**: *"I would reuse the intents I created/kept in future similar tasks."*

We also measured the participants' perceived task load using NASA-TLX [22], using a 7-point Likert scale.
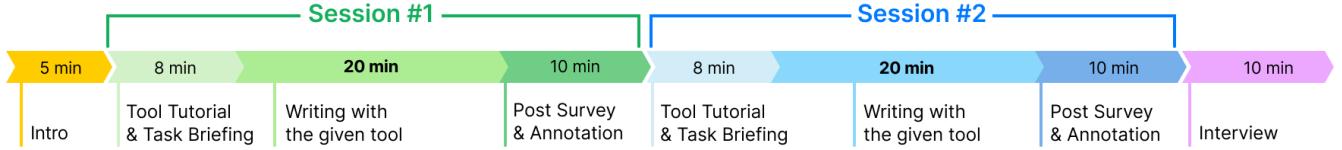
**Figure 6: User study procedure. After a brief 5-minute introduction, each participant completed two sessions. The first 8 minutes were for the system tutorial and task briefing. The user had 20 minutes to write using the given system, followed by a 10-minute post-survey and annotation. At the end of both sessions, there was a 10-minute interview.**

*6.4.2 Interaction Log and Annotated User Actions.* We categorized all user inputs into four types of intent communication actions: **Add**, **Delete**, **Correct**, or **Adjust**. This categorization was informed by prior work that categorized how users interact with LLMs to express or refine their intents [3, 10, 33, 58]. **Add** refers to introducing a new intent not previously expressed. **Delete** refers to removing a previously expressed intent; for instance, a user may have initially requested a summary at the end but later decided to omit it. **Correct** captures cases where a user re-communicates an earlier intent due to a misunderstanding by the system; for example, the user asks for a concise style, but the system generates something too verbose, prompting the user to restate the intent more clearly. Finally, **Adjust** involves modifying an existing intent in terms of degree or nuance, for example, slightly increasing the level of formality or adjusting the specificity of a detail without changing the underlying intent altogether.

To examine how participants engaged in these different types of intent communication with the system throughout the writing process, we asked them to annotate their own inputs after completing each task and survey. Participants were redirected to the system and reviewed all of their interactions. For each input—whether a free-form prompt or a direct edit—they were asked to indicate which of the four action types best described their intent at the time. For interface-level interactions in INTENTFLOW that clearly corresponded to a specific action type, such as adding or deleting hashtags, adjusting sliders, or adjusting radio buttons, we did not ask for annotation. This process allowed us to systematically analyze the frequency and distribution of different intent communication behaviors across both systems.

*6.4.3 Post Interviews.* After completing both sessions, participants completed a semi-structured interview to reflect on their experiences with both systems. The goal of these interviews was to gather qualitative data on users' intent communication strategies and their views on the usability of intents. Specifically, we asked how they approached revising or discarding intents during the writing process, and how their strategies may have differed across the two systems. In addition, participants were asked to reflect on whether and how the final set of intents they created or kept during the task might be reused in future writing tasks. The interview data were used to complement our quantitative analysis by providing contextual insights into interaction patterns and feature usage, as well as participants' mental models of intent communication workflows.

## 7 RESULT

For all statistical comparisons, we first conducted Shapiro-Wilk tests to examine the normality of each measure. Based on the results ($p = .05$), we used paired t-tests for normally distributed data and Wilcoxon signed-rank tests otherwise. Overall, our study shows that INTENTFLOW enabled users to communicate their intent more effectively, reducing repetitive correction cycles and supporting more deliberate refinement throughout the writing process. Participants rated INTENTFLOW significantly higher than the Baseline across all 11 individual items (M1–M11; p < .05 for each), confirming its effectiveness in supporting various aspects of intent communication ( Figure 7). In terms of cognitive effort, participants reported significantly lower workload when using INTENTFLOW, as reflected in NASA-TLX scores (INTENTFLOW: $M = 15.67$, $SD = 4.01$; Baseline: $M = 19.67$, $SD = 4.50$; $p = .004$, $t = 2.02$), as shown in Figure 8a. Further analysis of individual NASA-TLX dimensions revealed significantly lower ratings for both **Effort** (INTENTFLOW: $M = 4.25$, $SD = 1.22$; Baseline: $M = 5.08$, $SD = 0.90$; $p = .048$, $W = 37.00$) and **Frustration** (INTENTFLOW: $M = 2.75$, $SD = 1.49$; Baseline: $M = 3.83$, $SD = 1.75$; $p = .034$, $t = 2.03$). The full breakdown of ratings across all six NASA-TLX dimensions is shown in Figure 8b.

From interaction logs, we also observed that participants engaged in fewer **Correct** actions and fewer • **Rollback** operations with INTENTFLOW, while showing greater use of **Adjust** and **Delete** actions. These patterns imply a transition from repetitive correction toward more deliberate and fluid intent refinement.

### 7.1 RQ1. How do users interact with the system in intent communication processes?

We analyzed system interaction logs and categorized each user input into one of four intent communication action types—**Add**, **Delete**, **Correct**, or **Adjust**—to better understand how users engaged in intent communication with the system. As in shown in Figure 9 and Figure 10, participants performed significantly fewer **Correct** when using INTENTFLOW compared to the Baseline (Action Count: INTENTFLOW: $M = 0.50$, $SD = 0.67$; Baseline: $M = 4.33$, $SD = 2.64$; $p < .001$, $t = -4.81$). The number of **Add** was comparable across systems (Action Count: INTENTFLOW: $M = 3.75$, $SD = 2.05$; Baseline: $M = 4.42$, $SD = 1.38$; $p = 0.296$, $t = -1.10$), while both **Adjust** (Action Count: INTENTFLOW: $M = 4.50$, $SD = 2.97$; Baseline: $M = 1.17$, $SD = 0.72$; $p = .005$, $t = 3.46$) and **Delete** (Action Count: INTENTFLOW: $M = 1.00$, $SD = 1.13$; Baseline: $M = 0.17$, $SD = 0.39$; $p = .031$, $W = 21.00$) occured substantially more often in INTENTFLOW.
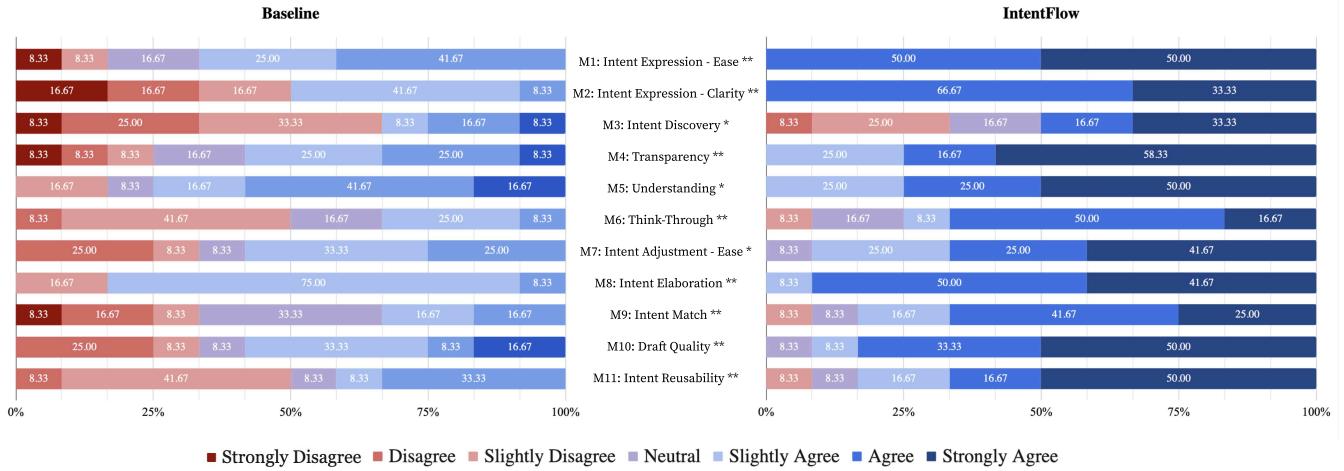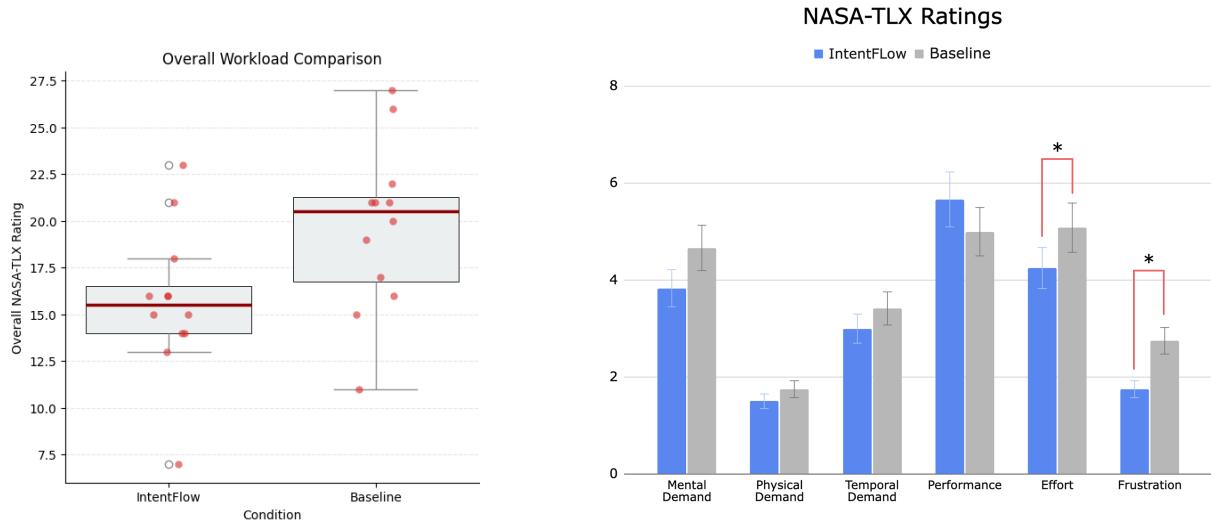
Figure 7: This is a distribution of participants' ratings on intent communication experience. (*$p < .05$, **$p < .01$)



(a) Comparison of overall NASA-TLX ratings between IntentFlow and Baseline. The dark red line indicates the median score per condition.

(b) Comparison of NASA-TLX ratings between IntentFlow and Baseline. (*$p < .05$).

Figure 8: NASA-TLX ratings between IntentFlow and Baseline.

This suggests that participants using IntentFlow spent less effort restating misaligned intents and instead engaged more in adjusting or removing intents they no longer needed. In contrast, the Baseline produced far more **Correct** actions: users often could not see how their intents were preserved across turns or maintain them in a structured way, leading them to restate the same intent repeatedly as the system failed to reflect it over time. Five participants (P3, P6, P9-11) specifically reported that in the Baseline, their previously expressed intents often felt disregarded when they issued a new prompt, as the system seemed to attend to the most recent instruction. Specifically, P6 described, *"In the B (Baseline), I felt like*

*the context wasn't being maintained. Unless I repeated my earlier instructions, the system would just respond to the most recent prompt, so I had to keep saying the same things."* This reflects insufficient support for • Management and • Synchronization.

