# Conversational AI Threads for Visualizing Multidimensional Datasets

MATT-HEUN HONG, University of North Carolina, USA
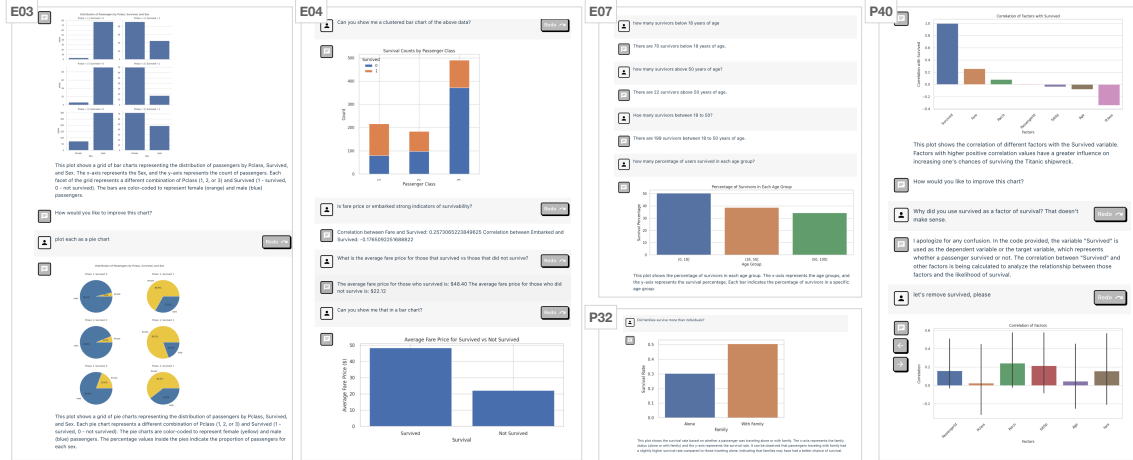
ANAMARIA CRISAN, Tableau Research, USA

Fig. 1. **Excerpts of participants' analytic conversations with AI Threads.** We show participant conversations in threads modifying the visual encoding (E03, E04, P40), as well as for sequential (E07) and individual (P32) utterances. Excerpts were derived from participants in crowd-sourced (P#) and interview (E#) studies that we conducted. AI Threads is capable of editing operations that add or remove attributes to encodings and change encoding types, among others. Overall, the chatbot proved itself capable, but we still encountered instances of incorrect responses, including within a continuous thread (P40).

Generative Large Language Models (LLMs) show potential in data analysis, yet their full capabilities remain uncharted. Our work explores the capabilities of LLMs for creating and refining visualizations via conversational interfaces. We used an LLM to conduct a re-analysis of a prior Wizard-of-Oz study examining the use of chatbots for conducting visual analysis. We surfaced the strengths and weaknesses of LLM-driven analytic chatbots, finding that they fell short in supporting progressive visualization refinements. From these findings, we developed AI Threads, a multi-threaded analytic chatbot that enables analysts to proactively manage conversational context and improve the efficacy of its outputs. We evaluate its usability through a crowdsourced study ($n = 40$) and in-depth interviews with expert analysts ($n = 10$). We further demonstrate the capabilities of AI Threads on a dataset outside the LLM's training corpus. Our findings show the potential of LLMs while also surfacing challenges and fruitful avenues for future research.

**Code & Supplemental Materials:** https://osf.io/6wxpa/

CCS Concepts: • **Human-centered computing** → **Empirical studies in visualization**; **Natural language interfaces**.

Authors' addresses: Matt-Heun Hong, University of North Carolina, USA, mhh@cs.unc.edu; Anamaria Crisan, Tableau Research, USA, acrisan@tableau.com.

## 1   INTRODUCTION

Generative Pre-trained Transformer (GPT) large language models (LLMs) [48] can assist in data-driven knowledge work by generating analytic code from natural language utterances [2, 14, 47]. These capabilities can support novice and expert analysts alike by reducing barriers to conducting preliminary investigations of data [8, 10, 40]. However, LLMs do not reflect the first attempts to facilitate analyses through natural language. Prior research on Visual Natural Language Interfaces (V-NLIs) and conversational agents have explored varying techniques for translating plain language utterances into corresponding data transformations and visualization specifications [22, 41, 57, 62, 66]. A recent head-to-head comparison of LLMs and several V-NLIs suggests that LLM performance is as capable, and potentially superior, to prior systems for resolving ambiguities around user intent behind an utterance [40]. However, this investigation was shallow and did not probe the potential for LLMs to operate as part of conversational interfaces for more dynamic analytic dialogues. As such, the capabilities of LLMs for conducting data analyses remain uncharted.

Analytic conversations have a complex structure that includes shifting conversation context [64] and sequential utterances that progressively modify analyses and visual encodings [13, 62]. Here, we explore how LLMs handle this shifting conversational context. Our research has three parts. First, we re-examine the results of a prior Wizard-of-Oz [64] study and summarize key findings on how an LLM-based analytic agent responds. A key feature of [64] is that they explicitly mark the shifts in conversational context, which we use as a ground truth in our assessments. Next, we developed AI Threads, a conversational interface that uses an LLM as an AI agent to produce both visual and textual outputs in response to user utterances. A novel aspect of AI Threads is the incorporation of multiple conversational threads, a common feature of messaging apps, to manage the context provided to the AI agent. In the third part, we use the data and objectives from [64] to conduct two usability studies that assessed the validity, robustness, and utility of AI Threads. The first is a crowd-sourced study with 40 participants, followed by a second in-depth study that conducted semi-structured interviews with 10 analytic experts. The findings from these user studies show that participants had an overall positive perception of AI Threads and that they found it to be quite capable. However, full reviews of the session logs also revealed its pitfalls and limitations. Expert participants provided further insight into the use of AI chatbots relative to their existing analytic processes.

Collectively, our research presents the following three contributions. **First, we discuss the structure of multi-threaded visual analytic conversations and describe their significance for LLM AI Agents.** We demonstrate how this contribution can help to proactively manage conversational context and improve the organization of the analysis as it progresses. **Second, we present the design and implementation of AI Threads**. We ground its development of in four design goals that aim to enable natural conversational flows for data analysis. **Finally, we present a set of user studies that showcase both the capabilities and limitations of LLM as analytic conversational agents.** We also identify how such conversational agents stand to influence data analysis.

As conversational agents supported by LLMs are poised to play a larger role in data analysis, including the generation of visualizations, our research presents a novel approach to supporting these complex and multi-faceted conversations.

## 2 RELATED WORK

We summarize related work pertinent to Visual Natural Language Interfaces, Conversational agents, and visualization recommendation systems.

### 2.1 Visual Natural Language Interfaces (V-NLIs)

V-NLIs have undergone an evolution of techniques to understand utterances to create appropriate mappings to dataset attributes, data transformation, and visual encodings. [40, 57]

Initial approaches, such as DataTone [16], Eviza [55], and FlowSense [71] used parse trees, grammar, rules, or heuristics to resolve these mappings. These prior works also uncovered the challenges of managing ambiguous utterances and capturing and modeling the semantics, intent, and pragmatics of these conversations [23]. The Orko system [63] reified these channels with a non-tabular dataset as well. While these techniques are still used in contemporary V-NLI systems [57], they are also critiques as being too brittle, which limits their efficacy [40].

More recent approaches have tried to learn these mappings directly from the data using machine learning approaches, including Deep Neural Networks. Examples for such systems include, Advisor [33], ncNet [38], Talk2Data [20], and Urania [19]. DNN approaches vary with regards to the specific architecture and labeled training corpora they use to learn mapping functions. More recently, the use of LLMs has been explored for V-NLIs by leveraging their program generating capabilities [2] to generate and execute code for transforming and visualizing data. Unlike DNN approaches, LLMs may only require prompt engineering and a few labeled examples to generate visualizations from natural language queries. As in the current work, Chat2Vis [40] uses GPT-3.5 to generate visualizations from natural language queries of a dataset. Their preliminary evidence suggests LLMs outperform other approaches for understanding a natural language utterance and producing an appropriate visualization. This transition from many early attempts at V-NLIs to the use of LLMs follows an arc from heuristic (e.g., [16, 55]) to data-driven approaches that many fields adjacent to NLP have similarly experienced.

We extend the prior research by conducting a deeper investigation into the use of LLMs to support the visual analysis of data. Compared to Maddigan *et al.* we conduct a more comprehensive assessment of a LLM's capabilities with data from actual human subjects, not just scenarios. Moreover, the AI Threads system we develop addresses many of the limitations that they identify.

### 2.2 Conversational Agents for Visual Data Analysis

Conversational agents, or chatbots, facilitate an interactive dialogue between the analyst and an automated agent [41]. V-NLI research has explored the potential of conversational interfaces powered by chatbots to conduct a data analysis [25, 56]. Iris [13] uses a single-threaded conversational interface that combination of multiple utterances to produce an appropriate response; notably, Iris produces both textual and visual responses. Kim *et. al.* [28] implement a pipeline that leverages natural language templates to facilitate a conversation for explaining visualizations. Both [13] and [28] touch upon the importance of composition operations (i.e. requiring multiple operations) to carry out these analytic conversations. Finally, to support participants in visual analysis, SNOWY [61] produces ranked suggestions for utterances that users may pose of the data and visualization.

Adjacent to these efforts of system development, prior research has also explored end-users expectations for analytic chatbots and the structure of analytic conversations. Hearst and Tory [21] show when conversing with analytic chatbots participants sometimes show a greater preference for text responses over visualizations. Setlur and Tory [64] use a Wizard-of-Oz study to study the relationship between participant utterances and the automatic generation of

visualizations, including the refinement of visualizations. They summarize their findings as a state transition model and discuss the importance of context and intent to facilitate these transitions and support both conversational and visual coherence. In follow-up research, they also demonstrate the importance of the input medium and modalities (voice or text) in influencing these conversations [56]. They found that text interactions in combination with conversational threading supported a deeper analysis compared to other modalities.

The utility of LLMs as analytic chatbots has not been extensively explored. Moreover, recent research has shown that LLMs can struggle with compositionality [12]. We present a detailed examination of LLMs for conducting analytic conversations that include compositional tasks.

## 2.3   Visualization Recommendation

Visualization recommendation (VisRec) techniques attempt to automatically generate visualizations primarily based on the attributes of the data. Unlike V-NLIs they are not tied specifically to natural language inputs, but have overlapping techniques to generate visualizations from user inputs.

Several techniques use both rule-based and heuristic approaches to general visualizations. ShowMe [39] uses a rule-based approach predicated on the field types (categorical or numeric) and best practices ranking. It requires the use of specifying one or more fields before it can automatically generate an encoding. VizDeck [27] similarly used data properties to recommend interactive dashboards, instead of single encodings. ShowMe was further extended in the Voyager and Voyager2 systems [69], which use statistical properties of the dataset to automatically generate and rank visualizations of a dataset while relaxing the requirement of an initial user specification. Draco [43] presented an approach to automatically incorporate perceptual efficacy experiments into these visualization recommendation paradigms. ChartSeer [72], Chart2Vec [7], and DeepEye [37] propose machine learning approaches as alternatives to the rule-based methods. Both ChartSeer [72] and Chart2Vec [7] learning a vector representation of visual encodings, albeit in different ways, to recommend visualizations, while DeepEye [37] uses a learning-to-rank approach. Recently, LIDA was developed as a way to recommend visualizations by using an LLM to generate specifications [10].

The previously mentioned approaches focus largely on specifications of the data attributes and mapping these to visual encodings. Other techniques can take additional contextual factors such as the analytic task, user intent, prior visualizations, and even the domain application. TaskVis [58] uses a pre-defined set of tasks and a rule-based approach to recommend visualization from a combination of data fields and analytic tasks. Frontier [31] similarly uses pre-defined task and intent classifications for data analysis to generate and organize visualization recommendations. The concepts of Frontier are reified in the Lux [32], a Jupyter notebook widget that interacts with Pandas, a widely used Python data analysis library. Medely extends these intent-based recommendations further to generate dashboard visualizations [45]. GraphScape [29] can propose modifications to existing encodings using a graph transition model. Finally, visualization recommendations can also take domain context into account. For example, GenoRec [44] and GEViTRec [9] both leverage data from the genomics domain to generate recommendations.

Our research combines many of the ideas from several disparate visual recommendation efforts, notably, using LLM (as LIDA does), visual sequencing to modify encodings, and the incorporation of contextual factors.

## 3   RE-EVALUATING ANALYTIC CONVERSATIONS WITH AI CHATBOTS

We conducted a re-analysis of study by Setlur and Tory [64] that examined the dialogue of visual analytic conversations. Our goal is not to reproduce the precise visualizations from their Wizard-of-Oz (WoZ) study, but instead, to evaluate *if*

and *how* a LLM is able to handle transitions in conversations from modify an existing visualization or generate a new one. We summarize the key findings here and provide a deeper analysis in Appendix A at the end of this manuscript.

