

Beyond Code Generation: LLM-supported Exploration of the Program Design Space

J.D. Zamfirescu-Pereira
zamfi@berkeley.edu
UC Berkeley
Berkeley, CA, USA

Eunice Jun
emjun@cs.ucla.edu
UCLA
Los Angeles, CA, USA

Michael Terry
michaelterry@google.com
Google DeepMind
Cambridge, MA, USA

Qian Yang
qianyang@cornell.edu
Cornell University
Ithaca, NY, USA

Björn Hartmann
bjoern@eecs.berkeley.edu
UC Berkeley
Berkeley, CA, USA

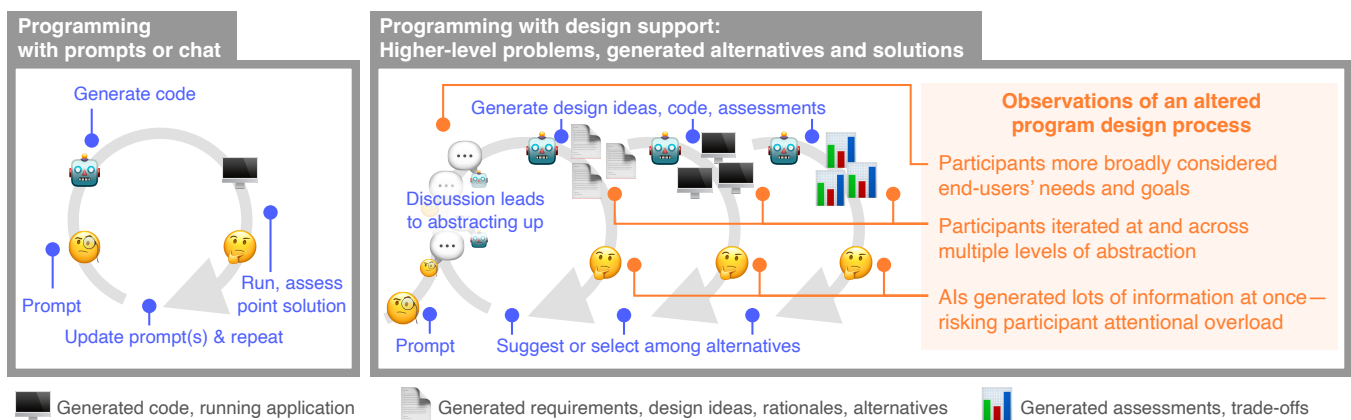


Figure 1: Today, LLM-powered programming assistance resembles a repeating cycle of “prompt”-to-“code” (left) by re-running generation with every prompt change. **PAIL**, our IDE (right), instead encourages developers to abstract up, pulling them towards a deeper understanding of the problem space through an exploration of alternative problem formulations and solutions. Throughout this exploration, **PAIL** tracks design goals, changing requirements, and surfaces implicit decisions—but the breadth and depth of information **PAIL** generates can be overwhelming.

Abstract

In this work, we explore explicit Large Language Model (LLM)-powered support for the iterative design of computer programs. Program design, like other design activity, is characterized by navigating a space of alternative problem formulations and associated solutions in an iterative fashion. LLMs are potentially powerful tools in helping this exploration; however, by default, code-generation LLMs deliver code that represents a particular point solution. This obscures the larger space of possible alternatives, many of which might be preferable to the LLM’s default interpretation and its generated code. We contribute an IDE that supports program design through generating and showing new ways to frame problems alongside alternative solutions, tracking design decisions, and identifying implicit decisions made by either the programmer or the

LLM. In a user study, we find that with our IDE, users combine and parallelize design phases to explore a broader design space—but also struggle to keep up with LLM-originated changes to code and other information overload. These findings suggest a core challenge for future IDEs that support program design through higher-level instructions given to LLM-based agents: carefully managing attention and deciding what information agents should surface to program designers and when.

CCS Concepts

• **Human-centered computing** → **Empirical studies in HCI**; **Natural language interfaces**; **Interaction design**; • **Computing methodologies** → **Artificial intelligence**.

Keywords

Program design, Code generation, Design space exploration, Generative AI, LLMs

ACM Reference Format:

J.D. Zamfirescu-Pereira, Eunice Jun, Michael Terry, Qian Yang, and Björn Hartmann. 2025. Beyond Code Generation: LLM-supported Exploration of



This work is licensed under a Creative Commons Attribution-NonCommercial-NoDerivs International 4.0 License.

the Program Design Space. In *CHI Conference on Human Factors in Computing Systems (CHI '25), April 26–May 1, 2025, Yokohama, Japan*. ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/3706598.3714154>

1 Introduction

A common vision of future computer programming relies on Large Language Models (LLMs) to do it all: identify requirements, write code and tests, perform simulated QA testing, deploy, etc. We are making rapid progress towards this vision on many fronts [1, 15, 20, 27, 57]. Yet the role of the human in this process is under-explored. In the predominant view, humans will be limited to setting an initial goal, possibly clarifying a few questions asked, and then rating the output. For some changes to the goal or generated code, the process will restart from, approximately, scratch. Many versions of this vision exist: some scale up the “autocomplete,” CoPilot-like [23] interactions relying on recognition-over-recall [49], others scale up the “chat” [1] or focus on providing point solutions [27].

While these visions best suit an understanding of program design as starting with one known goal, program design in fact involves a design process iterating on both the implementation details and design goals. For instance, a journalist may begin an analysis with the goal of assessing a relationship between two variables but then in the process of exploring possible visualizations realize that one of the variables has a lot of missing data. As a result, the journalist can pivot to assess another pair of relationships or consider alternative visualizations and statistical analyses to triangulate the initial relationship of interest. Similarly, a game creator faces questions of narrative, characters, interactions. Even a straightforward backend engineering task (e.g., storage system selection) requires consideration of trade-offs (e.g., latency, consistency, scale) that ultimately require refining the higher-level goals of the task (e.g., terms of a latency service agreement).

Viewing programming as a design activity illuminates the human’s need for tighter iterative loops and more fine-grained control than implied by the predominant vision. Design activity begins with consideration of different *problem* formulations before solutions are even discussed—as captured iconically by the *Double Diamond* (see, e.g., [17]), which we adapt to illustrate our targeted capabilities, in Figure 2. Designers test high-level hypotheses (e.g., around user needs, or problem urgency) with sketches and prototypes before expending engineering effort. All these explorations feed back into a better understanding of the problem and solution design spaces—but take time, and are often limited by human availability and cost. Temptingly, LLMs promise a near-free pool of “cognitive resources” that could in theory complement human cognition to improve the experience of designing and programming interactive artifacts and improve the final artifacts’ overall quality—by using LLMs to complement human limitations on time and cost, and using humans to overcome LLMs’ limitations on domain expertise and end-user understanding. We explore this promise in this work.

Through a formative study of creating small interactive sketches using ChatGPT¹ and Claude,² we find that (1) working through

a user-centered design process in chat alone leads to lost requirements and a lot of scrolling through to find prior questions, decisions, and discussions, and (2) individual point solutions—without consideration of alternatives—lead to significant anchoring bias. Combining these findings with knowledge from prior work in AI assistance for design, we develop the PAIL IDE.

PAIL embeds a chat agent with two other LLM-based agents into an integrated environment to facilitate higher-level consideration of design choices while programming with LLMs, enabling users to **design** and **iterate** on programs while working at higher levels of abstraction than the code itself. At the highest level, PAIL elicits design alternatives, rationales, and prompts to guide users toward user-centered design principles. PAIL also speculatively proposes and explores alternative problem and solution formulations, generating interactive prototypes to support epistemic goals through testing and other assessment. At the most granular level, PAIL tracks requirements and implicit decisions made by the (black-box) LLM without user input, surfacing abstractions to establish common ground and clarify understanding. As an IDE, PAIL also allows users to directly edit the code itself, providing full control over implementation details.

Through a lab study with 11 participants, we use PAIL as a probe to understand the benefits, limitations, and some future challenges for LLM-powered programming tools’ support for program design. PAIL helps participants consider their audience and communication goals; participants express appreciation for having a direct summary and for the ability to manipulate those decisions in situ, and 100% of participants find at least one unconsidered alternative that influenced their design work. Further, we find that rapid code generation allows actual executable interactive programs to fit into a *sketch* role: a cheap, “disposable,” [14] epistemological artifact, rather than a *work-in-progress* prototype representing an investment of time and resources—an ability that is suggestive of an altered program design process. Finally, our findings suggest an emerging issue in managing user attention: as more agents and UI affordances lay claim to being helpful and are integrated into developer environments and workflows, these tools will also need to competently balance user attentional capacities and desire for agency. We discuss these observations’ implications for future tools supporting programming and other design activities, and suggest opportunities for further research.

This paper makes three contributions: First, it describes PAIL, a prototype IDE that introduces a new set of interactions extending chat to support program design activity at higher levels of abstraction. Second, it offers a rich description of how and when PAIL’s design support is useful (and not) for creating interactive software prototypes. Lastly, it identifies core upcoming challenges for a future of programming in which programmers relay abstract instructions to LLM agents that synthesize code: these agents will, critically, need to manage and direct user attention, keeping users abreast of AI-initiated changes and carefully considering what information to show.

2 Related Work

Our work draws on empirical studies and theories of design and recent work on LLMs for creative tasks. Here, we review prior

¹<https://www.chatgpt.com/>

²<https://www.claude.ai/>

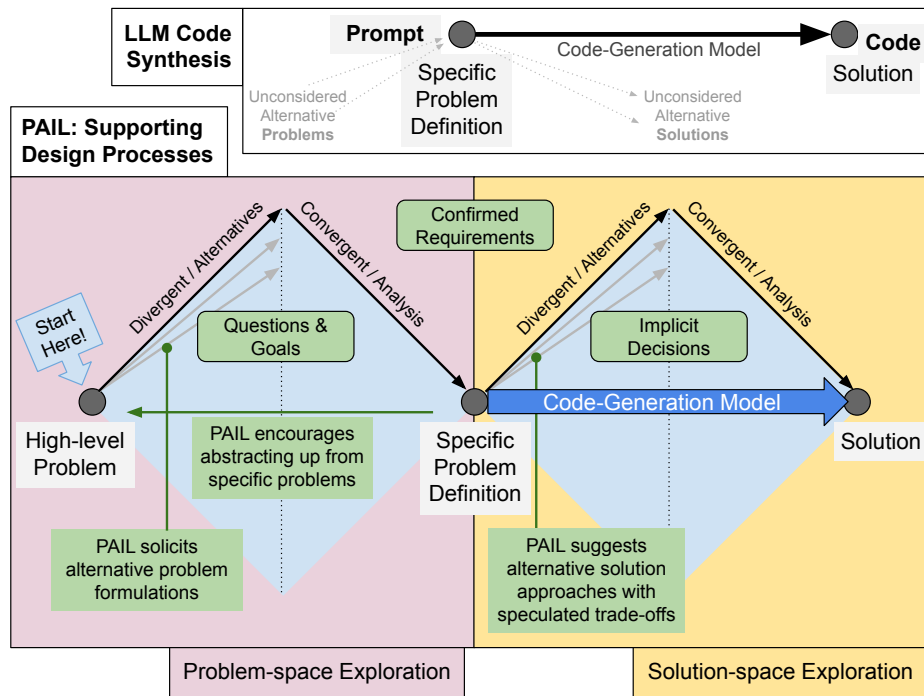


Figure 2: PAIL’s capabilities overlaid on the UK Design Council’s Double Diamond [17]. Building beyond a “prompt”-to-“code” view of LLM code generation (top), PAIL (bottom) embeds this capability within an IDE that supports developers through a larger design process: by abstracting up from a specific problem to a higher-level one, exploring alternative problems and alternative solutions, tracking design goals and requirements, and surfacing implicit decisions (green boxes).

work on systems for sketching and comparing alternatives within design processes; the nascent area of LLM-powered programming assistance as well as pre-LLM approaches; and, lastly, approaches to making better use of LLMs in structured complex tasks.