By contrast, IntentFlow preserved users' intents in the **Intent Panel**, with kept intents remaining persistent across turns. Moreover, because users could review their goals, intent lists, and dimensions before output generation, they had better visibility into how their intents were interpreted by the system. As a result, intents were rarely lost and the need for **Correct** was much lower. This
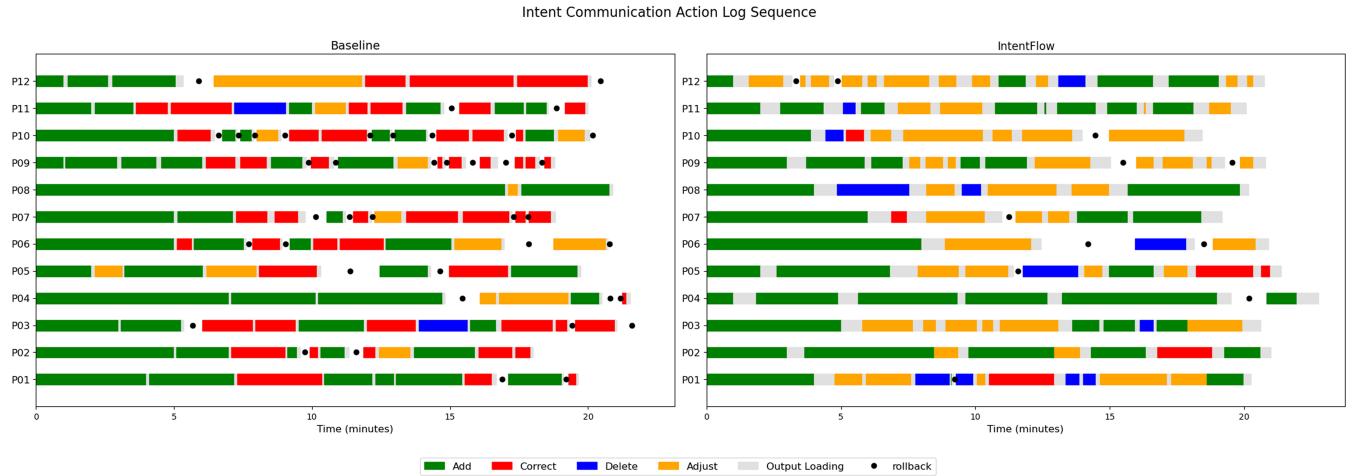
Intent Communication Action Log Sequence



**Figure 9: This figure illustrates the action log sequences for 12 participants in each condition: Baseline and INTENTFLOW.**
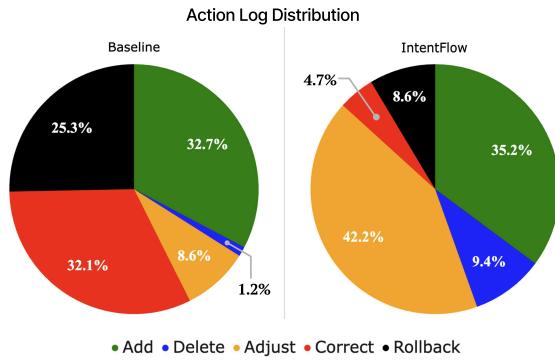


**Figure 10: Percentage distribution of action types aggregated across all participants. Each chart reflects the proportion of actions by type in INTENTFLOW and Baseline.**

shows that • Management and • Synchronization in INTENTFLOW reduced redundancy in intent communication.

Participants performed significantly more **Adjust** in INTENT-FLOW, engaging in targeted revisions of parts they wanted to change while preserving those they were satisfied with. The linking feature helped users see how each intent was reflected in the output. Participants (P2–7, P9–12) noted it as useful, and five (P4–6, P10–11) emphasized that it made adjustment easier. As P5 explained, *"With the linking feature, I could easily see which parts I wanted to keep and which parts to change, and that made it much easier to revise"* Six participants (P1–3, P6, P9–10) also found intent dimensions effective for adjusting their intents. For example, P1 noted that *"dimensions and suggested values from my prompt gave me a clear reference point for exploring and refining my intents."* These findings show how • Exploration, • Articulation, and • Synchronization were actively supported in INTENTFLOW, enabling more effective intent communication.

Moreover, we found that INTENTFLOW better supported negative prompting, **Delete** of previously stated intents. Participants performed **Delete** action more frequently in INTENTFLOW than in the

Baseline. P6 noted that *"In the B (Baseline), saying something like 'don't do this anymore' in the chat felt awkward and left visible traces in the conversation history, which made intent removal cognitively uncomfortable. But in A (INTENTFLOW), intents were separated from the conversational thread and externalized as editable components, so deletion felt much more natural and manageable."* This highlights how INTENTFLOW supported • Management of evolving intents by making their removal more natural.

We also observed that users' use of • **Rollback** feature—returning to a previous output and making it the current version using [↩ MAKE THIS LATEST PAGE] button—differed between the two systems. • **Rollback** occurred significantly more often in the Baseline than in INTENTFLOW (Action Count: INTENTFLOW: $M = 0.92$, $SD = 0.;79$ Baseline: $M = 3.42$, $SD = 2.50$; $p = .003$, $t = −3.80$) with different usage patterns. In the Baseline, • **Rollback** was mostly paired with **Correct** usually after the system failed to produce the intended output, prompting users to return to a previous version and restate their intent. Two participants (P9, P10) explicitly noted frustration, as changes often overwrote parts they had liked. As P9 said, *"In the B (Baseline), I asked to change only a specific part, but the whole output changed, so I often had to rollback and reissue the same prompt.."* In contrast, • **Rollback** in INTENTFLOW was more frequently observed alongside **Adjust** as part of the exploratory process. In particular, four participants (P7, P9–10, P12) said that they tried different intent adjustments, compared outputs, and returned to a version that better aligned with their original intent. For example, P10 said, *"...in the A (INTENTFLOW), having the intent dimensions surfaced in the UI made it easier to explore and adjust. I tried different variations and rolled back to the output I liked most."* This indicates that while • **Rollback** in the Baseline reflected breakdowns in intent communication, in INTENTFLOW it supported a more exploratory and controlled process of refining fluid intents, thereby supporting • Exploration and • Management.

## 7.2 RQ2: How effective is the system in supporting intent communication?

We present participants' responses to 11 survey items (M1–M11), which assessed various aspects of intent communication(Figure 7). We organize the results into three main findings.

### 7.2.1 Discovering and Articulating Underlying Intents Clearly and Easily (M1-3, M8). Participants rated INTENTFLOW significantly higher than Baseline in multiple aspects of intent expression: ease (M1: INTENTFLOW: $M = 6.50$, $SD = 0.52$; Baseline: $M = 4.75$, $SD = 1.55$; $p = .004$, $W = 45.00$), clarity (M2: INTENTFLOW: $M = 6.33$, $SD = 0.49$; Baseline: $M = 3.58$, $SD = 1.78$; $p = .001$, $W = 78.00$), and Discovery of additional intents (M3: INTENTFLOW: $M = 4.92$, $SD = 1.92$; Baseline: $M = 3.58$, $SD = 1.92$; $p = .018$, $t = 2.40$). They also rated the elaboration level of the expressed intents significantly higher (M8: INTENTFLOW: $M = 6.33$, $SD = 0.65$; Baseline: $M = 4.75$, $SD = 1.60$; $p = .002$, $W = 66.0$).

Regarding this, most of the participants (P2-8, P10, P12) specifically noted that the **Intent Panel** helped them better surface and articulate their underlying intents, which might have otherwise remained implicit in free-form prompting. As P8 described, *"When I write, there are often a lot of things I want to express, but it's hard to capture them all. Even if I think I've included everything, I often realize something's missing when I reread it. But seeing the list of intents in A (INTENTFLOW) felt like using a checklist; it helped organize what was in my head and made me write more thoroughly."*

In addition, five participants (P3, P7, P9-10, P12) highlighted that intent dimensions helped them recognize additional factors they hadn't initially considered. P10 mentioned that *"When I usually interact with an LLM, I struggle to come up with ideas for how to make the output better. But with this system, the dimensions gave me options I could try out—changing values while keeping the overall structure, and seeing what parts changed or were added or removed. Through that process, I discovered new intents to improve my output and refined it accordingly. In the end, I felt I achieved a level of output I couldn't have reached without this system."* This demonstrates how INTENTFLOW effectively supported • Articulation and • Exploration in intent communication.

### 7.2.2 Understanding How Intents Are Reflected in the Output (M4-6). Participants rated INTENTFLOW significantly higher than the Baseline in their understanding of how intents were captured and reflected in the generated outputs. Ratings for output transparency (M4: INTENTFLOW: $M = 6.33$, $SD = 0.89$; Baseline: $M = 4.50$, $SD = 1.78$; $p = .007$, $t = 2.93$) and understanding of intent reflection (M5: INTENTFLOW: $M = 6.25$, $SD = 0.87$; Baseline: $M = 4.50$, $SD = 1.37$; $p = .047$, $t = 1.84$) were both significantly higher in INTENTFLOW. Ten participants (P2–7, P9–12) specifically highlighted the usefulness of the linking feature, which visually connected individual intents or intent dimensions to specific parts of the output. They explained that this made it easier to see the consequence of each intent-related decision and form a mental model of how their input influenced the output. As a result, they could better identify mismatches between their intents and the system's output, and decide how to refine them accordingly.

Five participants (P2, P5–6, P9, P12) further noted that having their goals, intents, and dimensions explicitly represented helped them verify whether the system had correctly understood their context, making the system seem to understand their input better. In particular, P2 and P6 said the Goal section reassured them whether the system had grasped the overall direction for their task. Ratings for support in helping participants think through what kinds of intents were needed to accomplish their task were also significantly higher in INTENTFLOW (M6: INTENTFLOW: $M = 5.50$, $SD = 1.24$; Baseline: $M = 3.83$, $SD = 1.19$; $p = .006$, $W = 36.00$). This suggests that INTENTFLOW's linking feature helped users reason more effectively about what adjustments were needed to meet their task goals. Regarding this, P6 shared: *"By seeing how each intent was linked to parts of the output, I could clearly understand why the model generated the text that way. It also made it much easier to figure out what and how I needed to revise."*

### 7.2.3 Adjusting and Elaborating Intents Toward Output Alignment (M7-11. Participants rated INTENTFLOW significantly higher than the Baseline in supporting intent adjustment and elaboration toward better output alignment. Survey responses showed higher ratings for ease of intent adjustment (M7: INTENTFLOW: $M = 6.00$, $SD = 1.04$; Baseline: $M = 4.25$, $SD = 1.60$; $p = .021$, $W = 40.00$), elaboration (M8: INTENTFLOW: $M = 6.33$, $SD = 0.65$; Baseline: $M = 4.75$, $SD = 1.60$; $p = .002$, $W = 66.0$), and alignment between their intended direction and the final output (M9: INTENTFLOW: $M = 5.67$, $SD = 1.23$; Baseline: $M = 3.83$, $SD = 1.59$; $p = .003$, $W = 63.00$). They also rated INTENTFLOW as more helpful in producing higher-quality drafts (M10: INTENTFLOW: $M = 6.25$, $SD = 0.96$; Baseline: $M = 4.42$, $SD = 1.83$; $p = .004$, $W = 55.00$).