### 3.1 Overview of Tory and Setlur (2019)

Tory and Setlur [64] conducted a WoZ study to examine analytic conversations with an AI chatbot. Participants were tasked to explore the canonical Titanic dataset and answer the question *"What characteristics made it more likely that a passenger survived?"* Participants conducted their analysis by entering utterances (textual inputs) via a research prototype of Tableau's Ask Data feature. The Wizard generated a visualization in response and these were displayed to the participant. The study had four conditions. Two conditions treat each utterance as independent and produce a unique visualization for each, and two that take into account the previous visualization – these were called the context conditions. In the two context conditions, participants could elect to modify an existing visualization or begin a new analysis. The retention or shift in context was either explicitly stated (e.g., uttering "retain", or "start over") or implicitly determined (e.g., uttering "add to this", "remove this"). The Wizard would respond to such utterances according to the prescribed rules of the study.

We sought to *"replay"* the conversations from the context conditions using an LLM. The retention or shift in context observed in Tory and Setlur [64] delineate discourse chunks that split the analytic conversation between utterances signifying new analysis questions or goals and those intended to modify an existing visualization (see Appendix A.1.1 and Appendix A.1.3). We use these contextual transitions as ground truth to evaluate an LLM. To further examine the LLM's performance we also annotated utterances from this according to three phrasal classifications (interrogative, imperative, declarative) [30, 53] and assessed the ambiguity of these utterances according to the criteria proposed by Srinivasan *et al.* [60] (see Figure A.1.2 and Appendix A.1.2).

### 3.2 Replication with a ADA

To conduct this evaluation, we used OpenAI's ChatGPT Advanced Data Analysis (ADA) plugin[1]. ADA is an extension of ChatGPT for code generation in Python and is capable of conducting analytic conversations[2]. Managing the conversation context for an LLM is a nascent area of development. Currently, a tool like ADA takes an entire conversation as context – up to a set token limit. Once the token limit is reached, there are different strategies for managing context including summarizing the previous conversation, performing an embedding search to identify relevant parts of the conversation, or actively filtering prior utterances. We use ADA's default context management. Some or all of the aforementioned strategies and likely others may already be incorporated into ADA, but this information is not publicly available at present.

*3.2.1 Findings.* A summary of the findings between the original WoZ study and ADA is summarized in Figure 2. The default output for ADA is a textual, not visual, response but this behavior can be modified with the appropriate prompts (see Appendix A.2). Overall, we found that ADA was able to correctly reproduce the initial data visualization regardless of the attribute, data, or visual ambiguity of the utterance. For example, the interrogative utterance *"What percentage survived by age in 5-year bins?"* specifies the target attributes explicitly (survived and age), the data transformations (five-year bins and aggregation to a percentage), and leaves it to ADA to choose the visual encoding. ADA produces the appropriate bar chart. We observed that it can also respond effectively to non-explicit attribute references (e.g.,

---

[1]We used the version of ADA released before July 15, 2023
[2]ADA was initially released as the ChatGPT Code Interpreter Plugin https://openai.com/blog/chatgpt-plugins#code-interpreter

Fig. 2. **Comparison for visualization produced between the WoZ study and different systems.** Progressive utterances for modifying an initial visual encoding (top row) by different systems. Black boxes around the visualization indicate the creation of a new visualization—a failure to apply a progressive update when given an utterance with the intent to modify a previous visualization.

*"show only male"* or *"how much did each passenger pay to board the Titanic"*). Where we observed ADA struggle was transition context without an explicit reference to the previous visualization.

We observed that many of the modifying utterances in [64] were terse and either imperative or declarative (e.g., *"class color by fare bin"*, *"add trend line"*, or *"remove class"*; see also Appendix A.1.3). When an utterance had explicit anaphoric references, such as *this* or *add*, ADA was typically able to resolve the contextual dependency and modify the previous visualization. However, when there was no anaphoric reference, such as in the example *"class color by fare bin"*, ADA interpreted this utterance as independent and generated a new visualization even though the participant intended to modify the previous visualization. However, ADA's behavior was also inconsistent. After several anaphoric references, ADA appeared to learn that the participant intended to modify the previous visualization and interpreted subsequent utterances as intending to modify the previous visualization. This interpretation was not always correct and in practice could frustrate an end-user, as Tory and Setlur also observed [64]. Using utterances to start a new visualization (e.g., *"start over"*) also had unpredictable side effects because it was not clear to ADA at which point it should restart the

analysis. For example, its default response would be to reload the entire dataset instead of just beginning a new visualization. It was necessary to be explicit to produce an appropriate output (e.g., *"I am done modifying the previous visualization, let's start a new one"*).

Overall, we often needed to substantively modify the utterances from [64] to make them less ambiguous by adding anaphoric references to produce an appropriate response. This was a trial-and-error process and did not always produce the desired results. Additionally, the conversational thread became difficult to follow because it contained failed attempts to modify the utterance. Anaphoric references, while useful, are a manual and *ad hoc* approach for an end-user to manage context for the chatbot. While they are workable, they significantly impact the overall coherence of the conversation.

### 3.3 Implications: AI agents and analytic conversations

**Analysts do not precisely state their analytic goals at the outset, iterative refinement is the norm.** Our re-evaluation of Tory and Setlur [64] identifies several strengths and limitations of LLMs serving as conversational analytic agents. While we observed that ADA appeared to perform well at composing a singular visualization, it struggled to effectively manage a transitioning conversational context that is common in analytic dialogue. Conversational agents are developed for many purposes [41], the structure and nuances of analytic conversations merit their own set of unique approaches.

### 4 CONVERSATIONAL AI AGENT DESIGN GOALS

We draw on the findings from our re-evaluation of Tory and Setlur [64], the structure of analytic conversations (Section 2.2), V-NLI and conversational interfaces (Section 2.1), and finally research on analytic notebooks [26, 54] to derive a set of design criteria for an analytic conversational agent.

**DG1: Support Conversational Coherence.** The iterative nature of data analysis can result in a complex dialogue. Commonly this analytic dialogue is found not in conversational interfaces, but typically organized and presented within computational notebooks [26, 34, 54]. These notebooks interplay visualization, code, and text annotations that analysts structure over time into a more coherent analysis. In aggregate, this refinement process results in a branching structure [26, 50, 54] that resembles the conversational evolution and contextual shifts we observed in [64] (see Appendix A.1 and A.1.3). Transitioning a linear conversational flow, such as is used in ADA, to a multi-threaded one that resembles this branching structure will improve conversational coherence and consequently improve user experience. Prior research with chatbots [56] also shows that participants prefer this non-linear structure as well.

**DG2: Enable Naturalistic Conversation.** Prior research for V-NLIs and conversational agents specifies several criteria for effectively supporting naturalistic conversations. These include being able to resolve utterance ambiguities and incompleteness, sequential utterances, in addition to understanding the user's analytic intent (see Section 2.1). In response to end-user utterances, the chatbot should seamlessly support both textual and visual as appropriate [13, 21]. Finally, should the circumstance arise, the chatbot should fail gracefully with feedback to the user [1]. Elements of this design goal have already been broadly articulated in the literature and represent essentially table stakes for an analytic conversational agent. While prior systems have made progress on many of these goals there remain open challenges that are introduced by the combination of the above factors, the need to support conversational coherence, and the properties of LLMs.

**DG3: Allow Analysts to Proactively Manage Context.** Prior conversational V-NLIs use varying techniques for managing the state and conversational shift [41], but the novelty of using an LLM for data analysis creates new context

management challenges. Our findings suggest that managing context for analytic conversations may require new approaches compared to everyday conversations with chatbots. Analytic conversations often point to specific prior and require anaphoric references [13]. Moreover, explicit statements may also be required to shift context (e.g., uttering "retain," or "start over."). The prior WoZ study [64] only allowed participants to manage this context via utterances, but other approaches may offer more control without adding considerable overhead.

**DG4: Support Out-of-the-Box Visual Data Analysis.** LLMs often rely on a precise prompt in order to produce an appropriate response [5]. While analysts can explore ways to improve their prompts and utterances, for example via anaphoric references, end-users should be able to focus on conducting data analysis and not adhere to technology-specific vocabulary or need to repeatedly revise their utterances to produce desired results. Furthermore, effective prompt engineering requires specialized knowledge, and although recent efforts have aimed to help users explore automatically [73] or interactively [42, 59, 70] refine a prompt, these strategies interrupt the analysis and conversational flow. An effective conversational agent should manage such considerations for the end user without requiring them to delve into the technical details.

## 5  AI THREADS

We developed AI Threads, a prototype system for conducting multi-threaded analytic conversations with an LLM chatbot. AI Threads allows users to upload tabular datasets, prompt the LLM with analytic questions, review visualizations, and refine them. We take design inspiration from prior V-NLI systems in addition to common messaging apps (e.g., Slack) and analytic notebooks, while also incorporating the design goals articulated in Section 4.

### 5.1  Overview

AI Threads uses the GPT-3.5 model through OpenAI's ChatCompletion API to generate visualizations and captions based on descriptions of the dataset and user utterances. The system comprises a web-based application frontend to collect utterances and display the chatbot's textual and visual responses as well as a backend to generate and execute code that produces the responses. Figure 3 provides an overview of the AI Threads interface, which comprises two components : (A) a `Main Chat` area and (B) a `Thread` panel. The `Main Chat` area serves as a conversational interface that also shows the trail of the user's analytic process (**DG1, DG2**). The `Thread` panel allows for refinement and modifications to visualizations in the `Main Chat`. We use the separation between the `Main Chat` and `Thread` to proactively manage the conversational context (**DG3**). We refined the user experiences while iterating over early designs of the prototype with feedback from initial pilot studies.

We showcase AI Threads's capabilities primary through user studies (Sections 7 and 8) and usage scenarios (Section 9). However, in Appendix B we also show examples of AI Threads executing different edit operations (Figure 16) and for different utterance intents (Figure 17).

### 5.2  Main Chat Area

The `Main Chat` area includes utterance input area and retains an overview of the chat history. After the analysis of usability studies, we later included a ✿ button which automatically prompts the LLM to *"show me something interesting"*, an utterance entered by several participants. The analysis history in the `Main Chat` constitutes the diarized conversation between the chatbot and the end-user. The chatbot's response can be either text or static visualizations with accompanying explanatory captions.
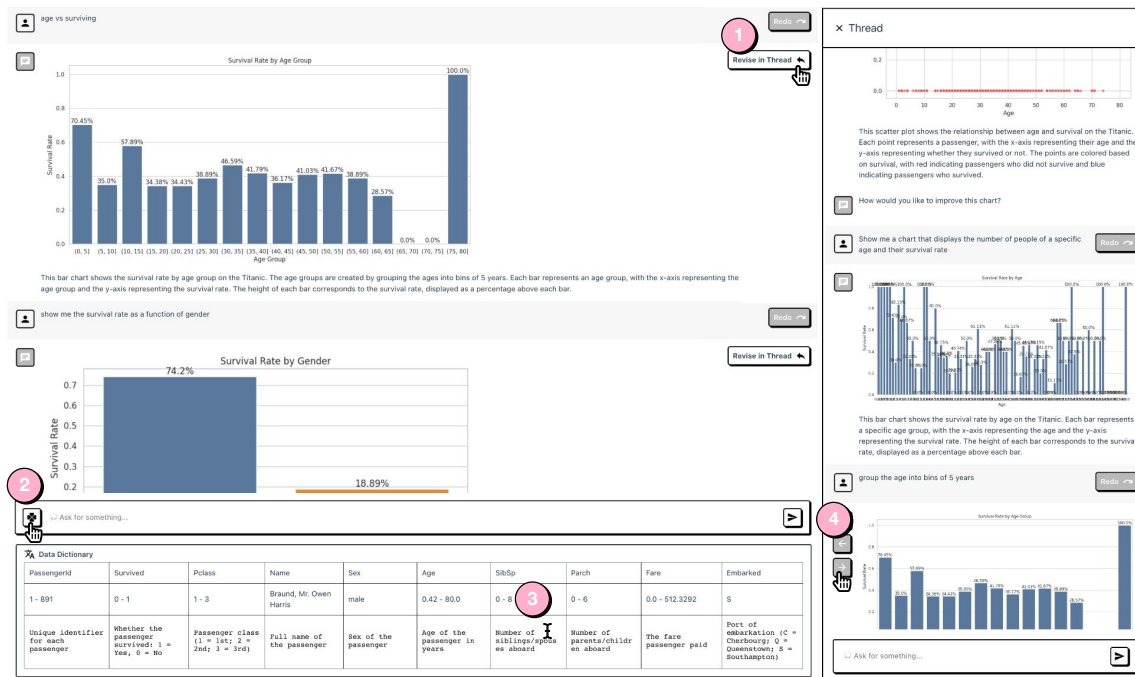
Fig. 3. **Overview of AI Threads interface.** Users begin their analysis in the `Main Chat` area (left), which shows the conversation history diarized by the user's and chatbot's responses. The chatbot's responses can be textual or data visualizations. To modify a visualization in collaboration with the chatbot the user can 1) open the `Thread` panel (right). The `Thread` panel retains a history of modifications that a user can return to and further refine at any time. Closing a `Thread` updates the visualization displayed in the `Main Chat` conversation with the latest version in the `Thread` panel. When in the `Main Chat` area, the user can 2) prompt the model to "*Show me something interesting*". The `Main Chat` also includes an 3) editable and persistent `Data Dictionary` to help orient both the user and AI during visual data exploration. When model outputs are re-generated after clicking a redo button, the user can 4) click the arrow keys to revert to previous generative outputs. Unlike ADA, reverting to a prior response using the arrow keys does not delete utterances or responses that have already occurred later in the conversation. This snapshot is from crowd-worker P13 (Section 7).