2.1 Program Design Processes; Sketching and Comparing Alternatives

In the fields of design and HCI, the process of designing is characterized by exploring the space of solutions and iteratively reformulating the problem to solve. Sketching, a form of rapid prototyping for quickly exploring the essential dimensions of a solution space, is a widely accepted, studied, and practiced tool for thought [14, 37]—serving in part to help designers explore alternatives, sometimes even in parallel [21]. Described by Buxton as distinct from *prototyping*, *sketching* aims to evoke, rather than describe. A prototype is often constructed for a specific goal, such as to explore how an artifact might work, look and feel, or relate to other components of a system [29], and can be considered a “work in progress” in the sense that it is meant to faithfully reproduce that aspect of a final artifact, perhaps even to confirm a particular design decision, and as such represents an investment of effort (and potentially materials). A sketch, in contrast, is meant to be quick, disposable and plentiful, suggestive, and serves more of a “cognitive offloading” role to build a shared understanding [14].

In divergent, exploratory phases of the program design process, programmers also often engage in *Exploratory Programming* [9, 56],

writing code as a *medium to prototype* with different ideas, with an open-ended goal and no existing specification. In such a phase, programmers may engage in *Opportunistic Programming* [12, 13], a form of programming in which speed and ease are prioritized over robustness, maintainability, or other engineering goals. Copied code found on the web is a common hallmark—and today’s LLMs’ code generation capabilities provide a new venue for finding code to incorporate into a prototype.

Considering and examining alternatives, in these phases an explicit activity forming part of the design process, has a long history in HCI too: recent examinations in the context of prompting text-to-image models like PROMPTIFY [11] and DREAMSHEETS [3] echo prior systems like Kery et al.’s VARIOLITE and VERDANT for computational notebooks [33, 34], Hartmann et al.’s JUXTAPOSE system for exploring alternatives in parallel in code [26], and Terry et al.’s *Parallel Paths* approach for parameterized vector art [63]. These systems in turn trace their lineages back to Marks et al.’s Design Galleries [46] for expensive-to-generate computer images, and are complemented by a set of theory-driven work aimed at understanding how and why these processes help designers [45, 52]. These systems and studies all point to parallel exploration and tracking of alternatives as a well-established, essential process in exploratory programming [32].

2.2 AI Assistance for Programming and Design

Our work here also builds on two major threads of research in AI assistance: assistance for programming, and assistance for design. On the programming side, the biggest impact comes from autocomplete-focused assistance through tools like CoPilot [23], which has been extensively studied [7, 19, 51] and shown to serve both “recognition-over-recall” [49] and epistemic (e.g., “oh, I didn’t know about list comprehension in Python!”) goals [7]. Other work has explored how LLMs can be used in the service of design processes for creative coding [5], data analysis [24, 36, 42, 47], and even, circularly, for the assessment of LLM outputs themselves [55].

Pre-LLM work in both AI assistance [16] and the usability of code synthesis [30, 44] is also relevant here, offering insight into the kinds of assistance programmers are looking for: support for writing mundane boilerplate or glue code, for reasoning, and for rapid iteration. Traditional code synthesis has found uses in a number of ways, most notably for HCI through a line of work on programming by demonstration—often repeated demonstrations following a feedback-driven design process, in the service of spreadsheet formula construction [25] and web scraping and automation [8, 40].

This empirical work is complemented by a set of theoretical, speculative, and design-oriented work around the opportunities and challenges of effectively instructing AI systems and designing both with and through them. These include studies of prompting [43, 66, 67], impacts of current LLMs on creative design processes [4, 59] speculative future design processes [64], new conceptual models (e.g., [58, 62]), and questions of agency (e.g., [38]) and perception (e.g., [35]). This body of work emphasizes the extent to which iteration is critical to design, especially when working with black-box models, that evaluating LLM outputs for correctness is challenging due to intrinsic unpredictability, and that steering LLMs benefits from understanding how LLMs go about performing the tasks a user asks them to perform.

2.3 Workflows Integrating LLMs in Complex Tasks

In order to take advantage of LLMs’ capabilities to perform simple tasks well and apply them to more complex tasks, researchers have explored workflows that integrate LLMs. Wu et al. proposed chaining as a technique for connecting LLM calls to each other [65]. Building on this interaction model, Arawjo et al. developed ChainForge [6], an interface for composing LLM prompts and assessing the results of prompts. Exploring alternatives to linear composition, Kazemitabaar et al. [31] compare two different forms of task composition: (i) phase-wise decomposition which batches steps together and (ii) step-wise decomposition which iterates on each step piecewise. Across this work, a common finding is that systems need to scaffold LLM usage in order for users *across experience levels* to make the most of an LLM’s capabilities.

This set of work also shows us that individual interventions to aid in design, in organization, in evaluation, and in grounding assumptions can all be helpful. Here, we aim to shed some light on what challenges will arise when we begin to build more complex cognitive support tools that integrate multiple of these affordances into a single system. Indeed, research in the AI community has shifted towards developing cognitive architectures for coordinating

multiple agents [60], even to complement each other’s strengths and weaknesses [27].

3 Designing PAIL

Through this work, we aim to understand how explicit design support can impact programmers’ design and prototyping processes. We adopt a Research through Design [68] approach and develop PAIL, a web-based IDE for programming p5.js sketches, to probe into what programming with LLMs could look like in the future.

3.1 Design Process

We developed PAIL iteratively based on the described prior work and a formative study. Throughout our design process, we discussed and incorporated feedback among the co-authors.

Two practical considerations constrained our tool and study design spaces upfront; these led us to focus primarily on experienced programmers making small interactive prototypes for personal or (non-deployment) professional use:

First, the current fidelity of LLM code synthesis limits the complexity of the programs that can be effectively constructed or iterated upon by these models with limited human intervention—but we expect new models to continue to develop these capabilities. To avoid having our observations unduly influenced by today’s model failure modes, we target smaller programs that can be fully included in LLM context windows, in a language (JavaScript) and using a framework (p5.js) that are well-represented in LLMs’ training data.

Second, we are interested specifically in *design support*, not necessarily in supporting end-users learning how to program—there is plenty of great work in that area already. Thus, we primarily target participants with a strong understanding of programming.

3.2 Formative Study

We conducted a formative study to identify challenges that emerge when following common design processes for program design using an LLM-powered chatbot. We recruited five participants from within the HCI and Design-focused groups at our institution, with a range of experiences creating interactive software prototypes. All participants were students in their final two years of undergraduate education or doctoral students.

Over each of six weeks, a subset of participants used ChatGPT or Claude to produce p5.js code prototypes for a variety of interactive programs, running those prototypes by copy-and-pasting any generated code into a standard p5.js environment.³ We chose ChatGPT and Claude for our study rather than GitHub CoPilot or other AI copilots because the latter tools do not provide design guidance nor engage in requirement elicitation, and require users to engage with code rather than higher-level abstractions.

Our participants’ interactive programs ranged from basic games and simulations (e.g., a fashion simulator that lets you try on—virtually—images of clothing) to showcases for artwork, aggregators for real estate listings, and other sensemaking tasks. Occasionally, participants were asked to create an interactive artifact *for someone else*, or asked to observe someone not related to the study using ChatGPT or Claude for the first time.

³<https://editor.p5js.org>

Our formative study participants observed that tracking the outcomes of design discussions was a major challenge, as these discussions rapidly disappeared into chat history and became very challenging to find or recall—and, as decisions receded into the chat history further, the models themselves were also less likely to consider them in future iterations of the design. Additionally, both ChatGPT and Claude limited their responses to one of (1) a single point solution paired with a description of that solution, without discussion of trade-offs or alternative solutions, or (2) several possible solutions, paired with a description, but not comparing solutions with each other or discussing trade-offs among them.

Participants expressed an appreciation for being able to quickly generate code to test out ideas, but found it challenging to later identify those experiments or revert to the code used for them. Participants also reported feeling unsure about the code the LLMs generated, looking at it relatively rarely—the common iterative loop consisted of running the generated code and asking the LLM directly to fix any observed errors, rather than reviewing the code directly for errors. Conspicuously lacking was support for an understanding of the generated code, beyond the high-level overviews produced by the LLMs—overviews which lacked mention or explanation of critical decisions or assumptions made by the LLM. Even when these were provided, they were typically buried within paragraphs of overview, requiring a close read to find, and participants almost never read these explanations closely because of a low signal-to-noise ratio.

Lastly, those participants who mediated a program generation exercise for someone else reported spending substantial time identifying what needs that user had, and then synthesizing solutions for those needs. These activities required more concrete epistemic goals in prototyping.

3.3 Design Goals

From our formative study, we distilled a set of four design goals.

DG1: PAIL should provide explicit support for generating, keeping track of, and comparing alternative designs.

DG2: PAIL should extract and track task requirements and decisions explicitly, outside of the primary conversation dialogue, to keep these salient and visible to both the human designer and the LLM that is generating code.

DG3: PAIL should take reasonable steps to help users avoid reading the program code, keeping them at their desired level of abstraction in communicating in terms of functionality and goals—including by surfacing, in natural language, decision points that exist only implicitly in the code.

DG4: PAIL should help users consider the needs and goals of the ultimate users of the programs being designed.

We considered several possible designs for showcasing alternatives (**DG1**) and version histories in early prototypes of PAIL, including a Git-like “commit” structure inspired by Litt [41]. We ultimately selected a lightweight “pull” model influenced by Kery et al.’s observation that, in data science workflows, “versioning” in the traditional software sense can be too heavyweight to be useful [34]. We also decided that, as suggested by Lunzer and Hornbæk [45], PAIL should offer *prospective* (speculative) comparisons of possible

alternatives, not only retrospective comparisons of selected alternatives. As a result, we implemented PAIL to suggest alternatives whenever possible (**DG1**).