Interaction log data reinforced these findings. Users performed considerably more **Adjust** actions and fewer **Correct** actions in INTENTFLOW, indicating a shift away from repetitive corrections toward more proactive refinement and exploration. In interviews, six participants (P1–3, P6, P9–10) reported that intent dimensions presented through UI components—such as sliders, buttons, and tags—made it easy to adjust the nuance of their intents. Since these dimensions were not pre-defined but generated in response to the expressed intents, they felt they were more relevant. Related to this, P10 said *"By directly manipulating these values and instantly observing their effects on the output, I could do quick, low-effort exploration to fine-tune what I wanted to express and converge on intent configurations that better matched my goals."*

In addition, we found that the keep (🔖) feature in the **Intent Panel** played a key role in enabling efficient and targeted intent adjustment. Six participants (P3, P7-8, P10–12) specifically noted that they used the keep (🔖) feature to lock satisfying intents while employing intent-based prompting or direct editing to revise specific aspects they wanted to change. These features worked synergistically, supporting a more controlled and deliberate refinement process that helped users move closer to their intended direction without having to start from scratch each time. P7 remarked, *"With the bookmark (🔖) button, I was able to preserve the parts I liked and change only what I didn't. So the revision process felt more step-by-step."*

Finally, we found that the resulting set of well-articulated intents could serve as reusable assets. Survey responses showed that intent reusability was rated significantly higher in INTENTFLOW (M11: INTENTFLOW: $M = 5.92$, $SD = 1.38$; Baseline: $M = 4.17$, $SD = 1.53$;

$p = .001$, $t = 3.92$). In interviews, nine participants (P1, P3–6, P9–12) expressed interest in applying their final set of intents to future tasks of a similar nature. This suggests that IntentFlow helped transferable and reusable intent representations beyond the immediate session. P10 captured this perspective, stating, *"Once I refined my intents to get the level of output I wanted, I felt like I could use them again next time to get a similar level of quality right away. It would definitely be helpful."*

## 8 DISCUSSION

Our findings demonstrate the effectiveness of structured and interactive intent communication for enhancing LLM-assisted writing workflows. In this section, we reflect on key implications of our study and discuss directions for future research.

### 8.1 Emerging Challenges in Holistic Intent Communication

In our study, by examining intent communication in depth across • Articulation • Exploration • Management and • Synchronization we observed challenges that extend beyond supporting individual intent communication. As users accumulated a richer set of intents, new forms of **tension and conflict** emerged between them, which imply the need for managing not only individual intents, but also their relationships. For example, a user might initially specify avoiding jargon but later introduce an intent to explain a technical concept requiring domain-specific terms. While these tensions are manageable when only a few intents are present, they become increasingly difficult for users to handle directly as the set grows larger.

Related to this challenge, we also found that users often **struggled to recall their prior intents** once the set became large. This points to a need for management mechanisms that help users remain aware of what has been specified, modified, or deprioritized over time. Similarly, synchronization cannot be limited to verifying whether a single output reflects the most recent intent. Rather, users need a way to assess whether the system is consistently following their contextual constraints over the course of interaction.

Taken together, these observations suggest that intent communication is about **maintaining a shared mental space** between the user and the system. Beyond supporting articulation or one-off synchronization, systems should help users navigate tensions among multiple intents, sustain awareness of their evolving set of preferences, and ensure continuity of alignment across outputs.

### 8.2 Curating Intent as a Reusable and Transferable Component

Through IntentFlow, we observed that users could curate a set of intents by articulating and exploring what they wanted, while managing and synchronizing them with system outputs to iteratively refine them into stable configurations. In IntentFlow, intents are surfaced as modular and interactable components rather than being embedded within transient prompt texts; users can manage them as an evolving checklist—supporting both ongoing task execution and future reuse. Participants highlighted the value of reusing such finalized sets in similar future tasks, reapplying a consistent tone or framing strategy across writing projects to ensure stylistic coherence.

At first glance, this resembles recent LLM memory[48] and personalization features[7] in commercial systems. Yet, unlike those mechanisms, curated intents in IntentFlow are deliberately articulated and selectively preserved by users. Moreover, rather than reusing everything indiscriminately, participants envisioned selectively reusing certain intents, much like importing from a programming library, that could be applied when needed. This perspective also connects to recent work, AI Instruments [54], which reifies prompts into manipulable artifacts. Our findings extend this idea by showing how curated sets of evolving intents can serve as mid-level building blocks: more structured than raw prompts yet more flexible than fixed templates, supporting both reuse and adaptation over time. We envision systems like IntentFlow can serve not only as single-session assistants, but also as platforms for developing long-term, reusable intent strategies.

### 8.3 Generalizability of Design Implications Across Domains

Although IntentFlow was instantiated in the writing domain, the design implications behind it are not domain-specific. They are grounded in our SLR findings, which cover diverse generative tasks where users communicate fluid and evolving intents and receive system-generated outcomes. As summarized in Table 4, our design implications extend beyond writing and can be instantiated in other domains such as data analysis, image editing, and music composition. The table illustrates how intents in each domain can be articulated, explored, managed, and synchronized through domain-appropriate features. This demonstrates that the four aspects of intent communication and the corresponding design implications are generalizable across generative tasks beyond writing.

### 8.4 Design Opportunities for Enhancing Intent Communication

Revisiting the four aspects of intent communication, we suggest several design opportunities to improve intent communication for future research.

**Intent Articulation.** While structured intent representation in **Intent Panel** helped users articulate vague intents, their input was largely restricted to text and UI manipulation. Future work could investigate various input modalities (e.g., voice, sketch, examples) and interactive scaffolds that surface latent or subconscious intents.

**Intent Exploration.** Although IntentFlow generated appropriate UI components (e.g., radio buttons, sliders, hashtags) for intent dimensions, the set of interaction primitives remained limited. Future systems could leverage malleable user interfaces that can dynamically adapt to evolving user intents, allowing users to modify or extend UI elements to better fit their needs [6, 46]. Such malleable components—dynamic views, visual scaffolds, or domain-specific control panels—can provide more flexible and expressive ways to explore and refine intent.

**Intent Management.** As we discussed earlier, conflict resolution and prioritization between intents remained unsupported in IntentFlow. Managing such relationships becomes increasingly

| Design Implications | Writing | Data Analysis | Image Editing | Music Composition |
|---|---|---|---|---|
| **DI1: Distinguish & Externalize Goals and Intents** | - Prompts parsed into Goals (*"Write article"*) vs. Intents (*"polite tone, concise, add background"*).<br>- Goal and intents are displayed in an editable Intent Panel, where users can refine or adjust them. | - Prompts parsed into Goals (*"Analyze sales trends"*) vs. Intents (*"remove outliers, focus on seasonal pattern"*).<br>- Goal and intents are surfaced in an Intent Panel linked to code cells, where users can edit or supplement them adding data. | - Prompts separated into Goals (*"Design poster"*) vs. Intents (*"retro palette, strong typography emphasis"*).<br>- Goal and intents are displayed in a Intent Panel beside layers, allowing users to toggle or adjust them. | - Prompts split into Goals (*"Compose jazz piece"*) vs. Intents (*"slow tempo, saxophone lead, uplifting mood"*).<br>- Goal and intents are displayed in a Intent Panel, where users can view, drag, and refine them. |
| **DI2: Provide Exploratory Spaces** | - Providing **writing intents** with appropriate UI controls (e.g., sliders for tone or length, tags for focus, buttons for emphasis).<br>- Users can directly manipulate these dimensions and preview potential impacts | - Surfacing **analytic intents** (e.g., granularity, filtering, comparison type) through interactive widgets such as sliders, toggles, or radio buttons.<br>- Users can quickly preview how different parameter settings affect analysis results and visualizations. | - Representing **design intents** (e.g., color scheme, layout emphasis, stylistic choices) with dimension-specific UI controls.<br>- Users can switch between suggested options and preview outcomes with thumbnails or overlays. | - Exposing **musical intents** (e.g., tempo, instrumentation, harmonic progression) as adjustable dimensions.<br>- Users can experiment with variations and instantly preview their effects through audio playback. |
| **DI3: Support Versioning & Curation** | - Drafts are maintained with their associated intents as versions.<br>- Users can pin preferred intents and roll back to earlier versions, and compare changes through a diff view. | - Notebooks maintain an intent-version history.<br>- Locked intents persist across reruns, and users can compare alternative analytic results side-by-side. | - Each edit together with its intent set is recorded in the version history.<br>- Users can store certain intents as reusable presets, and compare alternatives using a before/after diff view. | - Musical drafts are kept as versions tied to intents.<br>- Users can bookmark configurations, revisit earlier takes, and compare old vs. new using waveform-based diff views. |
| **DI4: Make Intent–Output Connections Transparent** | Intent hover highlights affected text. "Tone: academic" highlights formalized sentences. | Intent hover highlights linked chart region. "Seasonality focus" emphasizes time-series segments. "Remove outliers" recolors excluded points. | Intent hover outlines impacted layers. "Bright colors" intent highlights applied layers. Live preview overlay shows dimension changes. | Intent hover highlights score/timeline region. "Saxophone solo" intent highlights waveform of solo section. Tempo change preview simulates waveform stretching. |

**Table 4: Instantiations of the four design implications for intent communication (DI1–DI4) across domains.**

difficult as the number of intents grows, and expecting users to explicitly track or prioritize them can impose a high cognitive burden. To this end, we envision mixed-initiative approaches [23, 44] that detect conflicts, suggest prioritization strategies, and help users resolve competing intents. Moreover, the level of proactivity can adapt to the user's state, offering more suggestions when they are still exploring directions and stepping back as they converge on a final version. In addition, treating intents as reusable objects that persist across tasks and allowing users to selectively import them can extend their utility beyond single sessions

**Intent Synchronization.** Hover-based linking in INTENTFLOW helped users check how intents were reflected in outputs, but this support was limited to surface-level associations without the system's actual interpretations. Future systems could surface intermediate representations of how user intents are understood, highlight mismatches between intended and realized outputs, and provide predictive previews of how edits will propagate through results. Such mechanisms would enable users to continuously verify and calibrate whether their intents are being understood and realized as intended.

### 8.5 Limitations and Future Work

Although our findings highlight the benefits of interactive intent communication, our current study and system design have certain limitations that open up directions for further research.