The chatbot's responses can be modified in two ways. First, if an error is produced or if the user is not satisfied with a visualization, they can click the [Redo] button to automatically re-trigger another attempt (**DG2**). Both the current and prior responses are retained and the user can scroll through them using ← or → arrows. Unlike ADA, reverting to a prior response using the ← arrow does not delete utterances or responses that have already occurred later in the conversation, allowing the user to simply select which visualization should be part of the analytic context sent to the model (**DG3**). Second, a participant can actively modify a visualization by clicking on the [Revise in Thread] button that triggers the `Thread` panel. Once a participant completes their refinements in the `Thread`, the visualization in the `Main Chat` is also updated and becomes part of the main analytic context (**DG3**) without the intermediary refinement utterances to clutter it (**DG1**). As a last resort, users can also prompt the system to "undo" their previous prompt, which deletes the last user and assistant messages. This allows users to retype their utterances if their analysis has led to a dead end.

Below the utterance input area is a `Data Dictionary` [51], which is intended to help both the AI and the user understand the structure of their dataset (**DG4**). The `Data Dictionary` breaks down each attribute in the dataset according to its name, its value range or an example, and a brief description of the attribute. The attribute descriptions

are automatically generated by AI Threads after uploading the dataset and are user-editable. The inclusion of each column's possible values and editable descriptions is intended to help the LLM resolve potentially ambiguous attribute names as well as elicit higher-order reasoning about the dataset (**DG4**; details follow in Section 6.1).

### 5.3 `Thread` Panel

The `Thread` panel comprises the same elements as the `Main Chat` area (excluding the data dictionary). It is intended to support the user in refining analysis ideas and visualizations. There may be several reasons why a user may wish to initialize a thread. One that we have already described is the progressive elaboration that is common in analytic processes (**DG1**). Upon generating an initial visualization, the user may gain an insight that motivates them to further refine their utterance. This refinement can include modifying attribute bindings to different channels (by adding or removing them), modifying the data transformations (e.g., the bin width in a histogram), adding highlights or annotations, and adding reference or regression lines among other things [29, 52]. The content of a `Thread` persists even after the participant closes the panel, allowing them to continue their analysis or share with others how their analysis evolved (**DG1**).

Other reasons may arise from a misalignment between the user's intent or desired output and the response of the chatbot. For example, the chatbot may generate an encoding type that the user does not find effective (e.g., a pie chart when a bar chart is desirable). When utterances are ambiguous, the chatbot may also make incorrect choices of data transformations or visual encodings. The end-user can provide more specific details for the chatbot in follow-up utterances (**DG2**). AI Threads is able to proactively accept and incorporate user feedback that corrects its outputs, we provide examples of these from user studies and a concrete set of edit operations in Appendix B-Figure 16 . The use of the `Thread` panel also need not be limited to modifications. They can also be used to prompt the chatbot to explain the visualization it generated, provide suggestions for alternative visualizations, and extract insights from a visualization (for example, correlations, distributions, or trends [61]).

## 6 SYSTEM DESIGN

In this section we describe AI Threads's architecture that supports ingestion of user utterances and the display the result within the `Main Chat` and `Thread`. We provide an overview of AI Threads's pipeline in Figure 4.

### 6.1 Data Sources

Participants can upload a single CSV dataset to carry out their data analysis. A string referring to the file's location on the local system is also stored and passed to GPT-3.5 as part of its code generation. When a dataset is uploaded, using the procedure outlined in Figure 6, AI Threads will automatically generate a `Data Dictionary` which includes each column's name, data type, and a value range (for nominal or ordinal columns) or a non-null example value (for all other data types). The LLM also generates a natural language description of each field.

The `Data Dictionary` is used as part of future instructional system prompts to generate valid and executable code; it is also displayed to end users (Section 5.2) as an editable reference table, where they can alter descriptions to provide more accurate information for the AI (**DG4**). These descriptions help the LLM identify and use the appropriate attribute and reason about its values (e.g., the utterances "*Filter gender by male*" vs. "*show only male*"), which supports a more naturalistic conversation out-of-the-box (**DG2**). We believe that descriptions about columns can also elicit higher-order reasoning from the LLM when necessary, for example whenever the user prompts the model to "*show me something interesting,*" by leveraging knowledge about semantically related concepts captured within the LLMs embedding space.
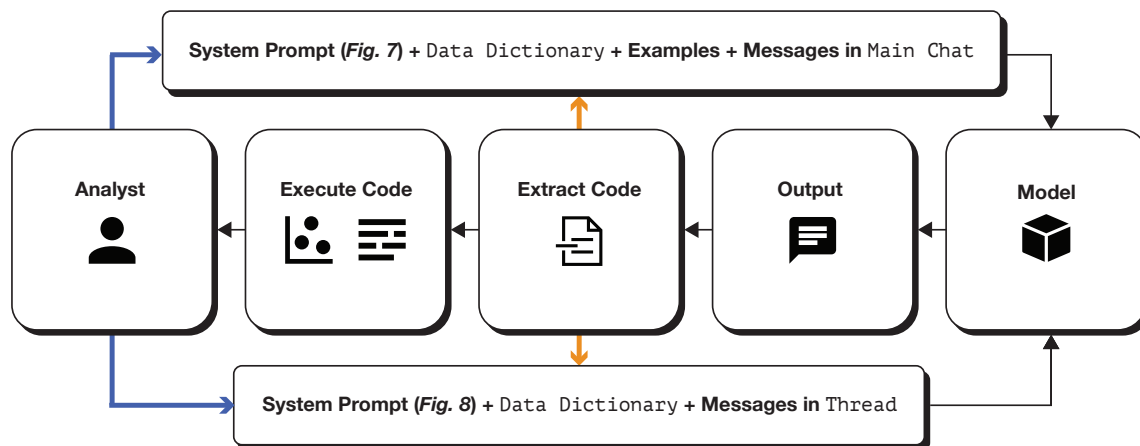
Fig. 4. **Overview of AI Threads's code generation pipeline.** When an analyst enters a **user utterance**, the message is sent to GPT-3.5 with different system prompts and conversational context depending on whether they are interacting with the `Main Chat` (top) or a `Thread` (bottom) panel. The corresponding system prompts are further described in proceeding figures. If the AI's **assistant outputs** contain executable code, the code becomes part of the subsequent conversational contexts. After the code executes, the user will see the rendered charts along with accompanying captions.

Future work presents interesting possibilities for extending this approach by augmenting the LLM with a knowledge graph [36].

AI Threads can leverage knowledge that is external to the uploaded dataset as part of GPT-3.5's training data[3] to augment their existing analysis. For example, with the Titanic dataset participants can ask *"Who was the captain on the Titanic?"*, which is not available in the dataset itself. Figure 5 shows an example where a participant from the usability interviews (Section 8) computes the total fare paid by participants and asks AI Threads to convert this value to US dollars in 2021. AI Threads produces an estimate with some guidance from the participant. To the best of our knowledge, such capabilities were not possible with prior V-NLI and conversational systems.

### 6.2 Code Generation

User utterances initiate the generation and execution of Python code to apply the appropriate set of data transformations using `pandas` and produce visualizations using `seaborn` [67] and `matplotlib` [24]. In addition to producing a visual encoding, GPT-3.5 is also instructed to produce a text caption describing the encoding choices. Depending on how the utterance is posed, GPT-3.5 sometimes also generates an interpretation of what is shown.

*6.2.1 Processing User Utterances.* Prior V-NLI and conversational interfaces have used a variety of methods, from parse trees to heuristics, to decompose a participant's utterance, deduce its intent, and generate an appropriate output. We use GPT-3.5 to perform these tasks, providing it with instructional prompts to guide its visualization code generation or refinement. Subsequently, the relevant visualization code is extracted from GPT-3.5's original responses to display charts and captions to the analyst.

---

[3]At the time of this writing, the details of GPT-3.5's training data have not been publicly disclosed.
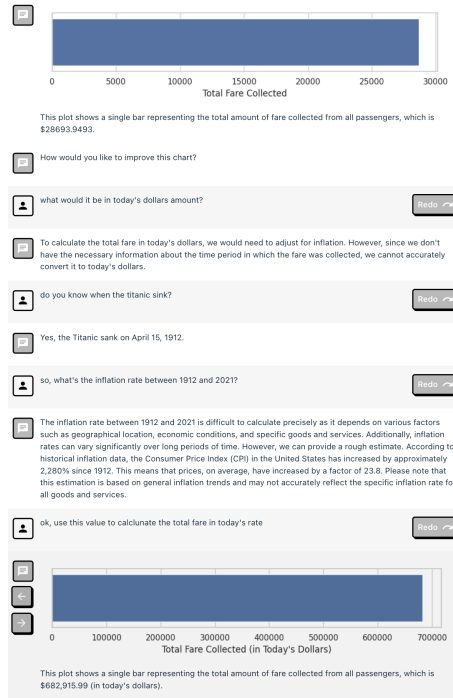
Fig. 5. **Example of the chatbot obtaining and applying external knowledge to a analysis of the Titanic dataset in a `Thread`.** A user prompts the chatbot to compute the total fare paid and convert the amount to 2021 USD. The chatbot initially failed to produce a response, but with guidance from the user (Participant E02 in from the usability interviews (Section 8)) the chatbot eventually converted $28,693.9493 to $682,915.99. The initial fare sum is correct, but the amount adjusted in 2021 USD is incorrect; the increase Consumer Price Index the chatbot used to compute this value is closer to a factor of 27.9 (according to the US Bureau of Labor Statistics), not the 23.8. However, the calculation ($28, 693.9493 \times 23.8 = 682, 915.99$) was performed correctly.

*6.2.2 System Instructional Prompts.* The GPT-3.5 ChatCompletion API delineates three role types to manage multi-turn conversations: system, assistant, and user. System instructional prompts can define an LLM's **role** and **desired behaviors** and augment the user messages to achieve a desired response.

In the `Main Chat`, the full instructional prompt provided to GPT-3.5 contains the agent's **role**, which is to answer questions about the dataset with visualizations, and **desired behaviors**, which is to render and display a chart with an accompanying caption). The system is also prompted to use **chain-of-thought (CoT) reasoning**, which has been shown to improve the performance of LLM code generation [68]. The CoT prompt is reinforced with a set of few-shot examples (See Section 6.2.3). Lastly, the system message includes the previously generated **data dictionary**. Crafting such instructional system prompts can require specialized knowledge to produce via a manual trial-and-error process [35], and thus is something we do not expose to the user. We refined our prompting strategy over time and found it to qualitatively improve the chatbot's responses without the need for the user to input additional instructions (**DG4**).

In the `Thread`, the agent is given a different **role**, which is to revise a previously generated visualization code to change visual encodings or make stylistic changes. A user might also ask the model to present additional data columns, to aggregate the data in a new way, or to filter the data, which again requires information from the **data dictionary**. However, the design of this instruction ensures that each user prompt is interpreted to be a question or a command that

Using a `pandas` script, AI Threads creates a markdown table from the uploaded dataset that includes column names, data types, and data values structured as follows:

```
|             | Data Type  | Range or Example       |
|:-----------|:-----------|:-----------------------|
...
```

This markdown table and the original filename are used to produce the full `Data Dictionary` using GPT-3.5.

Before the model is queried, we give it the following system prompt:

**System**

```
We will generate markdown tables that begin with the following columns:
|             | Data Type  | Range or Example       | Description               |
|:-----------|:-----------|:-----------------------|:--------------------------|
```

Finally, the following user message is sent to GPT-3.5:

**User**

```
Here is a markdown table summarizing columns in [FILENAME]:
[PANDAS TABLE]
Generate descriptions for what the data column values mean and concatenate the
descriptions to the markdown table.
```

Fig. 6. **Instruction prompt for generating a `Data Dictionary` of an uploaded dataset.**

continues or modifies a specific analysis (**the original code**), as opposed to starting a new analysis that will require the synthesis of fundamentally different code.