To help ensure that user communications at a high level of abstraction were likely to refer to concepts the LLM also had a handle on, and based on recent findings in cognitive science on the discovery and usage of abstractions in conversation, we incorporated an explicit mechanism for grounding specific jargon with a textual description (**DG3**). For example, consider a puzzle game with a specific win condition: explicitly calling out “Win Condition,” with a description of what that phrase means, offers users both a specific name to refer to this concept, and confirmation that both the user and PAIL are using that phrase to refer to the same sort of underlying concept.

A second affordance aimed at a maintaining a level of abstraction higher than code is a set of individual code change summaries, in natural language text, that accompany any code change PAIL suggests. Incorporating this explicit mechanism for tracking code changes at a higher level allows users to avoid the “context switch” of shifting their thinking mode from design considerations, to code, and back again.

Lastly, explicitly surfacing “implicit” decisions (necessarily) made by the LLM when it synthesizes code in response to an ambiguous request provides a mechanism for considering these decisions and generating possible alternatives (**DG2**, **DG3**). These decisions can be critical to a design, and include decisions about how to represent data, like user progress through an application, or how to operationalize certain constraints, like how to validate a “win condition” in a game. In our formative study, we found that unsurfaced decisions were discussed only when these decisions had a visible impact the user noticed and chose to inquire about, but constrained the future design space regardless—so we chose to make them explicit here.

3.4 System Design

What emerged from our design process was the need for two sets of entities: a set of “agents” that engage with the user directly (such as the CONVERSATIONAGENT) or indirectly through UI affordances (such as the DESIGNAGENT and REFLECTIONSAGENT). These agents are complemented with a set of “views” (see below) aggregated into a “design panel” on the right-hand side of PAIL (see Figure 3).

The design panel is PAIL’s primary differentiating feature offering a consistent interface to four design aids, represented as four subsections to the panel, shown in Figure 3:

- (1) **DQs Design Questions & Goals**: this section tracks the kinds of questions an interaction designer would ask when designing a new interactive piece, including problem formulations,
- (2) **REQs Confirmed Requirements**: this section tracks decisions that the human has made or confirmed.
- (3) **DECS Implicit Decisions**: this section identifies and surfaces decisions the LLM has implemented in the code without explicit confirmation, such as choices regarding how data is structured, or how displayed values are calculated.
- (4) **UABS Useful Abstractions**: this section offers grounding terminology, alongside one-sentence descriptions, for the

user and the AI to validate that they are referring to the same concepts when they use project-related language.

Each of these subsections contains a list of items generated by the REFLECTIONSAGENT in response to the ongoing conversation in the chat panel, described below. Each item in the design panel includes a justification (“rationale”) for that item as well as 2-3 possible “alternatives.” Any alternatives, except some that the REFLECTIONSAGENT identifies as important, are hidden behind an accordion-style view, revealed only when the enclosing item is clicked.

In one study participant’s project, within the DQs Design Questions & Goals section, appears the item “What themes or characters would be most engaging for a 6-year-old?” accompanied by the rationale “Engaging themes and characters can make learning more enjoyable and effective.” and the alternatives “Animals”, “Superheroes”, and “Fantasy worlds with magical creatures”. Finally, each of these alternatives, in turn, is then speculatively assessed by the DESIGNAGENT, which populates each alternative with a sentence about possible cost/benefit trade-offs of that choice. For “Fantasy worlds with magical creatures,” for example, the trade-off text is “Stimulates imagination; may distract from educational content.” A TRY button next to the alternative triggers the DESIGNAGENT to go ahead and make that change. A screenshot of the four sections of the design panel appears in Figure 4. (An additional REVERT button appears in response to clicking TRY, which restores the original code in the IDE; not shown.)

The design panel is located as the rightmost of PAIL’s four main panels (see Figure 3): a code panel, an output panel (with a console for errors and printed output), a chat panel, and the design panel. The code panel allows users to view and directly modify any project code, and uses a common in-browser code editor, Monaco,⁴ containing p5.js⁵ code that is then run and displayed in the output panel.

The chat panel is the interface to a prompted GPT-4o-based chatbot (CONVERSATIONAGENT) that can read the project code, patch it, or entirely replace it. Changes to code are summarized within the chat, displayed in “diff” form, and then propagated directly to the code panel, which shows a holistic “diff” over the prior iteration of the project. The CONVERSATIONAGENT itself is prompted to explicitly guide users through a focused User-Centered Design [2] process: identifying target users, evaluating their needs, assessing possible goals for the project given users’ communicated design ideas, and finally generating code for prototypes to test any hypotheses generated through this conversational design process.

The design panel and chat panel operate in a tightly integrated way: conversations in the chat trigger the REFLECTIONSAGENT to update the design panel’s contents, while manipulations in the design panel (i.e., trying a specific alternative) trigger code changes via speculative (i.e., uncommitted) executions through the CONVERSATIONAGENT by the DESIGNAGENT.

3.4.1 System Design Non-Goals. Equally important to the design of PAIL is what we did *not* include. Because our focus is on design process support, PAIL does not provide assistance with debugging or handling nonfunctional LLM-synthesized code, nor any kind of

automated QA or simulated user testing, two areas worth considering for an AI-assisted IDE. We leave investigation of these topics to future work.

3.5 PAIL Implementation

PAIL is implemented as a React single-page application, proxying calls to GPT-4o through a node.js backend that logs all requests. The code editor uses Microsoft’s Monaco⁶ editor configured to display code differences inline. The chat component is a custom-built turn-taking component for user communication with the CONVERSATIONAGENT; it shows chat messages as well as summaries of any code changes made by the CONVERSATIONAGENT alongside diffs with any changed code.

We chose the Monaco editor (a component from VS Code) specifically because of its overall flexibility and support for showing “diffs” of code inline. This code is run using the p5.js library run alongside an HTML Canvas element enclosed within an iframe element, using the loop-protect library⁷ to prevent infinite loops from interrupting event processing in the browser. The iframe’s console object (used to report errors to developers) is proxied through the browser’s postMessage API to allow p5.js errors and other console messages to pass through to the IDE.

A few technically challenging parts of PAIL’s implementation are worth noting.

3.5.1 Speculative Assessment. The code and design panel JSON structures are stored as separate “artifacts” in versioned flat files by the node.js server. To improve latency and offer an *a priori* indication of the impact a particular change might make the DESIGNAGENT speculatively assess a subset of identified alternatives to design panel items.

These speculative assessments vary in scope; depending on the nature of the alternative, this could include: (1) how that alternative might impact DQs design goals, such as considerations of target users; (2) what additional REQS requirements that alternative might reveal; and (3) any new DQs design questions that may emerge from consideration of that alternative. Figure 3 shows how agents and artifacts influence one another.

3.5.2 Differential Updates to Code. At the time PAIL was created, LLMs excelled at generating full files worth of code, but did not have a mechanism for incrementally updating existing code. Because PAIL is premised on rapid iteration across levels of abstraction, we decided it was important to have such a mechanism, so that users could ask for a small tweak (e.g., “can we make the small rectangles blue?”) and have that change implemented quickly. Other systems like Claude’s *Artifacts* and ChatGPT’s *Canvas* rewrite the full code from scratch each time, which is not only slow, but also often includes other, unrelated changes to the code—a complaint we also heard in our formative study.

Obvious solutions we tried, like asking an LLM to output a code change in *diff* or *unified diff* format, failed frequently: line numbers wouldn’t match up, edits were interleaved with both old and new code, and new code was sometimes inserted in the wrong place altogether. These errors align with some prior observations of

⁴<https://microsoft.github.io/monaco-editor/>

⁵<https://p5js.org>

⁶<https://microsoft.github.io/monaco-editor/>

⁷<https://www.npmjs.com/package/loop-protect>

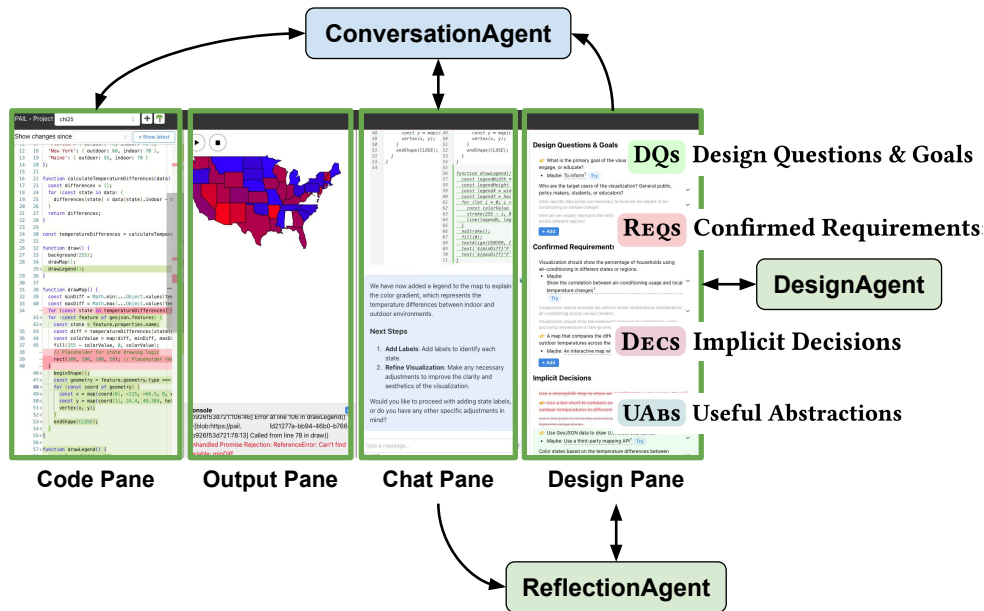


Figure 3: An overview of PAIL’s user interface and related agents. PAIL engages users in defining and refining (DQs) design goals, including target audience and desired impact, while tracking (REQS) confirmed requirements and (DECS) decisions implicit in an LLM’s synthesized code, surfacing (UABS) useful abstractions, and offering explanations and alternatives across the program design space. A full-scale reproduction of the design pane’s four views appears in Figure 4.

LLMs’ failure modes too: arithmetic is not a strong suit; text locality matters and disrupting it (by interleaving new and old code, for example) lowers performance.

PAIL instead provides all active code to the LLM with every request, transforming that code by prefixing L#### to each line, with #### replaced by each line’s consecutive line number—thus providing a built-in line-numbering scheme. Edits, implemented as a tool call, are represented as “old code” using the same L#### prefix, and “new code”, again with a line number prefix. This allows a straightforward extraction of which lines to replace (fuzzy-matched against line numbers and old code lines, for robustness), and locality in the production of new code. With this technique, we found that PAIL succeeds in patching code directly from LLM calls over 90% of the time, with failures readily discoverable—and full-rewrite as a backup. As a bonus, these incremental updates can be streamed to users as well.