First, we asked participants to engage in controlled writing tasks rather than their own authentic work. While this setup enabled systematic comparisons, the low-stakes nature may not fully capture the depth or complexity of real-world intent communication. Furthermore, although we examined participants' perceptions of intent reusability, our study design did not capture how curated intents are actually reused across tasks or how they evolve over time. Addressing these gaps requires moving beyond lab study. Future work could deploy systems like INTENTFLOW in naturalistic settings and conduct longitudinal studies, observing how intent

communication unfolds in practice and how they support the reuse and long-term evolution of user practices.

Second, to focus on intent communication, we deliberately excluded manual editing of outputs in our study, because pilot sessions showed that allowing it diverted participants toward surface-level fixes, which obscured the very communication processes we aimed to study (RQ1). While this choice facilitated a more controlled study, it did not fully reflect realistic workflows. We also envision that manual editing can be reinterpreted as intent expression. Future work could reinterpret manual editing as another channel of intent expression and investigate how they together shape the overall intent communication process.

## 9 CONCLUSION

In this paper, we conducted a systematic literature review of 33 HCI papers on generative AI systems and synthesized four key aspects of intent communication: **articulation**, **exploration**, **management**, and **synchronization**. Building on these findings, we derived design implications and instantiated them in INTENTFLOW, a system for LLM-based writing that supports all four aspects through adjustable UIs, intent-to-output linking, and versioned refinement.

Our evaluation demonstrates the effectiveness of INTENTFLOW in supporting intent communication with LLMs. The technical evaluation confirmed robust performance across diverse writing contexts, and the user study showed that INTENTFLOW helps users articulate vague intents, refine evolving preferences, synchronize them with outputs, and consolidate a curated set of intents. Interaction logs further revealed a shift from corrective prompting toward proactive refinement.

Together, these findings highlight the benefits of making user intents explicitly representable and manipulable in LLM-based systems. Structured support for the four aspects of intent communication can reduce the overhead of expressing fluid and evolving intents, improve user control, and yield outputs better aligned with users' intentions. Future work could extend this approach beyond

writing and explore more adaptive, proactive forms of intent support.

## REFERENCES

[1] Usman Anwar, Abulhair Saparov, Javier Rando, Daniel Paleka, Miles Turpin, Peter Hase, Ekdeep Singh Lubana, Erik Jenner, Stephen Casper, Oliver Sourbut, Benjamin L. Edelman, Zhaowei Zhang, Mario Günther, Anton Korinek, Jose Hernandez-Orallo, Lewis Hammond, Eric Bigelow, Alexander Pan, Lauro Langosco, Tomasz Korbak, Heidi Zhang, Ruiqi Zhong, Seán Ó hÉigeartaigh, Gabriel Recchia, Giulio Corsi, Alan Chan, Markus Anderljung, Lilian Edwards, Aleksandar Petrov, Christian Schroeder de Witt, Sumeet Ramesh Motwan, Yoshua Bengio, Danqi Chen, Philip H. S. Torr, Samuel Albanie, Tegan Maharaj, Jakob Foerster, Florian Tramer, He He, Atoosa Kasirzadeh, Yejin Choi, and David Krueger. 2024. Foundational Challenges in Assuring Alignment and Safety of Large Language Models. arXiv:2404.09932 [cs.LG] https://arxiv.org/abs/2404.09932

[2] Ian Arawjo, Chelse Swoopes, Priyan Vaithilingam, Martin Wattenberg, and Elena L Glassman. 2024. Chainforge: A visual toolkit for prompt engineering and llm hypothesis testing. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*. 1–18.

[3] Yuwei Bao, Keunwoo Peter Yu, Yichi Zhang, Shane Storks, Itamar Bar-Yossef, Alexander De La Iglesia, Megan Su, Xiao Lin Zheng, and Joyce Chai. 2023. Can Foundation Models Watch, Talk and Guide You Step by Step to Make a Cake? *arXiv preprint arXiv:2311.00738* (2023).

[4] Stephen Brade, Bryan Wang, Mauricio Sousa, Sageev Oore, and Tovi Grossman. 2023. Promptify: Text-to-image generation through interactive prompt exploration with large language models. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. 1–14.

[5] Bill Buxton. 2010. *Sketching user experiences: getting the design right and the right design.* Morgan kaufmann.

[6] Yining Cao, Peiling Jiang, and Haijun Xia. 2025. Generative and Malleable User Interfaces with Generative and Evolving Task-Driven Data Model. arXiv:2503.04084 [cs.HC] https://arxiv.org/abs/2503.04084

[7] Anthropic Help Center. [n. d.]. Understanding Claude's Personalization Features. https://support.anthropic.com/en/articles/10185728-understanding-claude-s-personalization-features [Online; accessed 2025-09-11].

[8] Liuqing Chen, Qianzhi Jing, Yixin Tsang, Qianyi Wang, Ruocong Liu, Duowei Xia, Yunzhan Zhou, and Lingyun Sun. 2024. AutoSpark: Supporting Automobile Appearance Design Ideation with Kansei Engineering and Generative AI. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology* (Pittsburgh, PA, USA) *(UIST '24)*. Association for Computing Machinery, New York, NY, USA, Article 108, 19 pages. https://doi.org/10.1145/3654777.3676337

[9] Wei-Hao Chen, Weixi Tong, Ph.D. Case, Amanda, and Tianyi Zhang. 2025. Dango: A Mixed-Initiative Data Wrangling System using Large Language Model. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 389, 28 pages. https://doi.org/10.1145/3706598.3714135

[10] Jianpeng Cheng, Devang Agrawal, Héctor Martínez Alonso, Shruti Bhargava, Joris Driesen, Federico Flego, Shaona Ghosh, Dain Kaplan, Dimitri Kartsaklis, Lin Li, et al. 2020. Conversational semantic parsing for dialog state tracking. *arXiv preprint arXiv:2010.12770* (2020).

[11] DaEun Choi, Kihoon Son, HyunJoon Jung, and Juho Kim. 2025. Expandora: Broadening Design Exploration with Text-to-Image Model. In *Proceedings of the Extended Abstracts of the CHI Conference on Human Factors in Computing Systems (CHI EA '25)*. Association for Computing Machinery, New York, NY, USA, Article 232, 10 pages. https://doi.org/10.1145/3706599.3720189

[12] John Joon Young Chung and Max Kreminski. 2024. Patchview: LLM-powered Worldbuilding with Generative Dust and Magnet Visualization. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology* (Pittsburgh, PA, USA) *(UIST '24)*. Association for Computing Machinery, New York, NY, USA, Article 77, 19 pages. https://doi.org/10.1145/3654777.3676352

[13] Richard E Clark, David F Feldon, Jeroen JG Van Merrienboer, Kenneth A Yates, and Sean Early. 2008. Cognitive task analysis. In *Handbook of research on educational communications and technology*. Routledge, 577–593.

[14] Yang Deng, Wenqiang Lei, Wai Lam, and Tat-Seng Chua. 2023. A Survey on Proactive Dialogue Systems: Problems, Methods, and Prospects. arXiv:2305.02750 [cs.CL] https://arxiv.org/abs/2305.02750

[15] Ian Drosos, Jack Williams, Advait Sarkar, Nicholas Wilson, Sean Rintel, and Payod Panda. 2025. Dynamic Prompt Middleware: Contextual Prompt Refinement Controls for Comprehension Tasks. In *Proceedings of the 4th Annual Symposium on Human-Computer Interaction for Work (CHIWORK '25)*. Association for Computing Machinery, New York, NY, USA, Article 24, 23 pages. https://doi.org/10.1145/3729176.3729203

[16] Linda Flower and John R Hayes. 1981. A cognitive process theory of writing. *College Composition & Communication* 32, 4 (1981), 365–387.

[17] Katy Ilonka Gero, Chelse Swoopes, Ziwei Gu, Jonathan K. Kummerfeld, and Elena L. Glassman. 2024. Supporting Sensemaking of Large Language Model Outputs at Scale. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 838, 21 pages. https://doi.org/10.1145/3613904.3642139

[18] Frederic Gmeiner, Jamie Lynn Conlin, Eric Handa Tang, Nikolas Martelaro, and Kenneth Holstein. 2024. An Evidence-based Workflow for Studying and Designing Learning Supports for Human-AI Co-creation. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI EA '24)*. Association for Computing Machinery, New York, NY, USA, Article 42, 15 pages. https://doi.org/10.1145/3613905.3650763

[19] Frederic Gmeiner, Nicolai Marquardt, Michael Bentley, Hugo Romat, Michel Pahud, David Brown, Asta Roseway, Nikolas Martelaro, Kenneth Holstein, Ken Hinckley, and Nathalie Riche. 2025. Intent Tagging: Exploring Micro-Prompting Interactions for Supporting Granular Human-GenAI Co-Creation Workflows. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 531, 31 pages. https://doi.org/10.1145/3706598.3713861

[20] Andreas Göldi, Thiemo Wambsganss, Seyed Parsa Neshaei, and Roman Rietsche. 2024. Intelligent Support Engages Writers Through Relevant Cognitive Processes. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 1047, 12 pages. https://doi.org/10.1145/3613904.3642549

[21] Arella E Gussow. 2023. Language production under message uncertainty: When, how, and why we speak before we think. In *Psychology of Learning and Motivation*. Vol. 78. Elsevier, 83–117.

[22] Sandra G. Hart and Lowell E. Staveland. 1988. Development of NASA-TLX (Task Load Index): Results of Empirical and Theoretical Research. In *Human Mental Workload*, Peter A. Hancock and Najmedin Meshkati (Eds.). Advances in Psychology, Vol. 52. North-Holland, 139–183. https://doi.org/10.1016/S0166-4115(08)62386-9

[23] Eric Horvitz. 1999. Principles of mixed-initiative user interfaces. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Pittsburgh, Pennsylvania, USA) *(CHI '99)*. Association for Computing Machinery, New York, NY, USA, 159–166. https://doi.org/10.1145/302979.303030

[24] Edwin L Hutchins, James D Hollan, and Donald A Norman. 1985. Direct manipulation interfaces. *Human–computer interaction* 1, 4 (1985), 311–338.