*6.2.3   Few-Shot Examples.* In the `Main Chat`, the initial system message is followed by few-shot examples of pre-designed **user** utterances and the **assistant's** code responses (Figure 7). These few-shot examples include how to transform data with `pandas` and use `seaborn` to achieve the desired image and text caption outputs. Utterances in the few-shot examples were: "*How many rows are in the dataset?*" which teaches AI Threads to display simple numeric responses as both a bar chart and within a text caption, and "*Plot the distribution of every continuous measure,*" which signals how it should display visualizations with multiple facets. Both questions are also agnostic to the semantics of specific datasets. Appendix B.2 contains full details about the examples.

*6.2.4   Code Extraction and Execution.* Responses from the GPT-3.5 ChatCompletion API are often not valid Python code excerpts that can be directly executed by AI Threads. Given that GPT-3.5 is a successor to both Codex and GPT-3 models, its responses may contain 'friendly' textual responses (usually beginning with "Sure!" [74]) or explanations of a step-by-step approach in addition to code. In order to prioritize analytic coherence (**DG1**), AI Threads only extracts a valid code block from model responses, which include the chart's caption. When valid code blocks can be extracted, they are automatically executed and outputs (visualization and text) are pushed to the front-end interface. If valid code blocks are not generated, AI Threads responds by displaying the initial model response to the user.

   For a variety of reasons, including API time-outs, the model may not generate executable code. In these cases, the participant can reattempt to generate the response from the front-end interface without the need to re-enter the utterance or provide a new one (Section 5.2). In other cases, the model may also output a simple text response, for example returning the average fare paid as a number instead of producing a visualization. In some instances, the text outputs of AI Threads pose clarifying questions to the analyst, for example in responding to the utterance *"tell me*

In the `Main Chat`, the following system instructions are provided to GPT-3.5:

**System**

```
You are a data visualization assistant working with a dataset called [FILENAME].
You will answer questions with this dataset.
When you do not use the dataset, say "I can't answer this question with the data,
but I found this:" You will use pandas and seaborn libraries to write Python
code which prints a single PNG image in a bytes64 representation. You will then
also print an alt-text caption for the chart.
Let's work responses out step by step to be sure you have the right answer.
The dataset: [DATASET DICTIONARY]
```

In the background, user conversations with the LLM are prepended with few-shot input-output examples as conversational context, e.g.:

**User**

```
How many rows are in the dataset?
```

**Assistant**

```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
...

# Load the dataset
df = pd.read_csv('./workspace/{filename}')

# Get the number of rows in the DataFrame
num_rows = df.shape[0]

# Create the bar chart using seaborn, adjusting the size of the chart as needed
width = 8
height = 2
plt.figure(figsize=(width, height))
sns.barplot(x=[num_rows])
...

# Print a caption
print("This plot shows a single bar representing the total number of rows in the dataset, 
which is " + str(num_rows) + ".")
```

The example response is truncated; see the full examples in Appendix B.2.

Fig. 7. **Pre-defined context in the `Main Chat`.** Although not exposed to the user, an exchange with the GPT-3.5 CodeCompletion API in the `Main Chat` will include a system prompt and few-shot examples along with previous exchanges.

**System**

```
Today we are being asked to revise the following code.
The revised code should print a new chart in a bytes64 PNG representation and
print a new alt-text caption in the same way as the original code.
The code: [ORIGINAL CODE]
The dataset [FILENAME] contains: [DATASET DICTIONARY]
```

Fig. 8. **Instructions in a `Thread`.** The role of GPT-3.5 in a `Thread` is to extend an analysis branch, using the **original visualization** as a reference for refinement or continuation.

*something that is cool"* an example response is *"as an AI assistant, I don't have access to personal information or preferences unless shared with me. If you have any specific preferences or questions, please feel free to share them, and I'll do my best to assist you".*

*6.2.5 Different Thread, Different Context.* AI Threads provides different system prompts and conversational contexts to GPT-3.5 depending on whether the user is interacting with the `Main Chat` and `Thread` panel (Figure 4). The conversational context is further tailored to a thread the user is inputting utterances to (Figure 8). Code that is generated within a thread is appended as prior context to subsequent utterances. Alternating between `Thread` and `Main Chat` allows the user can proactively manage context (**DG3**). This is because the `Thread` panel contains only the context relevant to that particular visualization and analysis while the `Main Chat` will only include the utterance and results of charts as refined in each thread.

*6.2.6 Alternative Visualization Libraries.* We explored other visualization libraries but found that they resulted in less reliable responses. In some instances, the libraries were too new to produce reliably executable code, as was the case for `altair` [65]. We also explored different languages for visualization specifications. Prior work [10, 47] has explored the use of JSON-formatted `Vega-Lite` specifications to produce visualizations using LLMs, but we did not use these approaches for two reasons. First, using `seaborn` allows for direct integration with its dependencies `statsmodels`, `scipy`, and `fastcluster` for advanced data modeling. Keeping code for the data operations and visualization in the same ecosystem improved AI Threads's ability to reliably generate executable code. Second, one advantage of using `seaborn` [67] or `matplotlib` was that GPT-3.5 would produce comments to document its data transformations and visualization processes, which is a useful exposition of its chain-of-thought reasoning [11] that we could leverage in the future. However, getting equivalent information from a JSON specification proved to be difficult.

## 7 USABILITY SURVEY

We conducted a survey of AI Threads's usability with 40 participants. The goals of our usability survey were to establish the diversity of use for AI Threads and participants' self-reported perceptions of their sessions, including the quality of responses and visualizations.

### 7.1 Session Protocol

All participants completed an initial intake that gathered additional information about their demographics, their programming, and data analysis background. Participants were then shown a video approximately two minutes in length[4] providing an overview of AI Threads , including examples of different queries and the use of threads. We use the Titanic dataset from [64] for this study. Participants were then directed to a the AI Threads and prompted to conduct an analysis that identified *Which factors increased one's chances of surviving the Titanic shipwreck?* We did not restrict participants from simply inputting this question to AI Threads. Participants received a completion code only after they spent at least three minutes conducting the analysis and after creating at least one thread. Following their session, participants were guided to a set of questions that were used as an attention check to establish whether or not they meaningfully explored the dataset. Finally, participants were asked to rate their session, chatbot responses, and visualization results on a series of five-point Likert scale questions while also providing justifications for their responses in free-form text.

---

[4]This video can be viewed here: https://www.youtube.com/watch?v=Ll0nPtD3kMc

### 7.2   Data Collection and Analysis

In addition to each participant's session log, we collected participant data from a survey that gathered Likert scale and free-form textual responses. We summarize all Likert scale questions using simple descriptive statistics. Participant utterances were analyzed using a ground theory approach. The session logs include a full transcript of the conversation, which constituted the participants' typed utterances as well as the system and AI Agent responses, including both the code generated for data transformations and visualization specifications and finally a PNG image of the resulting visualization presented to the user. During our analysis, we replayed these user sessions using AI Threads in read-only mode.

### 7.3   Participant Recruitment

Crowd-sourcing platforms are increasingly used by visualization researchers to recruit participants for a variety of tasks, including, as we do here, evaluating user interfaces [3]. We recruited participants from the Prolific platform[5], aiming to collect a balanced gender sampling. While our study does not require specialized analysis skills or knowledge of the Titanic dataset, our inclusion criteria aimed to increase participants' likelihood of having some analytic experience (for example, by having knowledge of certain programming languages or systems). We also restricted our study to English speakers in the U.S., with at least 100 approved tasks. We paid a rate of $3.50 for 20 minutes of effort, which equates to $10.50\hour—higher than the US federal minimum wage. We screened all participants' textual responses and session logs, removing participants that did not use AI Threads in a manner relevant to the study tasks.

### 7.4   Findings

*7.4.1   Participant & Session Overview.* Our study had a roughly even breakdown of participants that identified as female (n=19, 47.5%) or male (n=21; 52.5%). A total of 7 (17.5%) participants indicated they were currently students. Nearly half (n=17; 42.4%) of participants had background in one of either computer science, data science, machine learning, or natural language processing. Moreover, 35 (70%) reported routinely conducting data analysis, primarily with spreadsheet programs (Excel or Google Sheets; n=26; 65%) and/or scripting languages (Python or R; n=9; 22.5%); two participants reported having previously analyzed the Titanic dataset. Despite the technology's recency, 9 participants (22.5%) reported having previously used an LLM tool like ChatGPT to carry out data analysis.

The average session was approximately 24 minutes. The shortest session was 10 minutes and the longest was 65 minutes; the majority of sessions were closer to the average time than either of these extremes. Although there was nothing restricting what prompts participants could enter, it was interesting to note that approximately half began their session simply by pasting the task prompt into AI Threads. The rest began their session by asking other questions.

*7.4.2   Survey Responses.* Participants perception of the user interface, chatbot's responses, and the quality of the visualizations are summarized in Figure 9. Overall, the majority of participants agreed or strongly agreed that it was easy (77.5%) and enjoyable to use the Chatbot (72.5%). The majority also felt that the Chatbot understood what they were asking it (77.0%) and that the visualizations it produced were appropriate for the utterances they inputted (87.5%) and were easy to understand (71.8%). While some participants encountered errors while using the Chatbot, the majority (67.5%) did not find these errors to be onerous. While in aggregate participants responses report an overall positive experience, individual sessions still varied. For example, P11 reported that *"the answers seemed to be well thought out and professional. It seemed to know the exact answer I was looking for while providing additional information."* While P39
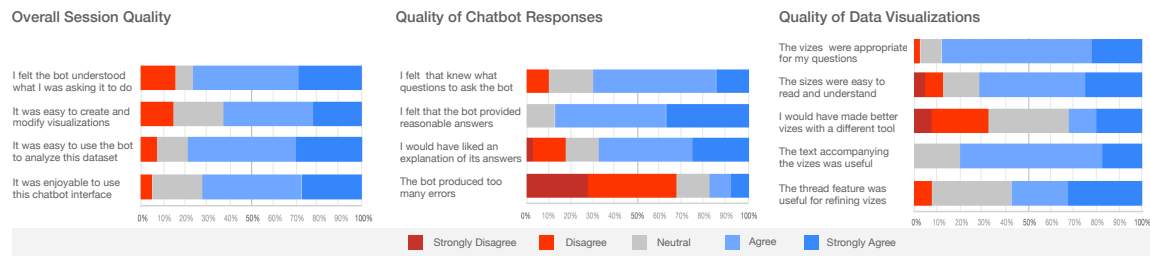
---

[5]https://app.prolific.co/

Fig. 9. **Summary of participants impressions of interacting with AI Threads**.

took the opposite view and found that *"it felt rather rough, like something still being worked on. It took a bit of effort to get the data I wanted."* From participants textual responses that justified their ratings, we identify several recurring themes.

**It was sometimes necessary to modify an utterance to achieve the correct response.** Several (n=6) participants indicated that the chatbot did immediately understand their utterance and they *"had to word things differently sometimes to get the results I was looking for"* (P26). At times the Chatbot interpreted an utterance too literally as P05 observed: *"the Chatbot will generally give you pretty much what you asked for, so it is important to understand how to talk to the Chatbot"* P06 noted that ambiguity of the utterance contributed to the importance of phrasing: *"I noticed that I had to word my questions carefully and specifically. For example, I asked the bot 'Which factors increased one's chances of surviving the Titanic shipwreck?', a very general and broad question, and the chart it came out with was basically empty. This was understandable considering the variety of factors that could be given as an answer to the question".* This observation is in conflict with **DG2**, but the majority of participants did not report having this experience. Overall this observation demonstrates how variable interactions with Chatbots can be.

**Participants thought visualizations were appropriate for their questions, but needed some improvements.** P18 expressed the general sentiment that *"visualizations were clear and a majority of the time were what I asked"*. Again, there was a spectrum of experience with P20 reporting that *"The visualizations were just okay. They often did not correlate to what I asked for, were not labeled well, and were not visually appealing"* Participants also had different experiences modifying visualizations. One participant appreciated that *"modifying them [the visualizations] only required simple commands."* (P25), while others reported that sometimes *"it [the chatbot] also pulled from incorrect data and had a hard time fixing it"* (P30). One participant provides a specific challenge: *"creating a chart breaking down survival rate by age, since it tried to do a bar graph for each individual age rather than breaking them into larger age ranges."*

**Threads were used to modify visualizations & ask for more information.** The majority (57.5%) of participants thought the threads features was useful, and from reviewing the sessions logs we observed its most common use was changing the encoding types. However, the overall number of ways participants used threads was quite varied. In Figure 10 we show examples of how participants modified visualizations through the useful of threads. P13 was enthusiastic about their experience: *"I really enjoyed the thread feature to edit the graphs and make them a bit easier with more information on what was wanted".* P31 describes their analysis process: *"When asking for the median age, it produced a box plot that was difficult to analyze. I opened threads and asked it to change it a histogram, which was easier to read."* Participants also used threads to get explanations about the chart (Figure 11). The explanations were useful to P05 who found that *"the chatbot helped further clear things up [about the chart] when I asked about it in the threads"* P18 had an interesting scenario of correcting the chatbot when it claimed that error bars were not present in the chart.