3.5.3 Differential Updates to Design Artifacts. In addition to incremental changes to the code, we implemented incremental changes to the four components of the Design Panel. Here, the challenge was primarily one of updating the visualization: these components are regenerated from scratch often, and in phases (first the summary line, then the alternatives, rationales, and highlights)—as with updates to code, we wish to begin displaying this data to users immediately as it comes in, especially when there is not already any data to display.

But when there is already-displayed data to update, PAIL needs to evaluate the new set of, e.g., **REQS** Requirements, which are streaming in as an incomplete JSON object, against a full set of

existing **REQS** Requirements complete with alternatives, rationales, and highlights. This is accomplished by fuzzy progressive prefix-matching of incoming data with existing data. For **REQS**, for example, individual requirements coming in through the stream are compared with existing requirements for equivalency up to a threshold; only when confidence is high that a requirement is new does the visualization update to display the new requirement as its contents stream in from the **REFLECTIONSAGENT**. Meanwhile, prefix-matched requirements (such as an incoming requirement with the incomplete text “Track the num” matching “Track the number of correct and incorrect answers.”) are not updated until changes are detected, avoiding re-rendering of existing requirements as they are streamed in. This technique is then applied, recursively, to the nested alternatives, rationales, and highlights.

4 PAIL Usage: An Example

To make the design and intended use of PAIL concrete, consider the following usage scenario. Sam, an experienced backend system software engineer and parent of two—kids 3 and 5 years of age—is eager to help her 5-year-old, Alex, learn how to read simple words. Sam’s tried what feels like all the iPad apps in the store, but none have clicked with Alex. Sam knows she could write a new app, but she’s not sure she can commit the time: programming in an unfamiliar environment (e.g., a game environment) often has a learning curve with high variance in time required, and there would likely be a lot of upfront preparation required before Sam could show Alex anything actually interactive—time ultimately wasted if Sam’s game ideas don’t appeal to Alex.

Design Questions & Goals

👉 How can we make the tool engaging for children in the age range of 6–10 years? ⌵

- *Maybe:* Gamification elements (e.g., points, badges)? Try

What kind of feedback will be most helpful for the child to learn from their mistakes? ⌵

How can we effectively track and analyze the child's performance to adjust the difficulty of questions? ⬆

Adaptive learning helps focus on areas where the child needs improvement.

- Use a simple correct/incorrect counter? Try
- Implement a detailed analytics dashboard? Try
- Crowdscore performance data for benchmarking? Try
- Replace ⬆

+ Add

Offers comparative insights; may raise privacy concerns.

(a) The **DQs** Design Questions & Goals section.

Confirmed Requirements

The tool should help a child learn and practice multiplication.

- *Maybe:*
The tool should help a child learn and practice all basic arithmetic operations. ⌵

Try

The tool should be engaging and interactive. ⌵

The tool should provide immediate feedback on answers. ⌵

👉 The tool should track the child's performance and adjust the difficulty of questions accordingly. ⌵

~~The tool should track performance for specific types of questions, such as multiplying by 5.~~

The tool should track performance for specific types of questions, such as multiplying by 1–10. ⌵

+ Add

(c) The **REQs** Confirmed Requirements section.

Implicit Decisions

Use a weighted random selection to generate multiplication problems.

- *Maybe:*
Generate problems using a machine learning model to predict the most beneficial problems. ⌵

Try

👉 Track the number of correct and incorrect answers for each multiplication fact. ⌵

Use canvas-based animations to represent multiplication problems. ⌵

+ Add

(b) The **DECS** Implicit Decisions section.

Useful Abstractions

The Child: the user who is learning multiplication.

- *Maybe:*
The User: a neutral term that could apply to any age group. ⌵

Try

Multiplication Problems: the set of questions generated for practice. ⌵

Feedback System: the mechanism that provides immediate feedback on answers. ⌵

👉 Performance Tracker: system that tracks the child's correct and incorrect answers to adjust question difficulty. ⌵

- *Maybe:* Crowdsourced performance data for benchmarking. ⌵

Try

~~Category Performance: tracks performance for specific types of questions, such as multiplying by 5.~~

Category Performance: tracks performance for specific types of questions, such as multiplying by 1–10. ⌵

+ Add

(d) The **UABs** Useful Abstractions section.

Figure 4: A screenshot of the contents of each of the four components of the design panel, from P9. The third question in **DQs** Design Questions & Goals (a) (“How can we effectively [...]?”) is expanded to show the rationale for that question along with three alternatives. The mouse hover over the third alternative (“Crowdscore [...]”) reveals a speculated trade-off for trying that option. In PAIL, these four sections are organized vertically in the design panel. Items with a finger are *important*, dotted underlines allow hover (see (a)) for trade-offs. Green background and red strikethrough indicate recent changes.

With PAIL, Sam engages in the following scenario:

- (1) Sam asks the CONVERSATIONAGENT for help making a game for Alex to practice reading simple words.
- (2) PAIL responds with a chat message containing a few questions; a few seconds later, these appear in the **DQs** Design Questions section of the design panel.
- (3) Sam scans the suggestions in the **DQs** list: about who the user is, what kind of help they need, and what her own goals are for the app, finding the following:
 - *What specific skills should the game focus on? (E.g., letter recognition, phonics, simple words)*—with an explicit suggestion she try *Focus on simple words*.

- *What types of activities or game mechanics would be most engaging?*—with the note that engagement will keep the child motivated.
- *Are there any themes or characters that your child particularly likes that we could incorporate?*

See Figure 4 for a screenshot of a similar set of **DQs** Design Goals.

- (4) Sam considers what she knows about Alex, what kind of support she needs, and more. Ultimately she clicks **TRY** next to the **DQs** Design Question alternative *Focus on simple words*.
- (5) PAIL's DESIGNAGENT notes this suggestion and follows up with another question: *Any kinds of activities your kid enjoys?*

- (6) This time, Sam responds in chat: *maybe a word-picture matching game?*
- (7) PAIL’s CONVERSATIONAGENT offers an idea for an implementation using a drag-and-drop interaction, asking *does this approach sound good to you?*
- (8) Sam finds this suggestion confusing. Drag-and-drop isn’t an interaction Alex is that familiar with, and in other apps it hasn’t worked very well. What are some alternatives? Sam scans the design panel and finds drag-and-drop interaction in the **DECS** Implicit Decisions section of the design panel. She clicks **TRY** next to *Use a tap-to-select interface for matching words to pictures.*
- (9) PAIL’s DESIGNAGENT has speculatively assessed this approach already, and identified that tap-to-select is *simpler for very young children, and can be more intuitive than drag-and-drop*—noting this when Sam was scanning through possible alternatives.
- (10) PAIL moves this item up to **REQS** Confirmed Requirements and produces an initial prototype of the matching interaction: simple words and simple emojis, in a single line. The CONVERSATIONAGENT offers a summary of code changes, and the REFLECTIONSAGENT records that **DECS** *Words and pictures are displayed in a simplified layout.*
- (11) Sam doesn’t love this choice, and clicks **TRY** next to the **DECS** alternative *Display words and pictures in a grid layout.*
- (12) In response to CONVERSATIONAGENT’s generated patch to use a grid layout, PAIL’s REFLECTIONSAGENT also moves this new requirement to the **REQS** Confirmed Requirements section.
- (13) At this point, Sam shows the prototype to Alex to gather feedback: are the pictures understandable? Are the words too complex, or too simple?
- (14) Sam and Alex can then work together to revise the project, incorporating Alex’s feedback about what he likes and doesn’t like about using the game.

Note that nothing here is beyond Sam’s capabilities as a software engineer. She knows how to write code, but is unlikely to find typical “paper prototyping” worthwhile for this use case. PAIL gives Sam the ability to “sketch” with code, overcoming the activation energy required to choose and initialize an environment, figure out a data model and rendering pipeline, or otherwise make a large number of *boilerplate decisions* that must be resolved to have a working program, but serve no other purpose towards Sam’s goals around understanding what will help Alex learn to read.

Instead, using PAIL, in 10 minutes Sam has already learned that Alex adores emoji iconography but the targets are too small for his fingers. Sam and Alex can continue using PAIL to iterate on the game, pursuing completely new directions without losing track of the lessons they learned from earlier prototypes. PAIL’s suggested alternatives enable Sam and Alex to break free of anchoring effects or the “sunk code” fallacy of energy invested in authoring code. Using PAIL, the energy invested is in considering alternatives, exploring the design space with potential users, and testing prototypes—epistemic actions resolving unknowns, which serve a useful design purpose whether they resolve positively or negatively.

5 User Study Procedure

We had two goals in our user study: (1) identify possible ways the various built affordances helped or did not help users (though not prove definitively); and (2) identify likely challenges inherent to AI support for design.

5.1 Interview Protocol

We began each interview by showing a demo use case for PAIL, with an emphasis on what we expect it to be helpful for: designing interactive applications. We showed how to use the chat functionality, explained the tight integration between the chat agent and the code (as mentioned in §3.5, the chat agent can directly modify the code), and then walked users through the various parts of the design panel, explaining each subsection. With each subsection, we showed examples of entries, complete with the rationales, lists of alternatives, and results of speculative execution, and, finally, demonstrated what happened when the user clicks on the **TRY** button. We finished our introduction by explaining our research question, answering any questions from participants, and then moving on to the first task.

We encouraged participants to think of PAIL not only as a *code generation* tool, but also as a *design* tool. We informed participants that they could ask for broad, high-level goals in the chat panel, such as “Help me design an interactive feature that goes along with this article: [pasted article]” (see §5.2, Task 1 below)—that the chatbot was designed to walk them through a design process, considering target users and the participant’s own high-level goals.

Because our primary goal is not to measure the effectiveness of PAIL from a typical “HCI system evaluation” perspective, but rather to learn what the point design that PAIL represents can tell us about the broader design space of AI-assisted programming, we actively encouraged participants to make use of PAIL’s affordances, and gave them suggestions on when and how to do so. This allowed us to observe a broader range of interactions and impacts than a more traditional lab study, or which might otherwise only be visible after participants develop an expertise using PAIL. Naturally, this choice comes with limitations on what we can thus claim.

5.2 Tasks

Our primary criteria for task selection were: ambitious and ambiguous goals; robustness to common LLM failure modes, such as tracking long (or multiple) code files; and enough design focus that we could reasonably observe a diversity of design approaches among our participants.

As mentioned in §3.3, because we are primarily interested in the impact of design support on program design, and not on LLMs’ code synthesis capabilities or ability to support end users in learning programming concepts, we chose tasks that were ambitious but achievable using a length of code that our choice of LLM could robustly handle. These interests also drove selection of our participant pool, discussed below.

We ultimately selected 3 tasks for our participants:

- Task 1: Create an interactive feature to go along with an article about the impacts of air conditioning on migration patterns in high-average-temperature areas.

Task 2: Create a game that helps a young child learn how to read or multiply.

Task 3: Create a simulation to show the effects of medical overtreatment, e.g., recommending screening for specific rare diseases for more people.