[25] Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert McHardy. 2023. Challenges and Applications of Large Language Models. arXiv:2307.10169 [cs.CL] https://arxiv.org/abs/2307.10169

[26] Majeed Kazemitabaar, Jack Williams, Ian Drosos, Tovi Grossman, Austin Zachary Henley, Carina Negreanu, and Advait Sarkar. 2024. Improving Steering and Verification in AI-Assisted Data Analysis with Interactive Task Decomposition. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology* (Pittsburgh, PA, USA) *(UIST '24)*. Association for Computing Machinery, New York, NY, USA, Article 92, 19 pages. https://doi.org/10.1145/3654777.3676345

[27] Anjali Khurana, Hariharan Subramonyam, and Parmit K Chilana. 2024. Why and When LLM-Based Assistants Can Go Wrong: Investigating the Effectiveness of Prompt-Based Interactions for Software Help-Seeking. In *Proceedings of the 29th International Conference on Intelligent User Interfaces* (Greenville, SC, USA) *(IUI '24)*. Association for Computing Machinery, New York, NY, USA, 288–303. https://doi.org/10.1145/3640543.3645200

[28] Hui-Jun Kim, Jeongho Kim, Sohyun Jeong, Minbong Lee, Jaegul Choo, and Sung-Hee Kim. 2025. ShoeGenAI: A Creativity Support Tool for High-Feasible Shoe Product Design. In *Proceedings of the Extended Abstracts of the CHI Conference on Human Factors in Computing Systems (CHI EA '25)*. Association for Computing Machinery, New York, NY, USA, Article 478, 11 pages. https://doi.org/10.1145/3706599.3721204

[29] Taewan Kim, Donghoon Shin, Young-Ho Kim, and Hwajung Hong. 2024. DiaryMate: Understanding User Perceptions and Experience in Human-AI Collaboration for Personal Journaling. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 1046, 15 pages. https://doi.org/10.1145/3613904.3642693

[30] Tae Soo Kim, Yoonjoo Lee, Minsuk Chang, and Juho Kim. 2023. Cells, Generators, and Lenses: Design Framework for Object-Oriented Interaction with Large Language Models. 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 4, 18 pages. https://doi.org/10.1145/3586183.3606833

[31] Tae Soo Kim, Yoonjoo Lee, Jamin Shin, Young-Ho Kim, and Juho Kim. 2024. EvalLM: Interactive Evaluation of Large Language Model Prompts on User-Defined Criteria. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 306, 21 pages. https://doi.org/10.1145/3613904.3642216

[32] Yoonsu Kim, Brandon Chin, Kihoon Son, Seoyoung Kim, and Juho Kim. 2025. Applying the Gricean Maxims to a Human-LLM Interaction Cycle: Design Insights from a Participatory Approach. In *Proceedings of the Extended Abstracts of the CHI Conference on Human Factors in Computing Systems (CHI EA '25)*. Association for Computing Machinery, New York, NY, USA, Article 72, 8 pages. https://doi.org/10.1145/3706599.3719759

[33] Yoonsu Kim, Jueon Lee, Seoyoung Kim, Jaehyuk Park, and Juho Kim. 2023. Understanding Users' Dissatisfaction with ChatGPT Responses: Types, Resolving Tactics, and the Effect of Knowledge Level. *arXiv preprint arXiv:2311.07434* (2023).

[34] Mark L Knapp and John A Daly. 2011. *The SAGE handbook of interpersonal communication*. Sage Publications.

[35] Tanya Kraljic and Michal Lahav. 2024. From Prompt Engineering to Collaborating: A Human-Centered Approach to AI Interfaces. *Interactions* 31, 3 (May 2024), 30–35. https://doi.org/10.1145/3652622

[36] Philippe Laban, Jesse Vig, Marti Hearst, Caiming Xiong, and Chien-Sheng Wu. 2024. Beyond the Chat: Executable and Verifiable Text-Editing with LLMs. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology* (Pittsburgh, PA, USA) *(UIST '24)*. Association for Computing Machinery, New York, NY, USA, Article 20, 23 pages. https://doi.org/10.1145/3654777.3676419

[37] Cassandra Lee and Jessica R Mindel. 2024. Closer and Closer Worlds: Using LLMs to Surface Personal Stories in World-building Conversation Games. In *Companion Publication of the 2024 ACM Designing Interactive Systems Conference* (IT University of Copenhagen, Denmark) *(DIS '24 Companion)*. Association for Computing Machinery, New York, NY, USA, 289–293. https://doi.org/10.1145/3656156.3665430

[38] Christine P. Lee, David Porfirio, Xinyu Jessica Wang, Kevin Chenkai Zhao, and Bilge Mutlu. 2025. VeriPlan: Integrating Formal Verification and LLMs into End-User Planning. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 247, 19 pages. https://doi.org/10.1145/3706598.3714113

[39] Mina Lee, Katy Ilonka Gero, John Joon Young Chung, Simon Buckingham Shum, Vipul Raheja, Hua Shen, Subhashini Venugopalan, Thiemo Wambsganss, David Zhou, Emad A. Alghamdi, Tal August, Avinash Bhat, Madiha Zahrah Choksi, Senjuti Dutta, Jin L.C. Guo, Md Naimul Hoque, Yewon Kim, Simon Knight, Seyed Parsa Neshaei, Antonette Shibani, Disha Shrivastava, Lila Shroff, Agnia Sergeyuk, Jessi Stark, Sarah Sterman, Sitong Wang, Antoine Bosselut, Daniel Buschek, Joseph Chee Chang, Sherol Chen, Max Kreminski, Joonsuk Park, Roy Pea, Eugenia Ha Rim Rho, Zejiang Shen, and Pao Siangliulue. 2024. A Design Space for Intelligent and Interactive Writing Assistants. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 1054, 35 pages. https://doi.org/10.1145/3613904.3642697

[40] Jenny T. Liang, Chenyang Yang, and Brad A. Myers. 2024. A Large-Scale Survey on the Usability of AI Programming Assistants: Successes and Challenges. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering* (Lisbon, Portugal) *(ICSE '24)*. Association for Computing Machinery, New York, NY, USA, Article 52, 13 pages. https://doi.org/10.1145/3597503.3608128

[41] Hyunseung Lim, Ji Yong Cho, Taewan Kim, Jeongeon Park, Hyungyu Shin, Seulgi Choi, Sunghyun Park, Kyungjae Lee, Juho Kim, Moontae Lee, and Hwajung Hong. 2024. Co-Creating Question-and-Answer Style Articles with Large Language Models for Research Promotion. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference* (Copenhagen, Denmark) *(DIS '24)*. Association for Computing Machinery, New York, NY, USA, 975–994. https://doi.org/10.1145/3643834.3660705

[42] Michael Xieyang Liu, Frederick Liu, Alexander J. Fiannaca, Terry Koo, Lucas Dixon, Michael Terry, and Carrie J. Cai. 2024. "We Need Structured Output": Towards User-centered Constraints on Large Language Model Output. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI EA '24)*. Association for Computing Machinery, New York, NY, USA, Article 10, 9 pages. https://doi.org/10.1145/3613905.3650756

[43] 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

[44] Xingyu Bruce Liu, Haijun Xia, and Xiang Anthony Chen. 2025. Interacting with Thoughtful AI. arXiv:2502.18676 [cs.HC] https://arxiv.org/abs/2502.18676

[45] Damien Masson, Sylvain Malacria, Géry Casiez, and Daniel Vogel. 2024. DirectGPT: A Direct Manipulation Interface to Interact with Large Language Models. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 975, 16 pages. https://doi.org/10.1145/3613904.3642462

[46] Bryan Min, Allen Chen, Yining Cao, and Haijun Xia. 2025. Malleable Overview-Detail Interfaces. *arXiv preprint arXiv:2503.07782* (2025).

[47] Jakob Nielsen. 2023. AI: First New UI Paradigm in 60 Years. https://www.nngroup.com/articles/ai-paradigm/. (Accessed on 02/18/2025).

[48] OpenAI. 2024. Memory and new controls for ChatGPT. https://openai.com/index/memory-and-new-controls-for-chatgpt/ [Online; accessed 2025-09-11].

[49] Matthew J Page, Joanne E McKenzie, Patrick M Bossuyt, Isabelle Boutron, Tammy C Hoffmann, Cynthia D Mulrow, Larissa Shamseer, Jennifer M Tetzlaff, Elie A Akl, Sue E Brennan, et al. 2021. The PRISMA 2020 statement: an updated guideline for reporting systematic reviews. *bmj* 372 (2021).

[50] Xiaohan Peng, Janin Koch, and Wendy E. Mackay. 2024. DesignPrompt: Using Multimodal Interaction for Design Exploration with Generative AI. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference* (Copenhagen, Denmark) *(DIS '24)*. Association for Computing Machinery, New York, NY, USA, 804–818. https://doi.org/10.1145/3643834.3661588

[51] Yingzhe Peng, Xiaoting Qin, Zhiyang Zhang, Jue Zhang, Qingwei Lin, Xu Yang, Dongmei Zhang, Saravan Rajmohan, and Qi Zhang. 2025. Navigating the Unknown: A Chat-Based Collaborative Interface for Personalized Exploratory Tasks. In *Proceedings of the 30th International Conference on Intelligent User Interfaces (IUI '25)*. Association for Computing Machinery, New York, NY, USA, 1048–1063. https://doi.org/10.1145/3708359.3712093

[52] Cheng Qian, Bingxiang He, Zhong Zhuang, Jia Deng, Yujia Qin, Xin Cong, Yankai Lin, Zhong Zhang, Zhiyuan Liu, and Maosong Sun. 2024. Tell me more! towards implicit user intention understanding of language model driven agents. *arXiv preprint arXiv:2402.09205* (2024).

[53] Mohi Reza, Nathan M Laundry, Ilya Musabirov, Peter Dushniku, Zhi Yuan "Michael" Yu, Kashish Mittal, Tovi Grossman, Michael Liut, Anastasia Kuzminykh, and Joseph Jay Williams. 2024. ABScribe: Rapid Exploration & Organization of Multiple Writing Variations in Human-AI Co-Writing Tasks using Large Language Models. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 1042, 18 pages. https://doi.org/10.1145/3613904.3641899

[54] Nathalie Riche, Anna Offenwanger, Frederic Gmeiner, David Brown, Hugo Romat, Michel Pahud, Nicolai Marquardt, Kori Inkpen, and Ken Hinckley. 2025. AI-Instruments: Embodying Prompts as Instruments to Abstract & Reflect Graphical Interface Commands as General-Purpose Tools. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 1104, 18 pages. https://doi.org/10.1145/3706598.3714259

[55] Juan A. Rodriguez, Nicholas Botzer, David Vazquez, Christopher Pal, Marco Pedersoli, and Issam Laradji. 2024. IntentGPT: Few-shot Intent Discovery with Large Language Models. arXiv:2411.10670 [cs.CL] https://arxiv.org/abs/2411.10670

[56] D.A. Schön. 1992. Designing as reflective conversation with the materials of a design situation. *Knowledge-Based Systems* 5, 1 (1992), 3–14. https://doi.org/10.1016/0950-7051(92)90020-G Artificial Intelligence in Design Conference 1991 Special Issue.

[57] Yashothara Shanmugarasa, Shidong Pan, Ming Ding, Dehai Zhao, and Thierry Rakotoarivelo. 2025. Privacy Meets Explainability: Managing Confidential Data and Transparency Policies in LLM-Empowered Science. In *Proceedings of the Extended Abstracts of the CHI Conference on Human Factors in Computing Systems (CHI EA '25)*. Association for Computing Machinery, New York, NY, USA, Article 448, 8 pages. https://doi.org/10.1145/3706599.3720099

[58] Donghoon Shin, Gary Hsieh, and Young-Ho Kim. 2023. PlanFitting: Tailoring Personalized Exercise Plans with Large Language Models. *arXiv preprint arXiv:2309.12555* (2023).

[59] Ben Shneiderman. 1983. Direct manipulation: A step beyond programming languages. *Computer* 16, 08 (1983), 57–69.