Fig. 10. **Example use of threads to modify analyses and charts.** We added grey pills to summarize the key refinements.

However, it was evident that the chatbot was not relying entirely on the chart and that data when creating explanations. In the P35's explanatory context, the chatbot provides an incorrect explanation of the chart; it hallucinated the response. We observed other instances of this kind of hallucination occurring, but it was mixed in with relevant and correctly generated responses – making it difficult to precisely identify factors that caused this to occur.

Overall, from this usability survey the use of threads was generally shallow. We noted that participants appearing to have stronger analytic backgrounds made more extensive and meaningful uses of threads, and were more likely to report positive experiences using them. This observation was the motivation for our follow-up study interviewing analytic experts.

### 7.5    Summary and Implications

In general, participants experiences align with our own investigation of the data from Tory and Setlur [64] and showcase the diversity of experience toward AI Chatbot use for data analysis. While individual sessions showcase a spectrum of experiences, the majority of participants reported being both pleasantly surprised and delighted by the accuracy of the chatbot's interpretation of their utterances and its responses. We used participant's responses to revise AI Threads by addressing errors these sessions surfaced and improving the overall responsiveness of the bot.
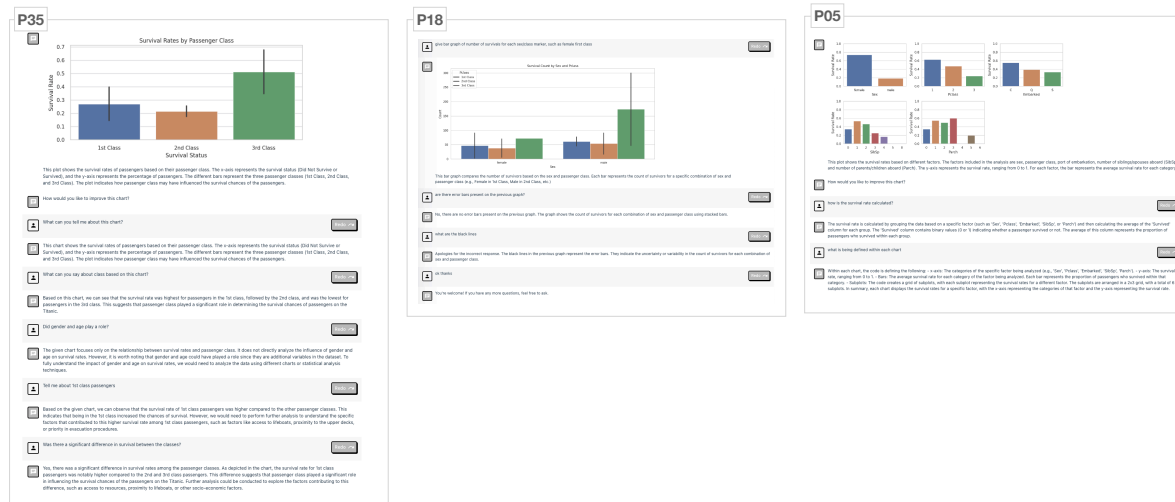
Fig. 11. **Example use of threads to explain information in the chart**

## 8 USABILITY INTERVIEWS

We follow up on the results of the usability survey by conducting semi-structured interviews with 10 individuals with data and visual analysis expertise.

### 8.1 Study Protocol

The participants completed an initial intake identical to the usability survey (See Section 7). Following this intake, participants were contacted to participate in a 30-minute follow-up session to review their analysis. In some instances, participants did not complete their analysis before the follow-up session, when this was the case they completed their analysis during the follow-up session, extending the total session time up to 45 minutes. Each session included one study administrator and a note-taker. The study administrator primarily followed and script to ask participants about their experience; where appropriate the administrator also followed up with improvised probing questions.

### 8.2 Participants

We set a recruitment goal of ten participants with data and visual analysis experience through social media. Participants who completed the intake form but that were ineligible to participate in this study, primarily because they did not have the requisite background, were not contacted for the follow-up study and their results were not used. We continued to leave the intake form open until ten suitable participants were identified. The participants we predominately male (n=9) with technical roles as a data scientist/engineer (n=2), data visualization experts (n=4), business analyst (n=4), and machine learning engineer (n=1); participants could select more than one role. Three participants also identified with other roles including engineer managers (n=2) and one technical writer. Participants reported using Python/R (n=2), Tableau/Lookr/PowerBI (n=4), and Excel/Google Sheets(n=4) to conduct data analysis. Participants also reported using Dataiku (n=1), Spotfire (n=1), and Smartsheet(n=1). Three participants said they had used LLMs to conduct data analysis. Two participants had previously used OpenAI's ADA to conduct analyses of different datasets.

Table 1. Categories of User Feedback

| Theme | Sub-themes | Example Quote |
| --- | --- | --- |
| Role of AI Threads in Analysis ($n = 8$) | Sociotechnical ($n = 8$) | "It's good to get familiar with a dataset and potentially be able to share data visualizations to their team or to somebody higher up" [E06] |
| | Data Engineering ($n = 7$) | "How do I kind of define these variables that I want to work with... and know that my chart is using that?" [E01] |
| Reflections on System Features ($n = 10$) | Ability to Organize Analysis ($n = 10$) | "Keeping the context of the conversation intact is actually a very useful feature... you can walk down this garden path but still not lose your train of thought." [E03] |
| | Having a Two-Way Conversation with AI ($n = 9$) | "I kind of learned I guess that maybe for more complex things you have to be a little more specific."[E08] |
| | Enchantment and Captivation ($n = 7$) | "I was surprised to get graphs... It was a nice surprise to actually get those." [E05] |

## 8.3 Data Collection & Analysis

As with the previous usability survey, all participants rated their experiences using Likert Scale data; participants did not have to provide textual justification for their scores as this information was collected during the follow-up interview. We again collected the full session logs for all participants that we reviewed with them during the follow-on session. Finally, all sessions were recorded and transcribed. We conducted a thematic analysis[4] to summarize participants' impressions of the tool's usability as well as questions related to their analysis, including their analytic intent behind prompts and the quality of AI's responses to their prompts. Themes and subthemes were synthesized using a ground theory open-coding approach [6].

## 8.4 Findings

We summarize participant impressions and suggestions in Table 1, which shows the themes and sub-themes that emerged from our thematic analysis. We identified two primary themes centering on 1) the role of AI Threads in data analysis and 2) reflections on existing features and suggestions for additional features. Participants reported an overall positive experience with only one reporting having an overall negative experience. This participant (E01) expressed wanting greater control of AI Threads's outputs and felt that natural language inputs did not enable such fine-grained refinements. Despite this, they still report having *"a magical moment [...]"* using AI Threads *"that's getting into Star Trek Communicator Land [because] I don't have to say, rebuild this chart or whatever. It just got it.".*

*8.4.1 The Role of AI Threads in Data Analysis.* Participants' overall impressions of AI Threads's usability largely align with the results previously reported in Section 7. From the interviews, we were able to get a better understanding of participant's overall analytics workflow and the ways the conversational agents fit into them. Two themes consistently emerged. First, nearly all (n=8) participants commented on the ability to share the analysis with others. Second, participants wanted to explore know the capabilities of AI Threads on larger and streaming datasets that were reflective of their common workflows. We expand on these themes here:

**Sociotechnical Considerations.** system's role was perceived predominantly as a useful tool for rapid exploratory analysis of a dataset. However, exploratory analysis exists as part of a socially structured data science workflow, and participants were quick to speculate how the data artifacts generated by AI Threads could be used to facilitate

collaboration. As E06 noted, AI Threads was *"good to get familiar with a dataset"* with the objective of *"potentially be able to share data visualizations to their team or to somebody higher up."* One participant (E02) speculated that in the near future, LLMs would be integrated into messaging applications such as Slack or Microsoft Teams where multiple analysts could collaboratively analyze a dataset.

On the other hand, two participants (E08 and E09) envisioned AI Threads as a separate tool that may serve as a replacement for Jupyter Notebooks, and the figures could be exported and shared with others. With these considerations in mind, it was natural for E04 to question *"will it take into account if it's for leadership, or in tech, or a data scientist?"* if they could control the persona of the LLM.

**Data and Processing.** Participants also considered what would be required before using AI Threads for data analysis. For example, E07 wondered *"How will it handle large datasets"*especially when the tables were in wide format (e.g., each column is a date). Particular considerations to the data engineering workflow were given when participants were used to making dashboards for tracking assimilated and streaming data. Even if the data were tidy and stable, E01 wondered *"how do I kind of define these variables that I want to work with... and know that my chart is using that?"* Although in E01's session, they were happy that AI Threads was able to automatically infer that they wanted to plot survival rate (instead of survivor count) without explicitly specifying the calculation, they were still concerned about the lack of provenance tracking with regards to the AI's data transformations.

*8.4.2 Reflections on System Features.* While using AI Threads participants also reflected on ways to expand its capabilities to further support them in data analysis. We break these into three sub themes:

**Ability to Organize Analysis.** Expert participants were overwhelming positive about the ability to organize their analysis using AI Threads, particularly with the threads feature. E07 noted that they wouldn't *"want this information in the thread to be used in the main chat."* and E08 reflected on how they effectively used the feature as they *"Switched to main chat when thread was different enough to start over."* The conversation in the `Main Chat` represented where the major data artefacts should be stored, and E03 noted that *"keeping the context of the conversation intact is actually a very useful feature... you can walk down this garden path but still not lose your train of thought"* **Collectively, these findings mean that AI Threads could be effectively used to keep track of the analysis process to maintain a coherent conversational memory not just for the AI, but also for the user.**

On the other hand, several participants noted that additional features could support the organization of data artifacts, such as bookmarking messages in both the `Main Chat` and the `Thread`. The ability to title Threads was also a frequently requested feature to support data sensemaking. E10 expressed that while the threads feature was useful for follow-up on analyses, but felt that it could still be limited as the analysis grew more complex. Moreover, they also pointed to issues with generative AI tools when *"trying to go back to a specific place in the conversation, it's not obvious where the entry points are to each thread."*. They further expanded that when analyzing data they *"got this exploding fractal mess, because that's what happens inside. And then ultimately, I'll go back and try and pick out the ones that are liked."* However, they acknowledge this was a common problem with all analysis tools and that extending threads is a promising way to address these issues. While E08 also speculated about a multi-level thread feature, they noted their concerns with *"going too far with threads and subthreads... it can be something like Slack, but something like Reddit would be too much."*

**Having a Two-Way Conversation with AI.** Most of the negative feedback about AI Threads concerned the passive role the AI seemed to play in the design of charts. E01 acknowledged that in order to analyze data effectively, people need to have expertise in both data analysis and the dataset itself. Before entering their first prompt, E09 took a significant

pause wondering what questions they should ask even though the dataset was about an infamous historical event with well-known narratives. Starting with AI suggestions would have been useful for them to get into the flow of analysis immediately.

Analysts would rather be analysts and not prompt engineers. Four out of nine expert users hoped to delegate aspects of the prompt design itself to the AI in the form of AI-drive chart recommendations prior to the start of data analysis. Such recommendations could take the form of a "suggestion box" similar to Show Me [39]. Experts will still have a hard time recollecting the names of charts (e.g., 100% stacked bar charts) and E08 mentioned that a chatbot making suggestions like *"Would you like to see a chart that combines [these two variables]?"* would be useful. Without such suggestions, as reflected in the behavior and utterances of 7 out of 10 participants, LLM users are *"mostly just playing around with what [the AI is] capable of"* at the beginning of sessions. To this end, E08 also noted that in-depth explanations of why AI generated its responses or threw unresolvable errors would be useful for users to *"form a mental model"* for how to prompt the AI model more effectively.

Both E04 and E05 noted that the model's text responses were "*very good.*" However, E05 speculated that an information bubble that reveals some out-of-the-box supporting narratives about why a generated visualization might be important in the data context may be useful to have.

**Enchantment and Captivation.** Generally, users were pleasantly surprised by AI Threads's ability to generate shareable visualizations from natural language alone. When comparing to existing public-facing AI systems, E05 was pleasantly surprised that AI Threads generated charts by default, remarking that generally *"graphs are a good way to"* get *"an answer as quickly as I could."* Concurrent with our analysis of ADA (Section 3) one participant stated that when they *"use ChatGPT and other similar tools and it's mostly text. Unless I specifically say: "Present this in the table," or "Give me steps," or something like that,"*. This highlights the importance of LLM-driven visualization interfaces to go beyond the ability of foundation models by automatically rendering code into desirable images, especially because visual data exploration is usually what is expected in earlier phases of data analysis. E02 remarked they saw interactive systems like AI Threads *"as the future of analytics."*

Whenever AI Threads was able to correctly assume design choices without explicit prompt engineering, users were captivated by the result. E01 noted that it was *"magical"* that they did not need to specify the number of bins in a histogram. They shared another such experience: *"When I did the survival by pclass [...] what I really meant was the rate, then it gave me exactly the rate. Those are perfect to me."* When E03 received a response to an utterance even though they did not specify a visual encoding, they remarked that *"In my mind I would have never gone with a stacked bar chart as my first option"* but that it was the appropriate response upon reflection. E07 noted that such automation highlighted the potential for NL-driven visual data analysis to empower them to get *"a lot of analysis done."* E01 speculated that AI could eliminate much of the time and effort traditionally needed to make complex visualizations and highlighted the potential for LLMs to overcome the traditional trade-offs between learnability and expressibility.