For participants who were parents, we started with task 2, and specifically asked for them to consider the needs of their own child and create a game to address one such need, so that we could observe whether having a very specific target user, who is well-known to the designer, has an impact on the design process. Parents were then asked to engage in task 1, if time permitted. All other participants were asked to engage in Task 1 first, then offered a choice of 2 or 3.

Several participants also requested time to engage in an open-ended exploratory task, in which they could experiment with PAIL in service of some personally meaningful goal; we supported this where time allowed.

5.3 Participants

Participants were a mix of 5 academics and 6 professionals, of whom 3 were parents. We recruited participants with varying levels of programming expertise, design experience, and prior LLM use; see Table 1. Our sample size was chosen in line with prior work around formative testing for usability [22, 53]: our goal is to explore the early benefits and likely challenges users encounter when engaged in our task using PAIL, and our experience with pilot users suggested that we would find benefits and challenges quite quickly.

Our participant pool is skewed towards, but not exclusively consisting of, professionals and graduate students in design- and STEM-related fields—and it is certainly not representative of the population at large. We are not making claims about how a specific population does or does not engage in specific behavior, but rather seek to identify PAIL benefits and forecast future challenges among a population that we expect is disproportionately likely to be early adopters of LLM-based tools.

We also skewed our pool towards participants with programming experience because we are most interested in design support, and wanted to avoid confounding factors caused by participant inexperience with code, programming systems, or other parts of PAIL that were not directly related to our research goals.

5.4 Analysis & Evaluation

Participants were instructed to think aloud while engaged in the tasks. We then undertook an exploratory data analysis, transcribing all videos and observing where participants directed their attention, where they made design decisions, and what factors appeared to influence them. We also noted where interviewers intervened. We then compared these observations across participants and categorized their approaches and uses of PAIL using a well-established affinity diagramming process [48], and through the use of service blueprints [10] that document participants' behavior and reflections. We elected to use these methods from HCI and service design because we seek to understand broadly how different participants engage in similar tasks with a new technology—one that does not have an existing set of users with pre-established work practices they are already engaged in.

6 Findings

In this section, we report findings from our participants' use of PAIL. In particular, we find that (1) participants frequently engage in rapid-fire repeated iteration, commonly at higher level of abstraction than code, and often (but not always) through use of PAIL's Design Panel. While in this state, participants treated iterations as disposable "sketches", demonstrating little attachment, and engaged in activities spanning all parts of the "4 D's" of design, from problem discovery and definition through solution development and delivery, sometimes within the same action.

We also find that (2) focused attention was spread quite thinly while participants used PAIL, and between changes to code, new messages from chat, and changes to the design panel contents, awareness of what was happening in PAIL was easily lost and was perceived as costly to regain—and that these perceived costs ultimately influenced where participants directed their attention.

In §6.1, we describe how participants use (or do not use) PAIL's design support to move across abstractions, consider alternatives, and explore the spaces of possible problems and user needs. Then, in §6.2 we examine where PAIL falls short in supporting users' design processes, and why.

6.1 Use of PAIL's Design Support

First, we report on ways in which participants used PAIL, drawn from in-interview reflections and post-session analysis of PAIL use. Across our open-ended program design tasks, participants were quite varied in their approaches.

A common (9/11) initial stumbling block was choosing the first action: though all participants were shown a "project in progress" within PAIL with an overview of its various components (see Fig. 3), most (7/11) were not sure at what level of abstraction to make their first request. Should they think of a solution first, and then request that solution? Should they simply provide a description of the design task (and in the case of article-associated interactive task, the article) directly to the CONVERSATIONAGENT? Only four participants began by asking the CONVERSATIONAGENT about the design task directly, even though we consistently provided guidance that PAIL could help them brainstorm ideas, too, and that they should not feel the need to wait until they had a concrete idea to start making requests.

6.1.1 Abstracting Up & Problem Exploration. Regardless of the nature of the initial request, the CONVERSATIONAGENT would then respond with questions aligned with the design process described above, asking about DQs design goals, target users and user needs, desired outcomes, and often offering some plausible responses for each. For some participants, this would be the first time they would step back to reckon with these factors explicitly; even participants who had thought about what specifically to design often did so without discussing goals or users, but rather brainstorming solutions directly based on what each solution could provide.

All participants interacted with the initial "design process" phase of the CONVERSATIONAGENT, and all were steered, to varying degrees, by being prompted with these questions. Participants who started by suggesting a specific point design would often use these questions as an opportunity to think more broadly. Even senior programmers with design expertise found something to consider

ID	Age	Parent?	Professional Background	Design Experience	Programming Experience	LLM Use
P1	40s	Yes	Design Researcher	Professional	Professional	Infrequent
P2	20s	-	HCI Graduate Student	Professional	Professional	Frequent
P3	30s	-	HCI Researcher	Amateur	Professional	Frequent
P4	30s	-	Software Eng Manager	None	Professional	Infrequent
P5	50s	-	Software Eng Manager	Professional	Professional	Moderate
P6	60s	-	Film & Media Artist	Professional	Amateur	Frequent
P7	20s	-	HCI Graduate Student	Professional	Amateur	Infrequent
P8	40s	Yes	Filmmaker & Academic	Limited	Limited	Infrequent
P9	50s	-	Sound Artist & Lecturer	Professional	Amateur	Infrequent
P10	60s	-	Software Engineer	Amateur	Professional	Moderate
P11	50s	Yes	Software Engineer	Limited	Professional	Moderate

Table 1: Study participants.

in CONVERSATIONAGENT’s prompting, expressing thoughts like “this really broke it down in an interesting way, who is your target audience” (P5). Later, P5 would reflect:

I didn’t have an idea of what I wanted to do when you first presented me with this problem, and so putting this in and then kind of exploring the “oh, i see, it came up with regions”—and I thought “that sounds reasonable let’s start with regions.

The questions could be an awkward fit for some participant-initiated tasks, however. For example, in tasks that were not targeted at user needs, the questions seemed not well-targeted to the participant’s goals. P6 wanted to use PAIL to make a particular “artistic” sketch they had in mind using p5.js, with which they were already familiar. Recalling their reckoning with the questions CONVERSATIONAGENT posed, P6 later reflected:

Yeah I guess the design element, meaning the “design” as the [...] medium that you’re working with here, and having these different questions, requirements, decisions, useful abstractions and things like that, is pretty interesting [long pause] But I like it, I like being turned on my ear, you know, it’s good.

Some participants also wanted to start directly with an example, rather than by thinking through user groups and user needs. One participant, before even using PAIL, rationalized this desire by noting that it was easier to iterate from a single design than to come up with one from scratch.

Typically, participants would then consider these questions, responding either directly to the CONVERSATIONAGENT, or scanning the summary of questions and possible answers in the design panel and then exploring a subset of REFLECTIONSAGENT-proposed alternatives through the DESIGNAGENT. Participants varied in how they responded to these questions when they did engage. For example, P11 treated the design questions and proposed answers as a checklist to select user groups and features from, clicking the **TRY** button on all the options that appeared, upfront, and only then ran the generated program.

Reflection. Recall that one of our major goals with PAIL was to encourage thinking about design goals and questions explicitly—on this count, we succeeded for many participants, but not all. It is clear from our study that some users are likely to want to start directly from a point solution, and only then reconsider design goals, target users, etc.—and future tools should consider how to best serve this population, perhaps by starting from a set of easily-comparable point solutions to reduce anchoring.

6.1.2 A Working Sketch: First Contact & Rapid Iteration. For most participants, the first prototype that instantiates a solution, for a sufficiently-formulated problem statement, is revelatory—exposing a number of mismatches between the participant’s understanding of the project and either the current state of the code, or the agents’ understanding of the project. For example, P9 expected a set of question and answer cards to be shuffled evenly across the canvas, but found them separated by category into distinct question and answer regions. Two other participants (P5, P7) converged on map views with the CONVERSATIONAGENT, but then saw initial prototypes that didn’t include maps *per se*, but rather stylized region diagrams with rectangles and triangles representing world or country regions.

These types of mismatches almost always resulted in a flurry of requests for low-level implementation fixes. Participants typically requested these fixes either directly in natural language from the CONVERSATIONAGENT or by finding a suitable alternative in the design panel and trying it. How long participants spent in this rapid-fire local iteration mode varied, from under a minute (P8) to more than 20 minutes (P7), during which their activities resembled a *flow state* [18]. These participants appeared to have deep concentration, could express what they wanted as next steps rapidly, noted an effortlessness to the repeated iteration, and with minimal rumination. Participants’ think-aloud would often pause in this state.

In this state, participants would often rapidly shift attention across the code, output, chat, and design panel, looking to make sense of what they were seeing, and for how to communicate desired changes most effectively. Those with greater programming expertise, or whose expertise was a closer match with the task

domain, unsurprisingly spent more time looking at the code, but not more time editing it directly—rather, seeking confirmation, in the code, of a “lack of surprises” (P3) to validate that the program was being authored in a way that met participants’ expectations. Meanwhile, those with less domain expertise (e.g., no familiarity with p5.js, or limited familiarity with authoring interactive artifacts) often ignored the code entirely past a certain point (7/11 participants).

Because none of the agents were designed (nor inclined) to break users out of this state, it would often continue until either some insurmountable “blocker” would interrupt, such as a bug the CONVERSATIONAGENT couldn’t fix, or the user reached a point where they achieved whatever prototype they had set out to achieve, and needed to reflect on what to do next. However interrupted, participants would next take stock of where they were, deciding whether to continue making iterative changes, or reconsider the current approach’s suitability towards higher-level goals.

Reflection. This rapid iteration was almost always quite broad, and not limited to any one of the traditional “4 D’s” of the Double Diamond: the *discover*, *define*, *develop*, and *deliver* phases of design. In fact, the same action might serve multiple goals: validating the defined problem while simultaneously progressing towards prototype development, for example, or revealing some new, unanticipated end-user need. To the extent that we expected to support design processes, it appears that rapid iteration with *interactive prototypes* as sketches enables a metaphorical *superposition* of the problem-formulation and solution-exploration phases of design, perhaps enabling more rapid iteration than paper sketches or other traditional *discovery* methods in design. We are not suggesting that these traditional methods of sketching and prototyping are obsolete, but rather noting that these interactive prototypes were being used in a manner similar to an ink-on-paper sketch in an architectural design setting: as disposable, low-investment *sketches* that allow the designer to investigate one facet of a design space without needing to simultaneously resolve decisions across the full design space.