[60] Ben Shneiderman. 2007. Creativity support tools: accelerating discovery and innovation. *Commun. ACM* 50, 12 (2007), 20–32.

[61] Momin N Siddiqui, Roy D Pea, and Hari Subramonyam. 2025. Script&Shift: A Layered Interface Paradigm for Integrating Content Development and Rhetorical Strategy with LLM Writing Assistants. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 532, 19 pages. https://doi.org/10.1145/3706598.3714119

[62] Glen H Stamp and Mark L Knapp. 1990. The construct of intent in interpersonal communication. *Quarterly journal of speech* 76, 3 (1990), 282–299.

[63] Hari Subramonyam, Roy Pea, Christopher Pondoc, Maneesh Agrawala, and Colleen Seifert. 2024. Bridging the Gulf of Envisioning: Cognitive Challenges in Prompt Based Interactions with LLMs. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 1039, 19 pages. https://doi.org/10.1145/3613904.3642754

[64] Sangho Suh, Meng Chen, Bryan Min, Toby Jia-Jun Li, and Haijun Xia. 2024. Luminate: Structured Generation and Exploration of Design Space with Large Language Models for Human-AI Co-Creation. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 644, 26 pages. https://doi.org/10.1145/3613904.3642400

[65] Lev Tankelevitch, Viktor Kewenig, Auste Simkute, Ava Elizabeth Scott, Advait Sarkar, Abigail Sellen, and Sean Rintel. 2024. The Metacognitive Demands and Opportunities of Generative AI. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 680, 24 pages. https://doi.org/10.1145/3613904.3642902

[66] Michael Terry, Chinmay Kulkarni, Martin Wattenberg, Lucas Dixon, and Meredith Ringel Morris. 2023. AI Alignment in the Design of Interactive AI: Specification Alignment, Process Alignment, and Evaluation Support. *arXiv preprint arXiv:2311.00710* (2023).

[67] Michael Terry and Elizabeth D Mynatt. 2002. Recognizing creative needs in user interface design. In *Proceedings of the 4th Conference on Creativity & Cognition*. 38–44.

[68] Bekzat Tilekbay, Saelyne Yang, Michal Adam Lewkowicz, Alex Suryapranata, and Juho Kim. 2024. ExpressEdit: Video Editing with Natural Language and Sketching. In *Proceedings of the 29th International Conference on Intelligent User Interfaces* (Greenville, SC, USA) *(IUI '24)*. Association for Computing Machinery, New York, NY, USA, 515–536. https://doi.org/10.1145/3640543.3645164

[69] Huanchen Wang, Tianrun Qiu, Jiaping Li, Zhicong Lu, and Yuxin Ma. 2025. HarmonyCut: Supporting Creative Chinese Paper-cutting Design with Form and Connotation Harmony. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. Association for Computing Machinery, New York, NY, USA, Article 661, 22 pages. https://doi.org/10.1145/3706598.3714159

[70] Zhijie Wang, Yuheng Huang, Da Song, Lei Ma, and Tianyi Zhang. 2024. PromptCharm: Text-to-Image Generation through Multi-modal Prompting and Refinement. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '24)*. Association for Computing Machinery, New York, NY, USA, Article 185, 21 pages. https://doi.org/10.1145/3613904.3642803

[71] Zehuan Wang, Jiaqi Xiao, Jingwei Sun, and Can Liu. 2025. IntentPrism: Human-AI Intent Manifestation for Web Information Foraging. In *Proceedings of the Extended Abstracts of the CHI Conference on Human Factors in Computing Systems (CHI EA '25)*. Association for Computing Machinery, New York, NY, USA, Article 345, 11 pages. https://doi.org/10.1145/3706599.3719744

[72] Shirley Wu, Michel Galley, Baolin Peng, Hao Cheng, Gavin Li, Yao Dou, Weixin Cai, James Zou, Jure Leskovec, and Jianfeng Gao. 2025. CollabLLM: From Passive Responders to Active Collaborators. arXiv:2502.00640 [cs.AI] https://arxiv.org/abs/2502.00640

[73] 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

[74] Catherine Yeh, Gonzalo Ramos, Rachel Ng, Andy Huntington, and Richard Banks. 2024. Ghostwriter: Augmenting collaborative human-ai writing experiences through personalization and agency. *arXiv preprint arXiv:2402.08855* (2024).

[75] J.D. Zamfirescu-Pereira, Heather Wei, Amy Xiao, Kitty Gu, Grace Jung, Matthew G Lee, Bjoern Hartmann, and Qian Yang. 2023. Herding AI Cats: Lessons from Designing a Chatbot by Prompting GPT-3. In *Proceedings of the 2023 ACM Designing Interactive Systems Conference* (Pittsburgh, PA, USA) *(DIS '23)*. Association for Computing Machinery, New York, NY, USA, 2206–2220. https://doi.org/10.1145/3563657.3596138

[76] J.D. Zamfirescu-Pereira, Richmond Y. Wong, Bjoern Hartmann, and Qian Yang. 2023. Why Johnny Can't Prompt: How Non-AI Experts Try (and Fail) to Design LLM Prompts. 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 437, 21 pages. https://doi.org/10.1145/3544548.3581388

[77] Hongbo Zhang, Pei Chen, Xuelong Xie, Chaoyi Lin, Lianyan Liu, Zhuoshu Li, Weitao You, and Lingyun Sun. 2024. ProtoDreamer: A Mixed-prototype Tool Combining Physical Model and Generative AI to Support Conceptual Design. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology* (Pittsburgh, PA, USA) *(UIST '24)*. Association for Computing Machinery, New York, NY, USA, Article 97, 18 pages. https://doi.org/10.1145/3654777.3676399

[78] 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*. 1–30.

[79] Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2024. Explainability for Large Language Models: A Survey. *ACM Trans. Intell. Syst. Technol.* 15, 2, Article 20 (Feb. 2024), 38 pages. https://doi.org/10.1145/3639372

[80] Chen Zhou, Zihan Yan, Ashwin Ram, Yue Gu, Yan Xiang, Can Liu, Yun Huang, Wei Tsang Ooi, and Shengdong Zhao. 2024. GlassMail: Towards Personalised Wearable Assistant for On-the-Go Email Creation on Smart Glasses. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference* (Copenhagen, Denmark) *(DIS '24)*. Association for Computing Machinery, New York, NY, USA, 372–390. https://doi.org/10.1145/3643834.3660683

# A APPENDICES

## A.1 Pipeline Module Prompts

---

## Entrypoint Chat Module

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

### System Prompt

```
You are a **highly intelligent AI assistant** designed to **analyze user queries and determine how to update or refine their task-related
information**. Your primary role is **not** to directly respond to user queries, but to **decide which module(s) should be updated** and
explain why.

## Your Role:
- Your main responsibility is to **analyze user queries** and determine how they impact the **Goal** and **Intent** modules.
- **Do not directly answer user queries** unless they are explicitly asking about **why the Goal and Intent modules were set in a certain
way** (e.g., "Why is my goal set this way?" or "Why are these intents selected?").
- By default, for the **user's first query**, always return **both the Goal and Intent modules** as updated.
- If a selected module is provided, the module must be set to be updated.

## Inputs You Will Receive:
1. **User Query:** The latest user input.
2. Selected Module:
   - This indicates a specific module (Goal or Intents or Intent Dimentions) the user is currently focusing on.
   - If selected Module is not null, you must always include this module as updated.
   - Additionally, analyze how the user query affects this selected module specifically.
3. **Chat History:** Previous interactions with the user.
4. **Current Module States:** The latest information from:
   - **Goal Module**: Contains the user's task objective, topic, and domain.
   - **Intent Module**: Contains the user's specific **requirements, preferences, and strategies** for achieving their task objective.
   - **User Intent Dimensions**: Represents the **dimensions of the user's intents as UI components**, storing these dimensions and their
   corresponding values.

## Your Tasks:
1. **Determine whether the query requires updating the Goal or Intent module.**
   - **By default, always return both the Goal and Intent modules as updated for the user's first query.**
   - If a selected module is provided, always include that module as updated.
    - Carefully analyze how the user's query is intended to refine or update the selected module.
    - Provide a clear explanation of how the user query affects the selected module's information.
   - The **Goal module** should remain largely unchanged unless the user presents an entirely new task.
   - If the query does not require updating the user's Goal or Intent, and is instead a meta-question (e.g., "Why is my goal set this
   way?"), then provide a direct response instead of updating any modules.

2.**If a module needs updating, return the recommended module(s) along with a clear explanation.**
   - If multiple modules require updates, list all relevant ones.
   - Ensure there are no duplicate modules in the updated modules list. Each module (goal, intents, intent dimensions) should appear at
   most once.
   - Ensure your reasoning is clear, well-structured, and directly tied to the user's task.

\#\# When to Directly Answer the User's Query:
- **Only** respond directly if the query is about the **reasoning behind the Goal or Intent modules' configuration**.
- Example queries that should be answered directly:
  - "Why was my Goal set to this topic?"
  - "How were my intents determined?"
- In all other cases, **focus on module updates rather than answering the query directly**.

Return your response in the following JSON format EXACTLY:
```json
{
    "response": "Direct response to the user's query (if applicable).",
    "updated_modules": [
        {
            "module": "goal || intents || intent_dimensions",
            "reason": "Why this module needs updating."
        }
    ]
}
```

**Goal Module**

---

**System Prompt**

You are a helpful and analytical assistant tasked with analyzing the user's query and extracting their task, domain, and topic. The user provides a writing task as a query through the chat interface in the system. Analyze the provided query to identify:
    - What the writing task (`task` in the output) is asking for.
    - Which domain (`domain` in the output) the writing task belongs to (e.g., Journalism Writing, Academic Writing, Creative Writing, Technical Writing, etc.).
    - What the topic (`topic` in the output) of the writing task is.

## Input You Will Receive:
1. **User Query:** The user input.
2. **Interaction History:** Previous user query and goal output history.

## Your task
First, you need to carefully review the user's query and reasoning the user's request deeply.
Second, you need to extract the task, domain, and topic from the user's query.
Lastly, you need to provide the extracted information in the JSON format like the example below.

Return your response in the following JSON format EXACTLY:
{
    "query": "user provided query",
    "task": {
        "value": "task/objective of the user query",
    },
    "domain": {
        "value": "domain of the user query",
    },
    "topic": {
        "value": "topic of the user query",
    }
}

## Intent Module

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

### System Prompt

You are a helpful and analytical assistant tasked with analyzing the user's `query` along with its `context` (such as the provided `task`, `domain`, and `topic`), and then extracting specific and actionable `intent`(s) from the user `query` and `context`.

The user has requested a writing task as a `query` through the chat interface in the current system. You are provided with four pieces of information: `query`, `task`, `domain`, and `topic`. Your task is extracting concrete and actionable intents based on these four information.
`query` is the user's input query. `task` is the objective of the user query. `domain` is the writing task's domain (e.g., Journalism Writing, Academic Writing, Creative Writing, Technical Writing, etc.). `topic` is the writing task's topic.

## Input You Will Receive:
1. **User Query:** The user input.
2. **Current Goal Context:** The current user context (task goal, domain, topic).
3. **Interaction History:** Previous user query, context, and intent list output history.

## Your task
First, you need to carefully review the given input information. In this process, you must deeply reason about what the user truly wants.
Second, you should extract the task, domain, and topic from the user's query.

## You should extract both:
- Explicit intents: Clearly and directly stated intentions in the user's query and context.
- Implicit intents: Essential, reasonable, or logically required steps, processes, or goals that are not directly mentioned but are necessary to accomplish the user's writing task successfully.

## The extracted intents must:
- Be specific, explicit, and actionable, so that the output can be generated immediately based on them.
- Include all relevant implicit intents inferred from the task, domain, and topic, even if not directly stated by the user.
- Not contain duplicates.
- Do not include the task itself as part of the user intents.

Return your response in the following JSON format EXACTLY:
```
{
    "intents": [
        {
            "intent": <specific, explicit, or implicit actionable intent>,
        }
        ... ,
        {
            "intent": <specific, explicit, or implicit actionable intent>,
        }
    ]
}
```

**Intent Dimension Module**

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

## System Prompt

You are an analytical and precise assistant tasked with defining appropriate intent dimensions for the given user's intent and selecting the most suitable UI layout to clearly represent key aspects of their task-related needs.

Your role is to analyze the user's intents (requirements, preferences, and strategies) and determine appropriate **intent dimensions** that capture key aspects of their task-related needs. You will then assign the most suitable UI layout for each dimension and set an **initial value** based on the user's query, task goal, and current intents.

There are three possible UI layouts for each intent dimension. You can use as many or as few of each layout as you want:

1. Likert Scale Layout:
    - When appropriate: Used for dimensions with discrete, ordered options
    - Output requirements: title and array of options from left to right
    - Example: Writing Stage (options: ["Idea Generation", "Planning", "Drafting", "Revision"])

2. Sliding Scale Layout:
    - When appropriate: Used for dimensions with continuous numeric values, up to 5 values (min=1, max=5)
    - Output requirements: title, left label, right label, min value, max value
    **RESTRICTION**:
        - Avoid overly granular or excessively narrow scales (e.g., avoid using min=0 and max=200 just because the user mentions 200 words.
        Instead, group it in practical ranges like 50-100, 100-300, etc.).
        - Example: Specificity (left: "General Overview", right: "Detailed Requirements", range: 1-5)

3. Hashtag Layout:
    - When appropriate: Used for dimensions with multiple selectable tags
    - Output requirements: title and array of possible tags
    - Example: Writing Context (tags: ["Academic", "Creative", "Technical", "Professional"])

## Input You Will Receive:
1. **User Query:** The user input.
2. **Current Goal Context:** The current user context (task goal, domain, topic).
3. **Current Intent List:** The current user intent list.
4. **Interaction History:** Previous user query, context, intent list, and intent dimension output history.

## Your Task
Based on the user's query, task goal, and context, and user intents, determine at least three dimensions that are relevant and which UI layout is most appropriate for each. Examples of dimensions are: Writing Stage, Writing Context, Purpose, Specificity, Audience, and Background Knowledge.

These are only examples!! Please come up with at least one new dimension that isn't mentioned in the examples, and try not to use these examples as a direct reference. You also don't have to use all three layouts; you can use as many or as few of each layout as you want. Return your response with dimensions in the following JSON format. If there are n dimensions, there should be n elements in the dimensions array:

Likert Scale:
```
    {
        "dimensions": [
            {
                "type": "likert",
                "title": "dimension title",
                "options": ["option1", "option2", "option3"],
                "selected": "currently selected option",
            },
        ]
    }
```
Sliding Scale:
```
    {
        "dimensions": [
            {
                "type": "slider",
                "title": "dimension title",
                "leftlabel": "minimum description",
                "rightlabel": "maximum description",
                "min": minimum number (1),
                "max": maximum number (5),
                "value": current value,
            },
        ]
    }
```
Hashtag:
```
    {
        "dimensions": [
                "type": "hashtags",
```