## 9   USE CASE

Finally, we demonstrate the use of AI Threads on a dataset collected after GPT-3.5's knowledge-cutoff date of September 2021. The Armed Forces of Ukraine has published data related to its military operations since the first week of the war, which began in 2022. A structured dataset was compiled from the Ukranian government's website as daily reports and continues to also be updated weekly on Kaggle. The dataset contains information about the cumulative number of losses per day sustained by Russian forces, including both military personnel and equipment. We consider a usage

| date | day | aircraft | helicopter | tank | APC | field artillery | MRL | military auto | fuel tank | drone | naval ship | anti-aircraft warfare | special equipment | mobile SRBM system | greatest losses direction | vehicles and fuel tanks | cruise missiles |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2022-02-25 | 2 - 543 | 10 - 315 | 7 - 316 | 80 - 4346 | 516 - 8435 | 49 - 5245 | 4 - 717 | 100.0 - 1701.0 | 60.0 - 76.0 | 0 - 4304 | 2 - 18 | 0 - 489 | 10.0 - 790.0 | 2.0 - 4.0 | Sievierodonetsk | 1796.0 - 7680.0 | 84.0 - 1406.0 |
| The date when the data was collected. | The numeric representation of the day when the data was collected. | The number of aircrafts lost by Russia. | The number of helicopters lost by Russia. | The number of tanks lost by Russia. | The number of Armored Personnel Carriers (APC) lost by Russia. | The number of field artillery lost by Russia. | The number of Multiple Rocket Launcher (MRL) systems lost by Russia. | The number of military vehicles lost by Russia. | The number of fuel tanks lost by Russia. | The number of drones lost by Russia. | The number of naval ships lost by Russia. | The number of anti-aircraft warfares lost by Russia. | The number of pieces of special equipment lost by Russia. | The number of mobile Short Range Ballistic Missile (SRBM) systems lost by Russia. | The direction with the greatest losses for Russia (given as an example here). | The combined number of vehicles and fuel tanks lost by Russia. | The number of cruise missiles lost by Russia. |

Fig. 12. **Ukraine-Russian war dataset**. The `Data Dictionary` for this dataset was automatically derived by AI Threads.

scenario of a data journalist, Sam, as they analyze this data using AI Threads. A summary of Sam's analysis is shown in Figure 13.

Preparing for a story pitch meeting, Sam wants to quickly explore the dataset to find compelling visualizations that can be used to support their investigation. They upload the dataset and AI Threads generates a data dictionary including the description of each column (Figure 12). Sam begins by typing: *"Show trend in each equipment."* AI Threads correctly responds with a line chart displaying the cumulative loss of each Russian piece of equipment throughout the war. However, all 15 equipment types are difficult to distinguish using a single color palette. Sam directs AI Threads to '*"Show trend in all aerial equipments"*. The simplified line chart now clearly communicates that Russian drone losses are rising day by day. This is in stark contrast to the stagnation in the cumulative losses of aircraft and helicopters, even though the war has escalated.

Sam understands that since the data is maintained by the Ukrainian forces, it may inflate the number of true Russian equipment losses. So Sam opens a `Thread` next to the line chart and asks: *"These numbers might be inflated, but it might be still useful to show equipment losses relative to each other. What can I visualize instead of absolute losses?"*. AI Threads responds by recommending a stacked area plot to present the data, so Sam further elaborates by first uttering *"Can you make a stacked area plot?"* and subsequently to *"Show as a 100% stacked area chart"* and finally to *"Show the original*
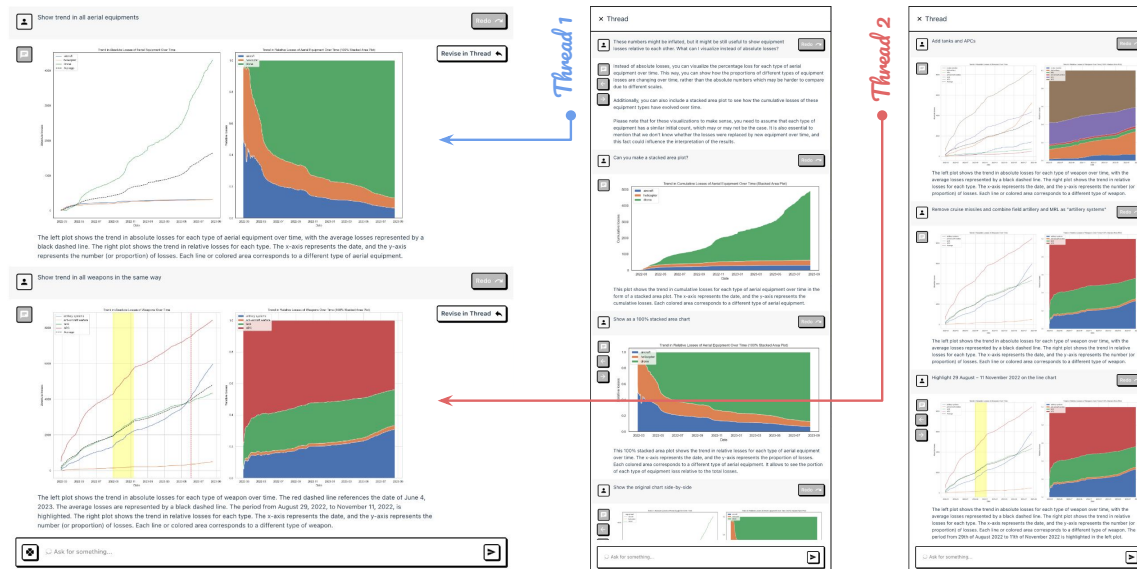


Fig. 13. **Example analysis for Ukraine-Russian war dataset**. Two visualizations and shown in the `Main Chat` area in addition to their modifications via two separate `Thread` panels.

*chart side-by-side"* – the result is combining line and area charts can highlight both absolute and proportional values. Finally, the line chart would seem more insightful if it depicted the average loss across all equipment types as a baseline to compare the trends of each piece of equipment. Sam directs AI Threads to draw the average loss as a dashed line on the line chart and exits the `Thread`.

Sam is pleased with their first visualization, which they see now in the `Main Chat` with the AI. However, the visualization only depicts a small subset of the data. Interested in building a narrative about well-known Ukrainian counteroffensives, Sam directs AI Threads to *"Show trend in all weapons in the same way"* (i.e., the same way as the previous result which Sam heavily refined in a separate `Thread`). Although AI Threads's new visualization correctly presents some weapon losses as a paired line and area chart, Sam is not entirely happy with the weapons shown.

Sam opens a `Thread` next to the figure. They direct AI Threads to *"Add tanks and APCs"* since tanks and APCs (Armored Personnel Carriers) are also used in offensive operations. Sam then directs the AI to *"Remove cruise missiles and combine field artillery and MRL as artillery systems."* Cruise missiles are single-use ammunition, unlike the rest which are multi-use equipment; field artillery and MRLs are both mid-range artillery systems. Sam now notices visual trends emerging in the line chart: two notable periods where there were relatively steeper inclines. They wonder if they correlate with the launch of major counteroffensives. Sam directs AI Threads to *"Highlight 29 August – 11 November 2022 on the line chart"* the duration of Ukraine's first counteroffensive, and is pleased to see the first dramatic increase across Russian weapon losses overlaying perfectly with the highlighted period.

Finally, the second counteroffensive was launched on the 4th of June, 2023 [6]. Sam directs AI Threads to *"Draw a reference line on June 4, 2023"* and notices that the dramatic rise in Russian artillery systems lost corresponds reasonably with the date. Sam exits the `Thread` and is ready to share the `Main Chat`, complete with charts plus titles and captions explaining their visual elements, with their colleagues in the pitch meeting.

## 10 DISCUSSION

The capabilities of Large Language Models (LLMs) for generating code and executing data analysis tasks have introduced a new paradigm, potentially overshadowing existing techniques for Visual Natural Language Interfaces (V-NLIs). However, it is vital to acknowledge the limitations of LLMs while appreciating their potential. In our exploration of these capabilities, we revisited the findings of a previous Wizard-of-Oz study [64]. Our analysis revealed that while existing LLM conversational agents excel at generating individual visualizations, they face challenges in accommodating shifts in conversational contexts, especially when transitioning from refining previous analyses to introducing new queries. We also noticed an intricate issue where the iterative refinement inherent to data analysis complicated the conversational history, making it difficult to navigate and comprehend

To address these challenges, we introduced a system, AI Threads, designed to harness the potential of LLMs as analytic agents. AI Threads facilitates multi-threaded conversations, enhancing a chatbot's contextual focus and improving conversational coherence. Through comprehensive user studies – one with crowd workers and another with analytic experts – we ascertained that AI Threads can adeptly handle a vast range of participant utterances, ensuring smooth transitions between threaded and non-threaded dialogues. Notably, AI Threads proved its competence even with datasets beyond the knowledge scope of LLMs. Collectively, our research underscores both the strengths and weaknesses of LLMs in their role as analytic agents. In the following sections, we delve into the limitations of our study and chart out potential avenues for future research.

---

[6]https://www.reuters.com/world/europe/russia-says-its-forces-thwarted-major-ukrainian-offensive-2023-06-04/

### 10.1  Limitations

Our research has several limitations. *First, a direct comparison with the closed-source ADA system remains challenging due to the undisclosed details of its implementation and update cadence.* Nonetheless, our open-source release of AI Threads offers a platform for the visualization and HCI communities to iterate and potentially match or surpass ADA's capabilities. Over time, the research community may develop a wider suite of open-source tools that not only emulate ADA but also pave the way for more transparent comparisons. *Second, the absence of quantitative benchmarks for visual analytic conversations led us to rely primarily on qualitative evaluations based on conversational shifts from [64].* While the lack of quantitative benchmarks limits our assessment of AI Threads's performance at scale, our approach offers invaluable insights into real-world analytic chatbot interactions. Future work could involve the establishment of robust benchmarks to assess the efficacy of such systems. *Third, we intentionally chose the Titanic dataset to maintain alignment with [64], but this dataset is likely part of GPT 3.5's training data.* It's worth noting that while AI Threads exhibits strong performance on this dataset, its proficiency may vary across diverse datasets. *Finally, the current AI Threads iteration does not delve into interactive visual encodings.* Our exploration revealed that blending data operations with visual analysis in different languages can be counterproductive. While we opted for static visualization libraries for consistency, future iterations of AI Threads might bridge this gap, tapping into the potential of interactive visual libraries.

### 10.2  Future Work

Our findings point to several fruitful areas of future research for visual analytic chatbots:

**Enhancing Conversational Context Management:** While extended token limits in LLMs allow for richer conversational contexts, they might not suffice for analytic dialogues, which exhibit unique patterns distinct from typical conversations. In AI Threads, we employed a user-driven transition between the `Main Chat` and `Thread` for context management. However, deep into an analysis users may also forget to open a thread. The potential exists for implicit transitions, similar to the Wizard in [64]. There are important design trade-offs between explicit (click on button) and implicit (inferring from utterance) initialization of a thread and focusing on the context provided to the chatbot. Exploring these different design trade-offs is out of the scope of our present work. Balancing the merits and challenges of explicit versus implicit context transitions, and exploring ways to refine these transitions based on user behavior, present fertile grounds for exploration.

**Evolving Conversational Interfaces for Analysis:** Feedback from usability interviews (Section 8) indicates a need to advance the conventional chat interface for data analysis. Participants sought enhanced control over the analytic dialogue flow and its evolution. Leveraging existing research on visual analytic provenance [18, 49] can improve chatbot user experiences. Further, integrating features like analysis suggestions [61] can further enhance the interactive dimension and collaborative potential of chatbots like AI Threads. Finally, participants in the usability interviews emphasized the importance of sharing their analyses with others. While AI Threads allows analyses to be stored and revisited, being able to effectively extract and organize parts of the conversation to share remains an open challenge. Situating AI chatbots with the larger sociotechnical considerations of their use is essential for guiding further explorations in this area [17].