6.1.3 The Design Panel, Alternatives, and Rationales. Nearly all (10/11) participants reported finding the design panel useful. Eight participants appreciated the summarization and tracking affordances that enabled quick scans of the (DQS, REQS) decisions made so far (and updates to those decisions) when new chat messages began to exceed in length how much participants desired to read. Seven participants expressed appreciation for the reporting of rationales and alternatives. Though no participants explicitly directed appreciation towards the DECS implicit decisions or UABS useful abstractions components of the design panel, we nonetheless observed several participants (P3, P5, P8, P9, P11) drawing insight from these components. P2, for example, explicitly drew a comparison with their prior ChatGPT experience:

I did try to use ChatGPT for that and it was...fine. [...] It looked [worse] than the one we just built, and it took me longer. [...] [In ChatGPT] I didn’t have the right mechanism to high-level changes at this level of abstraction.

Scanning the alternatives lists (and ultimately selecting an alternative to TRY) typically emerged after some initial trigger, as when

a participant would pause, unsure of what next step to take—in a broader sense than just making the next single decision. For example, at one point P10 recognized that their approach to visualizing one particular set of interrelated values (in their case, temperature and air conditioning usage) wasn’t going to work, and they weren’t sure where to take the project next. In this and other similar cases, the decisions and lists of alternatives offered a lower-cognitive-load path forward, by allowing participants to *select one of several possible options* rather than *generate a new idea entirely* as a next step, echoing the recognition-over-recall UX design heuristic [49]. Scanning these lists would often result in focus on a single choice, eliciting a reaction like “oh that seems like a good idea” (P4)—or would yield an alternative not directly on the lists, but still inspired by the scan, which participants would then suggest either in the chat or in the alternative lists’ open text field titled “Replace with...”.

A few participants (P5, P7, P11), in “flow state” (see §6.1.2), preferred to use the design panel’s alternatives almost exclusively to guide their exploration, avoiding chat. Asked why, P5 reported that it was easier to answer multiple-choice questions than to repeatedly write messages in chat—treating the lists of alternatives in the DQS design questions section as a sort of “design checklist,” a menu from which to select target users, goals, and more.

Reflection. To the extent that we expected the design panel to encourage participants to consider alternatives *prospectively*, our findings here and in the previous subsection suggest we could do better: participants only rarely actively sought out new problem formulations unless prompted either by PAIL, by an interviewer, or by a realization that their current design approach was not going to lead to success. Instead, they appear to consider alternatives only *retrospectively*, after exhausting a particular, though narrow, line of inquiry.

6.2 Stumbles, Mismatches in Participants’ Design Processes

Beyond participants’ use of PAIL’s specific affordances, we also observed behavior that bears on the design of future AI assistance for program design; here, we detail those observations.

6.2.1 Re-considering Design Problems. Though PAIL explicitly elicits design considerations from participants at the outset, there is no explicit support to bring users to *reconsider* design questions and goals during the implementation process. As a result, few participants explicitly reconsidered the highest-level design directions within PAIL, at least without explicit prompting from an interviewer. In fact, in the “flow state,” many participants would fixate and iterate on small details (e.g., colors, item position, text content) repeatedly. Meanwhile, the CONVERSATIONAGENT was happy to support this low-level iteration for at least as long as interviewers allowed.

But it was *not the case* that participants remained fixed in their beliefs about user needs and which is the right problem to solve—they simply did not reconsider their design goals *explicitly in conversation* with the CONVERSATIONAGENT. Instead, these realizations would come from specific prototypes that yielded specific forms of insight, such as whether particular problem formulation could be compellingly addressed. For example, P10 at one point reflected: can an end-user actually be convinced of a causal relationship between

air conditioning and climate change through a bar chart—or should the interactive feature instead focus on a different climate-related relationship?

For most participants, this reconsideration was more often voiced to the interviewer rather than to the CONVERSATIONAGENT, highlighting the need for explicit elicitation, at least initially.

6.2.2 Content Generation and Overwhelming Information. Nearly every participant at some point commented on the overload generated by PAIL. There is too much data, it is being generated too fast, it's too hard to look at everything, and it's not clear what one should be looking at. As a result, participants reported, much effort was spent figuring out where to look and how to evaluate changes. Running the project was almost always the top choice, but that did not always work. Sometimes, program behavior was not trivial to reproduce; bugs prevented visual output; or the program simply could not be run because the CONVERSATIONAGENT was in the process of updating it, which could take longer than a minute for large changes. While waiting, participants often scanned the design panel, or the “diff” view of the code, to understand the scope of recent changes or consider what steps to take after the code updates have completed. In P5's words: “Like, it's asking me so much in here, I'm not gonna read it every time.”

The main challenge causing a feeling of overwhelming information between the code, chat, and design panel affordances, participants reported, was simply too much data coming in at once. When participants stopped to read through the design panel contents, they could do so without a sense of overwhelm—and even the chat's relatively low signal-to-noise ratio content could be read in this way. Rather, the challenges arose when clicking a TRY button or making a request, which would trigger a cascade of changes in the interface; it was this cascade in particular that participants struggled with.

6.2.3 Application of Reflection and Design Agents. Third, participants did not appear to make substantial use of the rationales the REFLECTIONSAGENT and DESIGNAGENT provided for why certain decisions were made or what trade-offs would result from selecting a particular alternative. When asked, participants reported that they simply found it *easier to try an alternative* than to consider whether the provided rationale was valid and relevant—a version of “show, don't tell.” [61]

Participants also reported feeling little reason to *trust* these rationales, which were restricted in how much detail they provided. This lack of detail limited the expected epistemic value of the proclaimed alternative compared with manual testing, but it also made it challenging for participants to understand *on what basis* those rationales were generated, a critical input to participants' assessment of validity.

6.2.4 Attention, Expertise, and the Cost of Awareness. As the contents of the code, chat, and design panes updated “automatically” through agent updates, staying on top of the latest updates to any particular pane required substantial attention. Once lapsed, this “awareness” was costly to regain. The exact cost depended on expertise, all else being equal. For example, senior programmers rapidly lost awareness of the code, and their expertise helped them regain it quickly when needed. The cost of regaining awareness depended

on participants' choices for where to direct attention. If the code felt “hopeless” (P1), or “unfamiliar” (P5), regaining awareness became a priority only when a participant encountered a bug or issue. P5 described the experience of clicking TRY and watching the code change in response, in the following way:

Each time I click on something here [in the design panel] [...] I'm like “Ah! What part of this is important?”

Both domain and programming expertise played a role in mitigating those costs, and thus in how effectively participants could stay on top of changes to code and design. This effect manifested in a few ways: first, domain expertise helped participants more easily recognize the overall “shape” of code components as they came in, making it easier to stay on top of changes with lower self-reported cognitive demands. For example, one participant with data science expertise (P3) could recognize the boilerplate data formatting of sample data as it was generated, but found code for a simulation harder to stay on top of—while another participant with creative coding expertise (P7) found `p5.js` code more straightforward to retain awareness of.

We observed how participants expecting certain code to come from chat or design panel requests watched as that code streamed in from the CONVERSATIONAGENT and then engaged in reflection-in-action [54], expressing surprise (or dismay) when these results deviated from expectations. This behavior echoes the observations by Barke et al. [7] of experienced programmers using GitHub Copilot.

Lastly, some participants (P6, P7, P8) rapidly formed a clear, persistent vision for at least one requested task—and rarely, if ever, found themselves actively seeking design alternatives, feedback, or even an understanding of implicit decisions while in the “flow state” of trying to achieve that vision. In the case of P8, the interviewer switched one task's development context from PAIL to Anthropic's Claude AI system, which could generate and run code that was more aligned with the participant's vision at a more rapid pace than PAIL. This approach was an attempt to understand whether P8 would reach a “saturation point” where they were satisfied that their vision was achieved, and design support might be welcomed. Though a saturation point was reached, the desired design support was explicitly limited to “I'm not really interested in what the system might tell me, the only thing I'd want to do is try it with [my child]” (P8).

7 Discussion

One of our goals in conducting this work is in identifying the next set of challenges the HCI community is likely to face in helping programmers and other technologists design working programs with AI assistance. Based on our experience designing PAIL and the evaluation results, we identify open challenges for LLM-aided program design around (1) how users will define design goals and problems, and evaluate progress towards those goals; and (2) how tools will manage the trade-offs between providing *more complete* information and providing *more relevant* information from the very large set of information that LLMs and other AI systems can inexpensively and rapidly generate. These challenges point to a potential shift in the nature of programming, too, with an increased emphasis

on *interaction design challenges* like managing user attention and perhaps a decreased emphasis on humans relying on particular *programming language capabilities*—the latter increasingly mediated by the increasing use of higher-level natural language.

7.1 Defining Goals, Exploring Problem Spaces

In order for LLMs to effectively support program design, they must identify steps along the pathway towards a user’s ultimate design goal that match a user’s mental model of the problem space. However, as we saw in our evaluation, what steps are meaningful and how important they are depends on where users begin. If users have a point design to begin with, the appropriate next step may be to extrapolate key features before exploring alternatives. In contrast, if users do not know where to begin, a the right next step may be to name a design dimension before exploring concrete instances of it. Indeed, these approaches are characteristic of the double diamond of design as reflected in Figure 2.

But though the double diamond implies a certain linearity, we found that PAIL’s ability to rapidly generate code along many different directions enabled a very nonlinear approach to design. Participants would rapidly shift between exploring possible solutions and recognizing that their problem formulations may not have been addressable. For LLM-aided tools to be effective in helping users with design, they should facilitate *and recognize* this rapid movement between extrapolation and concretization that is enabled by their code generation capabilities.

In one sense, this reflects a shift from exploratory code as *prototype* to exploratory code as *sketch*. In many contexts, generating running code is considered an engineering project within a “solution-space” phase, rather than a “problem-space” phase. But, in its role as *sketch*, running code serves a concrete exploratory purpose in evoking what a solution could look like or work like, just as a hand sketch might.

7.2 Design Practices: Time Compression

While introducing design-related controls in PAIL helped users consider approaches more broadly, PAIL sometimes overwhelmed designers with too much incoming information with a low signal-to-noise ratio. This risks that users will ignore this information, even when structured to support design.

Our observations of this overload speak to one possible cause: a design process altered through compression in time by automation. A typical design process in program design relies on individual humans to write code, test applications with users, consider and enumerate alternatives and design rationales—and these activities as a result are paced at a human timescale. In PAIL, however, these activities are all accelerated, with agents providing information on multiple facets either concurrently or in quick succession. Given this acceleration, the potential for overload should be clear. One question, then, is how and why designers choose to direct their attention among this accelerated information stream, and what role expertise plays in choosing what to consider carefully, what to skim, and what to ignore.

Such an understanding of designer needs and behaviors would enable future systems to be created with an understanding of what feedback is useful, when, and through what mode of delivery. This

is not a new problem: Horvitz identified that a major challenge in automation is the selection of an “ideal action in light of costs, benefits, and uncertainties” [28] as early as 1999.

Ultimately, making sense of the large amount of *potential* information generated by LLM-based agents and other cognitive tools will require a new layer of interaction between human users and the underlying agents producing these insights. The research community is not short on approaches to handling large amounts of data, even when that data changes incrementally over time. But handling large amounts of data that change substantially over time, in contexts where it is hard to assess which of that data is critical to the user, is a challenge that we have only started to consider.