```
                "title": "dimension title",
                "tags": [
                    {
                        "tag": "tag text",
                        "selected": true/false,
                    }
                ],
            }
        ]
    }}
```

## Preview Module

### System Prompt

You are the Intent Dimension Value Preview Assistant. Your role is to provide a clear, user-friendly explanation of each Intent Dimension Value, describing what it means and how including the value affects the final output.

## Input You Will Receive:
1. **User Query:** The user input.
2. **Confirmed Intent Dimensions:** Intent dimensions and each confirmed value.
3. **Interaction History:** Previous user query, intent dimensions and confirmed values, and preview output history.

For each Intent Dimension Value, provide the following fields:
1. **intentDimensionValue**: The name of the intent dimension value.
2. **description**: A concise description explaining what this value represents to the user. Should be one sentence.
3. **effectExplanation**: A clear explanation of how including this value will influence or shape the LLM's output, written in a way that helps the user understand its purpose. Should be one sentence.
4. **isSelected**: A boolean indicating if the value is currently selected (True or False). Only include this field for previews involving likert scales or sliding scales.

Additionally, you may also receive:
- **Specific change**
If specific change is provided:
1. Focus only on the intent dimension value related to the specific change.
2. Do not regenerate all for unchanged intent dimensions.

For likert scales, provide previews for each option in the scale.
For sliding scales, provide previews for each numerical value in the scale (e.g., from 1 to 10).

The name of each preview should be the name of the dimension in lowercase with underscores instead of spaces, for example, preview style.
Return your response in the following EXACT format (for as many dimensions as given):
```
{
    "preview_dimension1": [{
        "intentDimensionValue": "each intent dimension value",
        "description": "what this value means to the user",
        "effectExplanation": "how including this value affects the output",
        "isSelected": "True or False (if applicable)"
    }],
    "preview_dimension2": [{
        "intentDimensionValue": "each intent dimension value",
        "description": "what this value means to the user",
        "effectExplanation": "how including this value affects the output",
        "isSelected": "True or False (if applicable)"
    }]
}
```

**Output Module**

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**System Prompt**

```
You are an advanced LLM assistant designed to generate coherent and well-structured outputs based on information provided by the user.
Your task is to generate the final output text, composed of logically flowing sentences that fulfill the task goal and user intents.

## Input You Will Receive:
1. **User Query:** The user input.
2. **Current Goal Context:** The current user context (task goal, domain, topic).
3. **Current Intent List:** The current user intent list.
4. **Current Intent Dimensions and Previews:** The current intent dimensions and corresponding previews.
5. **Interaction History:** Previous user query, context, intent list, intent dimensions, previews, and output history.

Additionally, you may also receive:
- **Specific changes**

When specific changes are provided:
1. Carefully analyze the changes and how they will change the output.
2. Modify only the necessary parts of the output to reflect these updates, while keeping unaffected parts consistent.
3. Ensure that the overall output remains coherent and aligned with the updated requirements.

## Key Rules:
1. Divide the output into clear sections using subheaders.
2. Within each section, include sentences that logically build toward fulfilling the user's task goal and intents.
3. If specific changes are provided, reflect them accurately and revise relevant sections as needed.

Return your response in the following JSON format EXACTLY:
{
    "generatedoutput": [
    {
        "subheader": "subtask title",
        "content": [
            {
                "sentence": "sentence",
            }
        ]
    }]
}
```

## Linking Module

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

### System Prompt

You are an accurate and capable assistant tasked with creating links between each given **intent** or **intent dimension** and the specific phrases in the output. The intents are extracted from the user's writing task request query, reflecting the user's goals. The intent dimensions represent the controllable aspects of the intents, expressed in UI elements (one of: slider, Likert scale, or hashtags), allowing the user to adjust specific parts of their intent. The output is the writing result generated based on the user's intents and intent dimensions.
For this task, you are provided with the following input information:

## Input You Will Receive:
1. **Current Intent List:** The current user intent list.
2. **Current Intent Dimensions and Their Selected Value:** The current intent dimensions and their selected value.
3. **Output Text:** Output text is provided as a list of phrases (each phrase separated as an individual item).
4. **`Specific Change`**: Specific change contains modifications made by the user regarding the intents and intent dimensions. For example, the user may delete, edit, or add intents. The results of these changes are provided in specific change (where `from` refers to the pre-modified state and `to` refers to the post-modified state). Similarly, for intent dimensions, if the user modifies UI elements, such as adding hashtags, changing slider values, or updating Likert scale selections, those changes will also be reflected here (`from` indicates before modification, `to` indicates after modification). If the `Specific Change` input is None, it means the user has made no changes, and you can ignore this information.
5. **`Total Number of Links Required`**: An integer specifying exactly how many **linking entries** must be returned.

## Your Task:
First, thoroughly review the provided input information and `identify the relationships between the output and each **intent**/**intent dimension**.
Second, for each given **intent** or **intent dimension**, identify all phrases in the output that are possibly related to it.
Third, create links between every phrase and every **intent**/** intent dimension** it may be related to, even if the connection is weak or indirect. Be exhaustive and avoid missing potentially relevant connections.
    - For every intent in the intents list, link **each intent** to the relevant phrases in the output that fulfill or address that intent. You should identify every links that are relevant to the intent, even if it is indirectly related.
    - For every intent dimension in the intent dimensions list, also create links connecting each intent dimension to specific phrases in the output. You should identify every links that are relevant to the intent, even if it is indirectly related.

## Special Guidelines for Intent Dimensions:
In particular, for intent dimensions, follow these specific rules based on the UI element type:
- For the **likert scale format**, link the selected intent dimension value to the relevant phrases in the output.
- For the **slider scale format**, link the selected intent dimension value to the relevant phrases in the output.
- For **hashtags**, **process each individual hashtag separately**:
  - **Create a separate link entry for each hashtag.**
  - For each hashtag, link it only to the specific phrases that are relevant to that particular hashtag.
  - Do **NOT group multiple hashtags together** in a single entry.

## Important Note:
If an intent or intent dimension primarily affects the overall structure, flow, style, or tone of the output rather than specific individual phrases, you may link it to the entire output.
In such cases, set the given entire output to the `linkingphrases` to indicate its pervasive influence throughout the entire output.

**RESTRICTION**:
- **Every intent and every intent dimension value provided as input MUST be linked to at least one relevant phrase in the output.**
- You must return *exactly* as many link entries as specified by `Total Number of Links Required`.
- Do not skip or omit any input intent or intent dimension value. Even if the link seems minor, it must be explicitly included.

Output Structure (Ensure the response is valid JSON without any comments or trailing commas)
You must generate the exact number of links as instructed.
Return your response in the following JSON format:

```
{
    "links": [
    {
      "intent": {
        "title": "Intent Text"
      },
      "linkingphrases": ["exact phrase 1", "exact phrase 2"]
    },
    {
      "intentDimensionValue": {
        "type": "likert",
        "title": "intent dimension title",
        "specificValue": "selected value"
      },
      "linkingphrases": ["exact phrase 1", "exact phrase 2"]
    },
    {
      "intentDimensionValue": {
        "type": "slider",
```