**Instituting Analytic Chatbot Guardrails:** The persuasive capabilities of computer technology [15] and visualization systems in particular [46] have long been known. However, they should be re-examined with greater urgency in light of AI analytic agents. Hallucinations with LLMs continue to be a challenge and these issues extend to data analysis and

visualization. AI Threads appeared to have reasonably high accuracy (e.g., analyzing the correct columns, executing appropriate data transformations, generating the appropriate encoding, and carrying out correct modifications). However, it also failed in ways that were not immediately obvious, such as changing the underlying data in the chart (Figure 1, P40) or providing an incorrect explanation (Figure 11, P35)). Instances of these errors are interleaved with reasonable and correct responses from AI Threads, and as such, could be easily overlooked. They also did not occur consistently, as some sessions had few or no errors while others had more. Developing guardrails that monitor these issues and address them is important for ensuring the accuracy of analyses and the trustworthiness of analytic chatbots.

## 11 CONCLUSION

The program synthesis capabilities of Large Language Models provide new opportunities and challenges for the construction of analytic conversational interfaces. We designed AI Threads to explore some of these preliminary challenges concerning the management of context and the preservation of conversational coherence. Through user studies and a usage scenario, we show that AI Threads is capable of supporting both the initial creation and refinement of visual encodings. Moreover, we show that its capabilities can support a richer conversational dialogue that alternates between text and visuals. While AI Threads and LLMs, in general, are not without their limitations, our findings point to their compelling potential and important avenues for future work.

## REFERENCES

[1] Saleema Amershi, Dan Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi Iqbal, Paul N. Bennett, Kori Inkpen, Jaime Teevan, Ruth Kikin-Gil, and Eric Horvitz. 2019. Guidelines for Human-AI Interaction. In *Proc. CHI'19*. 1–13. https://doi.org/10.1145/3290605.3300233

[2] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. 2021. Program Synthesis with Large Language Models. *ArXiv* abs/2108.07732 (2021). https://arxiv.org/abs/2108.07732

[3] R. Borgo, L. Micallef, B. Bach, F. McGee, and B. Lee. 2018. Information Visualization Evaluation Using Crowdsourcing. *Computer Graphics Forum* 37, 3 (2018), 573–595. https://doi.org/10.1111/cgf.13444

[4] Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative research in psychology* 3, 2 (2006), 77–101.

[5] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. *ArXiv* abs/2005.14165 (2020). https://api.semanticscholar.org/CorpusID:218971783

[6] Kathy Charmaz. 2006. *Constructing Grounded Theory*. Sage Publications, London; Thousand Oaks, Calif.

[7] Qing Chen, Ying Chen, Wei Shuai, Ruishi Zou, Yi Guo, Jiazhe Wang, and Nana Cao. 2023. Chart2Vec: A Universal Embedding of Context-Aware Visualizations. *ArXiv* abs/2306.08304 (2023). https://arxiv.org/abs/2306.08304

[8] Liying Cheng, Xingxuan Li, and Lidong Bing. 2023. Is GPT-4 a Good Data Analyst? https://arxiv.org/abs/2305.15038

[9] Anamaria Crisan, Shannah E. Fisher, Jennifer L. Gardy, and Tamara Munzner. 2022. GEViTRec: Data Reconnaissance Through Recommendation Using a Domain-Specific Visualization Prevalence Design Space. *IEEE Transactions on Visualization and Computer Graphics* 28, 12 (2022), 4855–4872. https://doi.org/10.1109/TVCG.2021.3107749

[10] Victor C. Dibia. 2023. LIDA: A Tool for Automatic Generation of Grammar-Agnostic Visualizations and Infographics using Large Language Models. *ArXiv* abs/2303.02927 (2023). https://arxiv.org/abs/2303.02927

[11] Iddo Drori, Sarah Zhang, Reece Shuttleworth, Leonard Tang, Albert Lu, Elizabeth Ke, Kevin Liu, Linda Chen, Sunny Tran, Newman Cheng, Roman Wang, Nikhil Singh, Taylor L. Patti, Jayson Lynch, Avi Shporer, Nakul Verma, Eugene Wu, and Gilbert Strang. 2022. A neural network solves, explains, and generates university math problems by program synthesis and few-shot learning at human level. *Proceedings of the National Academy of Sciences* 119, 32 (Aug. 2022). https://doi.org/10.1073/pnas.2123433119

[12] Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lorraine Li, Liwei Jian, Bill Yuchen Lin, Peter West, Chandra Bhagavatula, Ronan Le Bras, Jena D. Hwang, Soumya Sanyal, Sean Welleck, Xiang Ren, Allyson Ettinger, Zaïd Harchaoui, and Yejin Choi. 2023. Faith and Fate: Limits of Transformers on Compositionality. *ArXiv* abs/2305.18654 (2023). https://api.semanticscholar.org/CorpusID:258967391

[13] Ethan Fast, Binbin Chen, Julia Mendelsohn, Jonathan Bassen, and Michael S. Bernstein. 2018. Iris: A Conversational Agent for Complex Tasks. In *Proc. CHI'18*. 1–12. https://doi.org/10.1145/3173574.3174047

[14] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, and Ming Zhou. 2020. CodeBERT: A Pre-Trained Model for Programming and Natural Languages. *ArXiv* abs/2002.08155 (2020). https://arxiv.org/abs/2002.08155

[15] BJ Fogg. 1998. Persuasive Computers: Perspectives and Research Directions. In *Proc. CHI'98*. ACM Press/Addison-Wesley Publishing Co., USA, 225–232. https://doi.org/10.1145/274644.274677

[16] Tong Gao, Mira Dontcheva, Eytan Adar, Zhicheng Liu, and Karrie G. Karahalios. 2015. DataTone: Managing Ambiguity in Natural Language Interfaces for Data Visualization. In *Proc. UIST '15*. https://doi.org/10.1145/2807442.2807478

[17] Lorentsa Gkinko and Amany Elbanna. 2023. The appropriation of conversational AI in the workplace: A taxonomy of AI chatbot users. *International Journal of Information Management* 69 (2023), 102568. https://doi.org/10.1016/j.ijinfomgt.2022.102568

[18] David Gotz and Michelle X. Zhou. 2008. Characterizing users' visual analytic activity for insight provenance. In *2008 IEEE Symposium on Visual Analytics Science and Technology*. 123–130. https://doi.org/10.1109/VAST.2008.4677365

[19] Yi Guo, Nana Cao, Xiaoyu Qi, Haoyang Li, Danqing Shi, Jing Zhang, Qing Chen, and Daniel Weiskopf. 2023. Urania: Visualizing Data Analysis Pipelines for Natural Language-Based Data Exploration. *ArXiv* (2023). https://arxiv.org/abs/2306.07760

[20] Yi Guo, Danqing Shi, Mingjuan Guo, Yanqiu Wu, Qing Chen, and Nana Cao. 2021. Talk2Data: A Natural Language Interface for Exploratory Visual Analysis via Question Decomposition. https://api.semanticscholar.org/CorpusID:258685266

[21] Marti Hearst and Melanie Tory. 2019. Would You Like A Chart With That? Incorporating Visualizations into Conversational Interfaces. In *2019 IEEE Visualization Conference (VIS)*. 1–5. https://doi.org/10.1109/VISUAL.2019.8933766

[22] E. Hoque, P. Kavehzadeh, and A. Masry. 2022. Chart Question Answering: State of the Art and Future Directions. *Computer Graphics Forum* 41, 3 (2022), 555–572. https://doi.org/10.1111/cgf.14573

[23] Enamul Hoque, Vidya Setlur, Melanie Tory, and Isaac Dykeman. 2018. Applying Pragmatics Principles for Interaction with Visual Analytics. *IEEE Transactions on Visualization and Computer Graphics* 24, 1 (2018), 309–318. https://doi.org/10.1109/TVCG.2017.2744684

[24] J. D. Hunter. 2007. Matplotlib: A 2D graphics environment. *Computing in Science & Engineering* 9, 3 (2007), 90–95. https://doi.org/10.1109/MCSE.2007.55

[25] Ecem Kavaz, Anna Puig, and Inmaculada Rodríguez. 2023. Chatbot-Based Natural Language Interfaces for Data Visualisation: A Scoping Review. *Applied Sciences* 13, 12 (Jun 2023), 7025. https://doi.org/10.3390/app13127025

[26] Mary Beth Kery, Marissa Radensky, Mahima Arya, Bonnie E. John, and Brad A. Myers. 2018. The Story in the Notebook: Exploratory Data Science Using a Literate Programming Tool. In *Proc. CHI'18*. Association for Computing Machinery, New York, NY, USA, 1–11. https://doi.org/10.1145/3173574.3173748

[27] Alicia Key, Bill Howe, Daniel Perry, and Cecilia Aragon. 2012. VizDeck: Self-Organizing Dashboards for Visual Analytics. In *Proc. SIGMOD'12*. 681–684. https://doi.org/10.1145/2213836.2213931

[28] Dae Hyun Kim, Enamul Hoque, and Maneesh Agrawala. 2020. Answering Questions about Charts and Generating Visual Explanations. In *Proc. CHI'20*. 1–13. https://doi.org/10.1145/3313831.3376467

[29] Younghoon Kim, Kanit Wongsuphasawat, Jessica Hullman, and Jeffrey Heer. 2017. GraphScape: A Model for Automated Reasoning about Visualization Similarity and Sequencing. In *Proc. CHI'17*. https://doi.org/10.1145/3025453.3025866

[30] EVERT KUIJPERS and MICHAEL WILSON. 1992. A multi-modal interface for man machine interaction with knowledge based systems-mmi 2. *Knowledge Based Systems* 2 (1992), 3.

[31] Doris Jung-Lin Lee, Vidya Setlur, Melanie Tory, Karrie Karahalios, and Aditya Parameswaran. 2022. Deconstructing Categorization in Visualization Recommendation: A Taxonomy and Comparative Study. *IEEE Transactions on Visualization and Computer Graphics* 28, 12 (2022), 4225–4239. https://doi.org/10.1109/TVCG.2021.3085751

[32] Doris Jung-Lin Lee, Dixin Tang, Kunal Agarwal, Thyne Boonmark, Caitlyn Chen, Jake Kang, Ujjaini Mukhopadhyay, Jerry Song, Micah Yong, Marti A. Hearst, and Aditya G. Parameswaran. 2021. Lux: Always-on Visualization Recommendations for Exploratory Dataframe Workflows. *Proc. VLDB'21*. 15, 3 (nov 2021), 727–738. https://doi.org/10.14778/3494124.3494151

[33] Can Liu, Yun Han, Ruike Jiang, and Xiaoru Yuan. 2021. ADVISor: Automatic Visualization Answer for Natural-Language Question on Tabular Data. In *2021 IEEE 14th Pacific Visualization Symposium (PacificVis)*. 11–20. https://doi.org/10.1109/PacificVis52677.2021.00010

[34] Jiali Liu, Nadia Boukhelifa, and James R. Eagan. 2020. Understanding the Role of Alternatives in Data Analysis Practices. *IEEE Transactions on Visualization and Computer Graphics* 26, 1 (2020), 66–76. https://doi.org/10.1109/TVCG.2019.2934593

[35] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *Comput. Surveys* 55, 9 (2023), 1–35.

[36] Robert Logan, Nelson F. Liu, Matthew E. Peters, Matt Gardner, and Sameer Singh. 2019. Barack's Wife Hillary: Using Knowledge Graphs for Fact-Aware Language Modeling. In *Proc. ACL'19*. 5962–5971. https://doi.org/10.18653/v1/P19-1598

[37] Yuyu Luo, Xuedi Qin, Nan Tang, Guoliang Li, and Xinran Wang. 2018. DeepEye: Creating Good Data Visualizations by Keyword Search. In *Proceedings of the 2018 International Conference on Man. SIGMOD'18*. 1733–1736. https://doi.org/10.1145/3183713.3193545

[38] Yuyu Luo, Nan Tang, Guoliang Li, Jiawei Tang, Chengliang Chai, and Xuedi Qin. 2022. Natural Language to Visualization by Neural Machine Translation. *IEEE Transactions on Visualization and Computer Graphics* 28, 1 (2022), 217–226. https://doi.org/10.1109/TVCG.2021.3114848

[39] Jock Mackinlay, Pat Hanrahan, and Chris Stolte. 2007. Show Me: Automatic Presentation for Visual Analysis. *IEEE Transactions on Visualization and Computer Graphics* 13, 6 (2007), 1137–1144. https://doi.org/10.1109/TVCG.2007.70594

[40] Paula Maddigan and Teo Susnjak. 2023. Chat2VIS: Generating Data Visualizations via Natural Language Using ChatGPT, Codex and GPT-3 Large Language Models. *IEEE Access* 11 (2023), 45181–45193. https://doi.org/10.1109/ACCESS.2023.3274199

[41] Michael McTear. 2020. *Conversational ai: Dialogue systems, conversational agents, and chatbots.* Morgan & Claypool.