7.3 Managing “More Information” vs. “Better Information”

It’s also not clear who will wield the power that comes with controlling the information layer. Li et al. [39] have argued that tool designers wield a lot of power to shape thought and practices in the domain of creativity support tools. Vaithilingam et al., [64], meanwhile, suggest that LLM-assisted techniques like dynamic grounding can return some of the power to users by allowing tools to *adapt to where humans are*, rather than forcing humans to adapt their ways of thinking and practices to the tools’ capabilities as defined by their designers.

Our PAIL experiences raise the concern that we are on the brink of handing a lot of power to the IDEs, LLMs, and prompt developers building the next generation of tools, because each new LLM-based design affordance demands attention, and showing them all at once will require a major learning curve. Who directs attention, if not the tool? In a typical creativity support tool, the designers of the tools themselves wield direct control over tool behavior, but in LLM-powered systems, designers often cede varying degrees of control to black-box models they did not even have a hand in training. We may be moving from a formal system of rules and practices that are at least discoverable and interpretable, to an analog world of prompt-driven components whose very behavior is both inherently unpredictable and also dependent on the model it happens to be executed against. Our experiences with attention overload in other realms (e.g., social media [50]) suggest that this future may be less empowering for users and developers, rather than more.

A key design decision that will shape how systems wield this power goes beyond what information they surface to users and how, and to what information they *generate* and then choose to surface. While a major area of design has been and will continue to be principles for managing information better, our work with PAIL also asks if there *better information* to be managing in the first place.

8 Limitations and Future Work

PAIL represents one possible point in the design space of AI-supported program design. We did not explicitly compare PAIL with ChatGPT or Claude in our evaluation because it was not our goal to show that PAIL is more effective than those baselines, but rather to identify the challenges that arise when explicit design support is integrated into an IDE like PAIL. We hope our work here opens up a design

space with some initial, critical insights about what the community should tackle next.

That said, there are two key limitations of this work that future work should address: the generalizability of PAIL and transitions across levels of abstraction when programming with LLMs.

First, PAIL isn't intended for every programming task, nor for every programmer. Though PAIL helped get participants started regardless of design background, *some* level of design process literacy is required to make effective use of PAIL. Design process literacy can be taught, of course—or it can be designed into tools like PAIL. For example, recall that participants were most willing to engage with CONVERSATIONAGENT's questions about the end-user and their needs when the agent was *not* simultaneously producing changes to the program that a user could be testing—providing code in this context was actively counterproductive. This observation points to a compelling, but also concerning, mechanism by which a tool like PAIL could encourage more consideration of various parts of the design process: hold the project hostage by simply refusing to produce any code until the user has considered what the tool wants the user to consider.

Second, regardless of approach, designers need to operate across abstraction levels, knowing when to pop up to a high level from the weeds, and when to deep dive towards a point solution in order to better understand a particular neighborhood within the design space. PAIL currently supports a few levels within its hierarchy of abstraction: code, a layer of established program requirements and decisions embedded in that code, and a set of higher-level design-related concerns. As recently noted by Vaithilingam et al. [64], design decisions happen in a fractal pattern, with many decisions across many levels of abstraction. Future LLM-aided programming tools could provide finer-grained control over operation at different levels of abstraction to provide greater control over and understanding of generated code, as well as more targeted legibility for the immediate subtask at hand. We contribute here a deeper understanding that there are substantial challenges designers will face when transitioning across abstraction levels.

Supporting users identifying design problems and solutions, as they handle an information overload and transition across abstraction levels, requires that we tool designers reckon with several major open questions:

How much agency or initiative should automated systems be given to set direction? It's one thing to synthesize some code for a user to evaluate—it's another entirely to control when a user considers their high-level goals for a project instead of staying lost in the weeds, or to scope out a set of potentially-anchoring alternative points in the design space. And if programmers spend less time writing code and more time providing higher-level instructions, then managing *attention* in the face of *too much data* and *unknown signal-to-noise ratios* among that data will become critical. How can systems correctly decide what data to show users, and when?

9 Conclusion

Through this work, we explored the implications of LLM-aided program design, focused on support for problem formulation and assessing solution suitability. PAIL, our design probe, encourages

developers to follow a user-centered design process and tracks requirements discovered and decisions made through prototyping. Through our user study, we found some evidence that this kind of assistance can be helpful in broadly considering the program design space, but also uncover a set of challenges around managing attention and maintaining awareness of program updates, pointing to broader questions and trade-offs across generating and sharing information, ensuring information is relevant to users, and balancing agency between users and their new generation of tools.

10 Disclosure

The authors used ChatGPT for minor copyediting tasks.

Acknowledgments

The authors would like to thank the many individuals who supported this work, including: Chetan Goenka and Keira Swei, undergraduate researchers supported in their contributions to this work through Berkeley's Engineering Design Scholars summer program; Timothy J. Aveni, Shm Garanganao Almeda, James Smith, and Shreya Shankar, for many fruitful discussions on the concepts underlying this work, as well as feedback on drafts of this paper; our study participants, who graciously gave their time and feedback to this effort; and our anonymous reviewers, whose thoughtful comments shaped many revisions here.

J.D. Zamfirescu-Pereira was partially supported by a Google PhD Fellowship. Qian Yang's effort was supported in part by the AI2050 Early Career Fellowship program at Schmidt Sciences. This material is also based upon work supported by the National Science Foundation under Grant No. 2313078. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

References

- [1] 2024. Artifacts are now generally available. <https://www.anthropic.com/news/artifacts>
- [2] 2024. What is User Centered Design (UCD)? <https://www.interaction-design.org/literature/topics/user-centered-design>
- [3] Shm Garanganao Almeda, J.D. Zamfirescu-Pereira, Kyu Won Kim, Pradeep Mani Rathnam, and Bjoern Hartmann. 2024. Prompting for Discovery: Flexible Sense-Making for AI Art-Making with Dreamsheets. In *Proceedings of the CHI Conference on Human Factors in Computing Systems (CHI '24)*. Association for Computing Machinery, New York, NY, USA, 1–17. <https://doi.org/10.1145/3613904.3642858>
- [4] Barrett R Anderson, Jash Hemant Shah, and Max Kreminski. 2024. Homogenization Effects of Large Language Models on Human Creative Ideation. In *Creativity and Cognition*. ACM, Chicago IL USA, 413–425. <https://doi.org/10.1145/3635636.3656204>
- [5] Tyler Angert, Miroslav Suzara, Jenny Han, Christopher Pondoc, and Hariharan Subramonyam. 2023. Spellburst: A Node-based Interface for Exploratory Creative Coding with Natural Language Prompts. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (UIST '23)*. Association for Computing Machinery, New York, NY, USA, 1–22. <https://doi.org/10.1145/3586183.3606719>
- [6] Ian Arawjo, Priyan Vaithilingam, Martin Wattenberg, and Elena Glassman. 2023. ChainForge: An open-source visual programming environment for prompt engineering. 1–3.
- [7] Shraddha Barke, Michael B James, and Nadia Polikarpova. 2023. Grounded copilot: How programmers interact with code-generating models. *Proceedings of the ACM on Programming Languages* 7, OOPSLA1 (2023), 85–111. Publisher: ACM New York, NY, USA.
- [8] Shaon Barman, Sarah Chasins, Rastislav Bodik, and Sumit Gulwani. 2016. Ringer: web automation by demonstration. 748–764.