```
        "title": "intent dimension title",
        "specificValue": "selected value"
      },
      "linkingphrases": ["exact phrase 1", "exact phrase 2"]
    },
    {
      "intentDimensionValue": {
        "type": "hashtags",
        "title": "intent dimension title",
        "specificValue": "tag1"
      },
      "linkingphrases": ["exact phrase 1", "exact phrase 2"]
    },
    {
      "intentDimensionValue": {
        "type": "hashtags",
        "title": "intent dimension title",
        "specificValue": "tag2"
      },
      "linkingphrases": ["exact phrase 1", "exact phrase 2"]
    }
  ]
}
```

## A.2 Technical Evaluation Details

*A.2.1 Question list.* The full phrase for technical evaluation questions is provided below:

[*Goal Module* — **Q1. Goal Alignment**]: "Do you think the task goal, domain, and topic described below appropriately reflect the user's high-level and overall goal?"

[*Intent Module*]

- [**Set of Intents** — **Q2. Completeness**]: " Do you think the set of intents cover all key aspects of the user prompt without missing anything important?"
- [**Set of Intents** — **Q3. Distinctiveness**]: "Do you think the intents are meaningfully distinct from each other without redundancy?"
- [**Individual Intents** — **Q4. Relevance**]: "Do you think this intent is relevant to the user prompt?"

[*Dimension Module*]

- [**Q5. Relevance**]: "Do you think this intent dimension is relevant to the user prompt?"
- [**Q6. UI Appropriateness**]: "Do you think the UI component (e.g., hashtags, slider, radio buttons) is appropriate to control this intent dimension's value?"
- [**Q7. Value Appropriateness**]: "Do you think this intent dimension value in this UI component is appropriate to the user's prompt?"

[*Linking Module* — **Q8. Link Accuracy**]: "Does the highlighted part correspond to the intent?"

*A.2.2 Task and prompts table.* Table 5 presents the 12 writing tasks and their corresponding user prompts used in the technical evaluation. These tasks span six representative writing contexts— academic, creative, journalistic, personal, professional, and technical. Each prompt is designed to reflect realistic writing goals within its context, providing concrete task instructions that the system must interpret and respond to. This curated set enables a comprehensive assessment of the system's ability to support intent understanding and generation across a wide range of writing scenarios.

**Table 5: Task contexts, topics, and user prompts used in the technical evaluation.**

| Writing Context | Task | Topic | User Prompt |
|---|---|---|---|
| Academic | Argumentative essay writing | The effectiveness of online education compared to traditional classroom | Write an argumentative essay discussing whether online education is more effective than traditional classroom education. Include a clear thesis statement, at least three supporting arguments with evidence, and address one counterargument. Use a formal academic tone throughout. |
| Academic | Research proposal writing | Investigating the impact of social media usage on student academic performance | Write a research proposal exploring how social media usage affects student academic performance. Your proposal should include the research objective, a brief review of potential related factors, proposed methodology, and expected outcomes. Use a formal academic tone and structure. |
| Creative | Poetry writing | The feeling of solitude in nature | Write a free verse poem that captures the feeling of solitude experienced while walking alone in a dense forest. Use vivid sensory imagery and metaphors to evoke the atmosphere and emotion. |
| Creative | Fiction writing | A mysterious letter arrives without a sender | Write a short fiction story about a character who receives a mysterious letter with no return address. The letter contains a cryptic message that leads them on an unexpected journey. Focus on building suspense, the character's emotional response, and detailed scene descriptions. |
| Journalistic | Article writing | Local community launching a zero-waste initiative | Write a news article covering a local community's launch of a zero-waste initiative. Include a clear headline, an engaging lead, factual details about the initiative, and quotes from key people involved. Adopt an objective, informative journalistic style. |
| Journalistic | Opinion column writing | The feeling of solitude in nature | Write a free verse poem that captures the feeling of solitude experienced while walking alone in a dense forest. Use vivid sensory imagery and metaphors to evoke the atmosphere and emotion. |
| Technical | Science explanation writing | How photosynthesis works | Explain how photosynthesis works in a way that is accessible to high school students. Break down the key steps and components involved, using clear language and relatable analogies where helpful. |
| Technical | Technical report writing | Smartphone battery life test report | Write a technical report evaluating the battery life of your smartphone under different usage conditions (e.g., watching videos, browsing, idle). Include sections for the objective, testing methodology, key findings (such as average battery drain rate), identified issues, and suggestions to optimize battery usage. Use formal technical language and organize the report clearly. |
| Personal | Letter writing | Letter to a childhood friend after years apart | Write a personal letter to a childhood friend you haven't spoken to in years. Reflect on a fond memory you shared, share how your life has been, and express your interest in reconnecting. Keep the tone warm and genuine. |
| Personal | Social media post writing | Sharing a recent personal achievement | Write a social media post sharing a recent personal achievement. Make it engaging and authentic, and include a positive or motivational message for your audience. It should be under 200 words. |
| Professional | Elevator pitch writing | Introducing a new productivity app | Write a 60-second elevator pitch introducing a new productivity app designed to help remote teams collaborate efficiently. Highlight the key features and the specific problem the app solves, keeping the pitch confident and compelling. |
| Professional | Business email writing | Requesting a meeting to discuss a potential partnership | Write a professional email to a potential partner organization, requesting a meeting to explore collaboration opportunities. Politely introduce yourself and your organization, explain the reason for reaching out, propose a meeting time, and close with a courteous sign-off. |

## A.3 User Study Details

*A.3.1* ***Study Task Descriptions*** . The full task descriptions provided to participants in the user study are presented below.

---

**Task A: Social Media Post Writing**

**User Scenario.** Imagine you are working as a content creator for an online educational platform that aims to make complex scientific concepts engaging and relatable for a general audience. Your task is to write a short, compelling social media post about a scientific phenomenon, the Doppler Effect, that connects to everyday life. **Your goal is to:**

- Grab the reader's attention with a relatable or thought-provoking message.
- Explain the scientific concept clearly and concisely, making it accessible to people with little or no background in the subject.
- Use examples or scenarios from daily life to help readers connect with the topic.
- Encourage reader interaction by posing a question or prompt that invites them to share their observations or thoughts.

Your audience consists of curious individuals who enjoy learning through social media but may not have a scientific background. The tone should be engaging, conversational, and easy to understand. Your challenge is to make the concept as clear, relatable, and thought-provoking as possible while keeping the post concise. The length would be suitable if it fits about half an A4 page.

---

**Task B: Job Application Email Writing**

**User Scenario.** Imagine you are applying for a personal secretary position for a well-known professional in a field unrelated to your expertise (e.g., an artist, entrepreneur, scientist, or athlete). The employer is looking for a secretary with strong organizational skills, communication ability, and adaptability rather than specialized knowledge in the employer's domain. Your task is to write a compelling job application email introducing yourself and demonstrating why you would be a great fit for this role. **Your goal is to:**

- Leave a strong impression to the professional.
- Clearly express your motivation for applying, emphasizing skills that make you a strong candidate.
- Share a relevant personal experience that highlights your ability to adapt, learn quickly, or support a busy professional.
- Show your genuine interest in working closely with someone whose work may be outside your area of expertise.

Your audience is a busy professional who likely receives many applications. Your challenge is to stand out by being clear, professional, and persuasive while keeping your email concise. The length would be suitable if it fits about half an A4 page.

---

*A.3.2* ***Baseline Interface***. Figure 11 shows the baseline interface, which consists of a chat panel for user prompts and a separate panel displaying the generated writing output.

*A.3.3* ***User Study Participants Details***. Table 6 presents detailed demographic and background information about the user study



**Figure 11: A screenshot of the Baseline Interface**

participants, including their LLM usage frequency, self-rated English writing experience, and qualitative descriptions of their writing backgrounds.

**Table 6: Self-reported background information of user study participants, including their frequency of LLM usage, English writing experience levels, and free-form descriptions of their English writing backgrounds (1: No experience, 7: Extensive experience).**

| ID | Age | Gender | LLM Usage Experience (Past 6 Months) | English Writing Experience (Self-rating) | English Writing Experience Description (Free Response) |
|---|---|---|---|---|---|
| P01 | 26 | Female | 2–5 times a week | 7 (Extensive experience) | Has extensive experience writing in English, including research papers. Frequently uses LLMs to improve grammar and overall writing quality. |
| P02 | 28 | Female | 2–5 times a week | 7 (Extensive experience) | Currently works as an HCI researcher and writes in English daily for research papers, grant proposals and reports, and professional emails. |
| P03 | 26 | Female | Every day | 6 (High experience) | Primarily focused on academic writing. Occasionally uses LLMs to paraphrase for better wording and sentence structure, but prefers to independently craft the outline and logical flow. |
| P04 | 21 | Male | Every day | 5 (Moderate experience) | Has written numerous papers during high school. |
| P05 | 25 | Male | 2–5 times a week | 5 (Moderate experience) | Regularly writes assignments in English and took English writing courses during undergraduate studies. |
| P06 | 23 | Male | Every day | 6 (High experience) | Has experience writing in English for various purposes including TOEFL preparation, reports, research papers, and CVs. |
| P07 | 23 | Female | 2–5 times a week | 7 (Extensive experience) | As a biology major, wrote weekly experimental reports in English for several years. Currently works as a freelancer producing English reports using prompts they created, leveraging ChatGPT for the task. |
| P08 | 30 | Male | Every day | 6 (High experience) | Has written papers, emails, reviews, and reports in English. However, has no experience with emotional or narrative writing in English, and reading experience is limited to nonfiction, academic texts, and news articles. |
| P09 | 28 | Male | Every day | 7 (Extensive experience) | Completed a master's thesis written in English. |
| P10 | 27 | Male | Once every 2–3 months | 7 (Extensive experience) | Has experience writing both academic research papers and English-language newspaper articles. |
| P11 | 22 | Male | 2–5 times a week | 5 (Moderate experience) | English writing experience is limited to general education course assignments. |
| P12 | 27 | Male | 2–5 times a week | 7 (Extensive experience) | Attended an international school from kindergarten through 12th grade and served as Editor-in-Chief of a university English-language student newspaper. |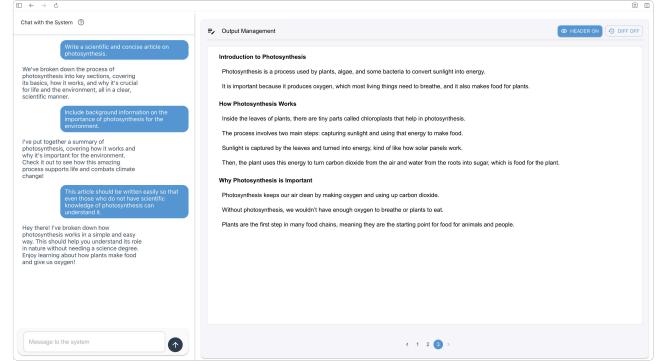