- [9] Mary Beth Kery and Brad A. Myers. 2017. Exploring exploratory programming. In *2017 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)*. 25–29. <https://doi.org/10.1109/VLHCC.2017.8103446> ISSN: 1943-6106.
- [10] Mary Jo Bitner, Amy L Ostrom, and Felicia N Morgan. 2008. Service blueprinting: a practical technique for service innovation. *California management review* 50, 3 (2008), 66–94. Publisher: SAGE Publications Sage CA: Los Angeles, CA.
- [11] 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 (UIST '23)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3586183.3606725>
- [12] J. Brandt, P.J. Guo, J. Lewenstein, S.R. Klemmer, and M. Dontcheva. 2009. Writing Code to Prototype, Ideate, and Discover. *IEEE Software* 26, 5 (Sept. 2009), 18–24. <https://doi.org/10.1109/MS.2009.147>
- [13] Joel Brandt, Philip J. Guo, Joel Lewenstein, Mira Dontcheva, and Scott R. Klemmer. 2009. Two studies of opportunistic programming: interleaving web foraging, learning, and writing code. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, Boston MA USA, 1589–1598. <https://doi.org/10.1145/1518701.1518944>
- [14] Bill Buxton. 2010. *Sketching user experiences: getting the design right and the right design*. Morgan kaufmann.
- [15] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukas Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating Large Language Models Trained on Code. <https://doi.org/10.48550/arXiv.2107.03374> arXiv:2107.03374 [cs] version: 2.
- [16] Yan Chen, Steve Oney, and Walter S. Lasecki. 2016. Towards Providing On-Demand Expert Support for Software Developers. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems (CHI '16)*. Association for Computing Machinery, New York, NY, USA, 3192–3203. <https://doi.org/10.1145/2858036.2858512>
- [17] Design Council. 2004. Framework for Innovation. <https://www.designcouncil.org.uk/our-resources/framework-for-innovation/>
- [18] Mihaly Csikszentmihalyi. 1990. *Flow: The psychology of optimal experience*. New York: Harper & Row.
- [19] Paul Denny, Viraj Kumar, and Nasser Giacaman. 2023. Conversing with copilot: Exploring prompt engineering for solving cs1 problems using natural language. 1136–1142.
- [20] Giulia Di Fede, Davide Rocchesso, Steven P Dow, and Salvatore Andolina. 2022. The Idea Machine: LLM-based Expansion, Rewriting, Combination, and Suggestion of Ideas. In *Creativity and Cognition*. 623–627.
- [21] Steven P. Dow, Alana Glasco, Jonathan Kass, Melissa Schwarz, Daniel L. Schwartz, and Scott R. Klemmer. 2010. Parallel prototyping leads to better design results, more divergence, and increased self-efficacy. *ACM Transactions on Computer-Human Interaction* 17, 4 (Dec. 2010), 1–24. <https://doi.org/10.1145/1879831.1879836>
- [22] Laura Faulkner. 2003. Beyond the five-user assumption: Benefits of increased sample sizes in usability testing. *Behavior Research Methods, Instruments, & Computers* 35 (2003), 379–383. Publisher: Springer.
- [23] GitHub. 2021. GitHub Copilot • Your AI Pair Programmer. <http://github.com/features/copilot>
- [24] Ken Gu, Madeleine Grunde-McLaughlin, Andrew McNutt, Jeffrey Heer, and Tim Althoff. 2024. How Do Data Analysts Respond to AI Assistance? A Wizard-of-Oz Study. In *Proceedings of the CHI Conference on Human Factors in Computing Systems (CHI '24)*. Association for Computing Machinery, New York, NY, USA, 1–22. <https://doi.org/10.1145/3613904.3641891>
- [25] Sumit Gulwani. 2011. Automating string processing in spreadsheets using input-output examples. *ACM Sigplan Notices* 46, 1 (2011), 317–330. Publisher: ACM New York, NY, USA.
- [26] Björn Hartmann, Loren Yu, Abel Allison, Yeonsoo Yang, and Scott R. Klemmer. 2008. Design as Exploration: Creating Interface Alternatives through Parallel Authoring and Runtime Tuning. In *Proceedings of the 21st Annual ACM Symposium on User Interface Software and Technology (UIST '08)*. Association for Computing Machinery, New York, NY, USA, 91–100. <https://doi.org/10.1145/1449715.1449732> event-place: Monterey, CA, USA.
- [27] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiaowu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. 2023. MetaGPT: Meta Programming for A Multi-Agent Collaborative Framework. <https://openreview.net/forum?id=VtmBAGCN7o>
- [28] Eric Horvitz. 1999. Principles of mixed-initiative user interfaces. In *Proceedings of the SIGCHI conference on Human factors in computing systems the CHI is the limit - CHI '99*. ACM Press, Pittsburgh, Pennsylvania, United States, 159–166. <https://doi.org/10.1145/302979.303030>
- [29] Stephanie Houde and Charles Hill. 1997. What do Prototypes Prototype? In *Handbook of Human-Computer Interaction (Second Edition)*, Marting G. Helander, Thomas K. Landauer, and Prasad V. Prabhu (Eds.). North-Holland, Amsterdam, 367–381. <https://doi.org/10.1016/B978-044481862-1.50082-0>
- [30] Dhanya Jayagopal, Justin Lubin, and Sarah E Chasins. 2022. Exploring the learnability of program synthesizers by novice programmers. 1–15.
- [31] Majeed Kazemitabaar, Jack Williams, Ian Drosos, Tovi Grossman, Austin Henley, Carina Negreanu, and Advait Sarkar. 2024. Improving Steering and Verification in AI-Assisted Data Analysis with Interactive Task Decomposition. <https://doi.org/10.1145/3654777.3676345> arXiv:2407.02651 [cs].
- [32] Mary Beth Kery. 2017. Designing Effective History Support for Exploratory Programming Data Work. (2017).
- [33] Mary Beth Kery, Amber Horvath, and Brad Myers. 2017. Variolite: Supporting Exploratory Programming by Data Scientists. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems (CHI '17)*. Association for Computing Machinery, New York, NY, USA, 1265–1276. <https://doi.org/10.1145/3025453.3025626>
- [34] Mary Beth Kery, Bonnie E. John, Patrick O'Flaherty, Amber Horvath, and Brad A. Myers. 2019. Towards Effective Foraging by Data Scientists to Find Past Analysis Choices. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, Glasgow Scotland Uk, 1–13. <https://doi.org/10.1145/3290605.3300322>
- [35] Pranav Khadpe, Ranjay Krishna, Li Fei-Fei, Jeffrey T Hancock, and Michael S Bernstein. 2020. Conceptual metaphors impact perceptions of human-ai collaboration. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW2 (2020), 1–26. Publisher: ACM New York, NY, USA.
- [36] Michelle S. Lam, Janice Teoh, James A. Landay, Jeffrey Heer, and Michael S. Bernstein. 2021. Concept Induction: Analyzing Unstructured Text with High-Level Concepts Using LLoM. In *Proceedings of the CHI Conference on Human Factors in Computing Systems (CHI '24)*. Association for Computing Machinery, New York, NY, USA, 1–28. <https://doi.org/10.1145/3613904.3642830>
- [37] James A. Landay and Brad A. Myers. 1995. Interactive sketching for the early stages of user interface design. In *Proceedings of the SIGCHI conference on Human factors in computing systems - CHI '95*. ACM Press, Denver, Colorado, United States, 43–50. <https://doi.org/10.1145/223904.223910>
- [38] Tomas Lawton, Kazjon Grace, and Francisco J Ibarrola. 2023. When is a Tool a Tool? User Perceptions of System Agency in Human-AI Co-Creative Drawing. In *Proceedings of the 2023 ACM Designing Interactive Systems Conference*. 1978–1996.
- [39] Jingyi Li, Eric Rawn, Jacob Ritchie, Jasper Tran O'Leary, and Sean Follmer. 2023. Beyond the Artifact: Power as a Lens for Creativity Support Tools. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (UIST '23)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3586183.3606831>
- [40] Toby Jia-Jun Li, Amos Azaria, and Brad A Myers. 2017. SUGLITE: creating multimodal smartphone automation by demonstration. In *Proceedings of the 2017 CHI conference on human factors in computing systems*. 6038–6049.
- [41] Geoffrey Litt. 2024. Patchwork lab notebook: Version control for everything. <https://www.inkandswitch.com/patchwork/notebook/>
- [42] 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. 1–31.
- [43] Vivian Liu and Lydia B Chilton. 2022. Design Guidelines for Prompt Engineering Text-to-Image Generative Models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems (CHI '22)*. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3491102.3501825>
- [44] Justin Lubin, Nick Collins, Cyrus Omar, and Ravi Chugh. 2019. Program Synthesis with Live Bidirectional Evaluation. *arXiv preprint arXiv:1911.00583* (2019).
- [45] Aran Lunzer and Kasper Hornbæk. 2008. Subjunctive Interfaces: Extending Applications to Support Parallel Setup, Viewing and Control of Alternative Scenarios. *ACM Trans. Comput.-Hum. Interact.* 14, 4 (Jan. 2008). <https://doi.org/10.1145/1314683.1314685> Place: New York, NY, USA Publisher: Association for Computing Machinery.
- [46] J. Marks, B. Andalman, P. A. Beardsley, W. Freeman, S. Gibson, J. Hodgins, T. Kang, B. Mirtich, H. Pfister, W. Ruml, K. Ryall, J. Seims, and S. Shieber. 1997. Design Galleries: A General Approach to Setting Parameters for Computer Graphics and Animation. In *Proceedings of the 24th Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH '97)*. ACM Press/Addison-Wesley Publishing Co., USA, 389–400. <https://doi.org/10.1145/258734.258887>
- [47] Andrew M. McNutt, Chenglong Wang, Robert A. DeLine, and Steven M. Drucker. 2023. On the Design of AI-powered Code Assistants for Notebooks. <http://arxiv.org/abs/2301.11178> arXiv:2301.11178 [cs].
- [48] Bill Moggridge. 2006. *Designing Interactions*. The MIT Press.

- [49] Jakob Nielsen. [n. d.]. 10 Usability Heuristics for User Interface Design. <https://www.nngroup.com/articles/ten-usability-heuristics/>
- [50] Amy Orben, Andrew K. Przybylski, Sarah-Jayne Blakemore, and Rogier A. Kievit. 2022. Windows of developmental sensitivity to social media. *Nature Communications* 13, 1 (March 2022), 1649. <https://doi.org/10.1038/s41467-022-29296-3> Publisher: Nature Publishing Group.
- [51] James Prather, Brent N Reeves, Paul Denny, Brett A Becker, Juho Leinonen, Andrew Luxton-Reilly, Garrett Powell, James Finnie-Ansley, and Eddie Antonio Santos. 2023. "It's Weird That it Knows What I Want": Usability and Interactions with Copilot for Novice Programmers. *ACM Transactions on Computer-Human Interaction* 31, 1 (2023), 1–31. Publisher: ACM New York, NY.
- [52] Eric Rawn, Jingyi Li, Eric Paulos, and Sarah Chasins. 2023. Understanding Version Control as Material Interaction with Quickpose. In *Proceedings of the SIGCHI conference on human factors in computing systems*.
- [53] Jeff Sauro and James R Lewis. 2016. *Quantifying the user experience: Practical statistics for user research*. Morgan Kaufmann.
- [54] Donald A Schon. 1984. *The Reflective Practitioner: How Professionals Think In Action*. Vol. 5126. Basic Books.
- [55] Shreya Shankar, J.D. Zamfirescu-Pereira, Bjoern Hartmann, Aditya Parameswaran, and Ian Arawjo. 2024. Who Validates the Validators? Aligning LLM-Assisted Evaluation of LLM Outputs with Human Preferences. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*. ACM, Pittsburgh PA USA, 1–14. <https://doi.org/10.1145/3654777.3676450>
- [56] Beau Sheil. 1986. DATAMATION®: POWER TOOLS FOR PROGRAMMERS. In *Readings in Artificial Intelligence and Software Engineering*, Charles Rich and Richard C. Waters (Eds.). Morgan Kaufmann, 573–580. <https://doi.org/10.1016/B978-0-934613-12-5.50048-3>
- [57] Chenglei Si, Yanzhe Zhang, Zhengyuan Yang, Ruibo Liu, and Diyi Yang. 2024. Design2Code: How Far Are We From Automating Front-End Engineering? <http://arxiv.org/abs/2403.03163> arXiv:2403.03163 [cs].
- [58] 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. 1–19.
- [59] Sangho Suh, Bryan Min, Srishti Palani, and Haijun Xia. 2023. Sensecape: Enabling Multilevel Exploration and Sensemaking with Large Language Models. *arXiv preprint arXiv:2305.11483* (2023).
- [60] Theodore R. Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L. Griffiths. 2024. Cognitive Architectures for Language Agents. <https://doi.org/10.48550/arXiv.2309.02427> arXiv:2309.02427 [cs].
- [61] M. Swan. 1927. *How You Can Write Plays: A Practical Guide-book*. S. French. <https://books.google.com/books?id=UyYrAAAAIAAJ>
- [62] 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. <http://arxiv.org/abs/2311.00710> arXiv:2311.00710 [cs].
- [63] Michael Terry, Elizabeth D. Mynatt, Kumiyo Nakakoji, and Yasuhiro Yamamoto. 2004. Variation in element and action: supporting simultaneous development of alternative solutions. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, Vienna Austria, 711–718. <https://doi.org/10.1145/985692.985782>
- [64] Priyan Vaithilingam, Ian Arawjo, and Elena L. Glassman. 2024. Imagining a Future of Designing with AI: Dynamic Grounding, Constructive Negotiation, and Sustainable Motivation. In *Designing Interactive Systems Conference*. ACM, IT University of Copenhagen Denmark, 289–300. <https://doi.org/10.1145/3643834.3661525>
- [65] Tongshuang Wu, Michael Terry, and Carrie J 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*.
- [66] 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*. 2206–2220.
- [67] 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*. 1–21.
- [68] John Zimmerman, Jodi Forlizzi, and Shelley Evenson. 2007. Research through design as a method for interaction design research in HCI. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '07)*. Association for Computing Machinery, New York, NY, USA, 493–502. <https://doi.org/10.1145/1240624.1240704>