[42] Aditi Mishra, Utkarsh Soni, Anjana Arunkumar, Jinbin Huang, Bum Chul Kwon, and Chris Bryan. 2023. PromptAid: Prompt Exploration, Perturbation, Testing and Iteration using Visual Analytics for Large Language Models. *ArXiv* abs/2304.01964 (2023). https://arxiv.org/abs/2304.01964

[43] Dominik Moritz, Chenglong Wang, Greg L. Nelson, Halden Lin, Adam M. Smith, Bill Howe, and Jeffrey Heer. 2019. Formalizing Visualization Design Knowledge as Constraints: Actionable and Extensible Models in Draco. *IEEE Transactions on Visualization and Computer Graphics* 25, 1 (2019), 438–448. https://doi.org/10.1109/TVCG.2018.2865240

[44] Aditeya Pandey, Sehi L'Yi, Qianwen Wang, Michelle A. Borkin, and Nils Gehlenborg. 2023. GenoREC: A Recommendation System for Interactive Genomics Data Visualization. *IEEE Transactions on Visualization and Computer Graphics* 29, 1 (2023), 570–580. https://doi.org/10.1109/TVCG.2022.3209407

[45] A. Pandey, A. Srinivasan, and V. Setlur. 2023. MEDLEY: Intent-based Recommendations to Support Dashboard Composition. *IEEE Transactions on Visualization & Computer Graphics* 29, 01 (jan 2023), 1135–1145. https://doi.org/10.1109/TVCG.2022.3209421

[46] Anshul Vikram Pandey, Anjali Manivannan, Oded Nov, Margaret Satterthwaite, and Enrico Bertini. 2014. The Persuasive Power of Data Visualization. *IEEE Transactions on Visualization and Computer Graphics* 20, 12 (2014), 2211–2220. https://doi.org/10.1109/TVCG.2014.2346419

[47] Gabriel Poesia, Oleksandr Polozov, Vu Le, Ashish Tiwari, Gustavo Soares, Christopher Meek, and Sumit Gulwani. 2022. Synchromesh: Reliable code generation from pre-trained language models. *ArXiv* abs/2201.11227 (2022). https://arxiv.org/abs/2201.11227

[48] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language Models are Unsupervised Multitask Learners. (2019).

[49] Eric D. Ragan, Alex Endert, Jibonananda Sanyal, and Jian Chen. 2016. Characterizing Provenance in Visualization and Data Analysis: An Organizational Framework of Provenance Types and Purposes. *IEEE Transactions on Visualization and Computer Graphics* 22, 1 (2016), 31–40. https://doi.org/10.1109/TVCG.2015.2467551

[50] Deepthi Raghunandan, Aayushi Roy, Shenzhi Shi, Niklas Elmqvist, and Leilani Battle. 2023. Code Code Evolution: Understanding How People Change Data Science Notebooks Over Time. In *Proc. CHI'23 (CHI '23).* Article 863, 12 pages. https://doi.org/10.1145/3544548.3580997

[51] Sabbir M. Rashid, James P. McCusker, Paulo Pinheiro, Marcello P. Bax, Henrique O. Santos, Jeanette A. Stingone, Amar K. Das, and Deborah L. McGuinness. 2020. The Semantic Data Dictionary – An Approach for Describing and Annotating Data. *Data Intelligence* 2, 4 (10 2020), 443–486. https://doi.org/10.1162/dint_a_00058

[52] Donghao Ren, Matthew Brehmer, Bongshin Lee, Tobias Höllerer, and Eun Kyoung Choe. 2017. ChartAccent: Annotation for data-driven storytelling. In *Proc. PacificVis'17.* 230–239. https://doi.org/10.1109/PACIFICVIS.2017.8031599

[53] Craige Roberts. 2018. Speech Acts in Discourse Context. In *New Work on Speech Acts.* Oxford University Press. https://doi.org/10.1093/oso/9780198738831.003.0012

[54] Adam Rule, Aurélien Tabard, and James D. Hollan. 2018. Exploration and Explanation in Computational Notebooks. In *Proc. CHI'18.* 1–12. https://doi.org/10.1145/3173574.3173606

[55] Vidya Setlur, Sarah E. Battersby, Melanie Tory, Rich Gossweiler, and Angel X. Chang. 2016. Eviza: A Natural Language Interface for Visual Analysis. In *Proc. UIST'16* (Tokyo, Japan). Association for Computing Machinery, New York, NY, USA, 365–377. https://doi.org/10.1145/2984511.2984588

[56] Vidya Setlur and Melanie Tory. 2022. How Do You Converse with an Analytical Chatbot? Revisiting Gricean Maxims for Designing Analytical Conversational Behavior. In *Proc. CHI'22.* Article 29, 17 pages. https://doi.org/10.1145/3491102.3501972

[57] Leixian Shen, Enya Shen, Yuyu Luo, Xiaocong Yang, Xuming Hu, Xiongshuai Zhang, Zhiwei Tai, and Jianmin Wang. 2023. Towards Natural Language Interfaces for Data Visualization: A Survey. *IEEE Transactions on Visualization and Computer Graphics* 29, 6 (2023), 3121–3144. https://doi.org/10.1109/TVCG.2022.3148007

[58] Leixian Shen, Enya Shen, Zhiwei Tai, Yiran Song, and Jianmin Wang. 2021. TaskVis: Task-oriented Visualization Recommendation. In *EuroVis 2021 - Short Papers*, Marco Agus, Christoph Garth, and Andreas Kerren (Eds.). The Eurographics Association. https://doi.org/10.2312/evs.20211061

[59] Chandan Singh, John X. Morris, Jyoti Aneja, Alexander M. Rush, and Jianfeng Gao. 2022. Explaining Patterns in Data with Language Models via Interpretable Autoprompting. *ArXiv* abs/2210.01848 (2022). https://arxiv.org/abs/2210.01848

[60] Arjun Srinivasan, Nikhila Nyapathy, Bongshin Lee, Steven M. Drucker, and John Stasko. 2021. Collecting and Characterizing Natural Language Utterances for Specifying Data Visualizations. In *Proc. CHI'21.* Article 464, 10 pages. https://doi.org/10.1145/3411764.3445400

[61] Arjun Srinivasan and Vidya Setlur. 2021. Snowy: Recommending Utterances for Conversational Visual Analysis. In *Proc. UIST'21 (UIST '21).* Association for Computing Machinery, New York, NY, USA, 864–880. https://doi.org/10.1145/3472749.3474792

[62] Arjun Srinivasan and John Stasko. 2017. Natural Language Interfaces for Data Analysis with Visualization: Considering What Has and Could Be Asked. In *Proc EuroVis'17 Short Papers.* 55–59. https://doi.org/10.2312/eurovisshort.20171133

[63] Arjun Srinivasan and John Stasko. 2018. Orko: Facilitating Multimodal Interaction for Visual Exploration and Analysis of Networks. *IEEE Transactions on Visualization and Computer Graphics* 24, 1 (2018), 511–521. https://doi.org/10.1109/TVCG.2017.2745219

[64] Melanie Tory and Vidya Setlur. 2019. Do What I Mean, Not What I Say! Design Considerations for Supporting Intent and Context in Analytical Conversation. In *2019 IEEE Conference on Visual Analytics Science and Technology (VAST).* 93–103. https://doi.org/10.1109/VAST47406.2019.8986918

[65] Jacob VanderPlas, Brian Granger, Jeffrey Heer, Dominik Moritz, Kanit Wongsuphasawat, Arvind Satyanarayan, Eitan Lees, Ilia Timofeev, Ben Welsh, and Scott Sievert. 2018. Altair: Interactive statistical visualizations for python. *Journal of open source software* 3, 32 (2018), 1057.

[66] Henrik Voigt, Ozge Alacam, Monique Meuschke, Kai Lawonn, and Sina Zarrieß. 2022. The Why and The How: A Survey on Natural Language Interaction in Visualization. In *Proc. NAACL'22*. 348–374. https://doi.org/10.18653/v1/2022.naacl-main.27

[67] Michael L. Waskom. 2021. seaborn: statistical data visualization. *Journal of Open Source Software* 6, 60 (2021), 3021. https://doi.org/10.21105/joss.03021

[68] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. https://arxiv.org/abs/2201.11903

[69] Kanit Wongsuphasawat, Dominik Moritz, Anushka Anand, Jock Mackinlay, Bill Howe, and Jeffrey Heer. 2016. Voyager: Exploratory Analysis via Faceted Browsing of Visualization Recommendations. *IEEE Transactions on Visualization and Computer Graphics* 22, 1 (2016), 649–658. https://doi.org/10.1109/TVCG.2015.2467191

[70] Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022. AI Chains: Transparent and Controllable Human-AI Interaction by Chaining Large Language Model Prompts. In *Proc. CHI'22*. Article 385, 22 pages. https://doi.org/10.1145/3491102.3517582

[71] Bowen Yu and Cláudio T. Silva. 2020. FlowSense: A Natural Language Interface for Visual Data Exploration within a Dataflow System. *IEEE Transactions on Visualization and Computer Graphics* 26, 1 (2020), 1–11. https://doi.org/10.1109/TVCG.2019.2934668

[72] Jian Zhao, Mingming Fan, and Mi Feng. 2022. ChartSeer: Interactive Steering Exploratory Visual Analysis With Machine Intelligence. *IEEE Transactions on Visualization and Computer Graphics* 28, 3 (2022), 1500–1513. https://doi.org/10.1109/TVCG.2020.3018724

[73] Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2022. Large Language Models Are Human-Level Prompt Engineers. *ArXiv* abs/2211.01910 (2022). https://api.semanticscholar.org/CorpusID:253265328

[74] Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. 2023. Universal and Transferable Adversarial Attacks on Aligned Language Models. arXiv:2307.15043 [cs.CL]

## A   ANALYTIC CONVERSATIONS WITH CHATBOTS: ADDITIONAL DETAILS

We re-analyzed the results for 21 participants pertaining to the *context-only* and *context + intent* conditions from by Tory and Setlur (2019) [64] to establish a baseline of visual analytic conversations with AI chatbots. Tory and Setlur summarize their results as a state transition model for analytic conversations and the importance of these transitions to maintain visual and conversational coherence. Compared to their analysis, we focus more on understanding the overall structure and evolution of these conversations to assess the responses AI agent that is based on a GPT language model.

### A.1   The Structure of Analytic Conversations

*A.1.1   Visualizing Discourse Chunks.* We translated these linear conversations into a branching tree structure based upon utterance markers, such as 'start over' or 'reset', that separate discourse chunks. In Figure 15 we show the structure conversational structure of four participants; additional conversation trees are in the online supplemental materials. Each branch is one discourse chunk and the start of each branch is a new analytic question. These conversations contain analytic questions with branching structures that representing refinements of the data attributes (e.g., adding or removing attributes), data operations (e.g., filtering), or visualization (e.g., changing the encoding type; binding attributes to different channels).

| Utterance | Classification | Attributes | Data Operations | Vis Specifications |
|---|---|---|---|---|
| How many people survived? | Interrogative | survived (e) | count (s) | - |
| Class colored by fare bin | Imperative | pclass, fare (e)      (e) | bin (e) | color |
| What are the top 5 age bins that survived? | Interrogative | age, survived (e)      (e) | bin, top n (e)    (e) | - |
| Show only male | Imperative | sex (v) | filter (s) | - |
| Survivors by age | Declarative | age, survived (e)      (e) | - | - |

Fig. 14.   Examples utterances according to their phrasal type and ambiguity. We use the definitions from [62] to define reference ambiguity to attributes, data operations, or encoding choices. We also extend the definition of ambiguity to include reference to a prior visualization

*A.1.2   Classifying Utterances.* Prior research by Gao *et al.* [16] and Srinivasan *et a.* [62] have proposed classifications for conversational utterances and their ambiguity. We find the utterances from this study resemble findings align with the classifications (questions, commands, queries, and others) proposed in [62], however, we recast these definitions here to align with linguistic classifications of utterances in discourse [30, 53]. We propose using *interrogative* in lieu of questions and *imperative* in lieu of commands. While less uncommon, some utterances took the form of declarative utterances (queries or other), such as *"Women in first class who survived by age"* or *"survivors by age"*, that are posed neither as a question nor a specific instruction to do something. Even though the imperative "show me" command could be implied (e.g., *"show me survivors by age"*), we classify the utterance as declarative because "show me" is not explicitly stated. Similar to prior research, we define ambiguity according to the utterance's specification of the attributes, data operations (including aggregations and filtering), and visualization specification (including chart types and encoding choices). Again, we use the classifications of Srinivasan *et al.* [60] to indicate whether attribute references are explicitly

matched all or portion of an attribute's name (e.g.,*plot class by age*), are semantically similar (e.g., *"on average how much did people pay?"*; paid is semantically similar to fare), based upon values (e.g., *"Show only male"*, male it a value in the sex attribute), or is implicit based upon the response (e.g., *"bin by 10-year intervals"* when age is the chart's x-axis).
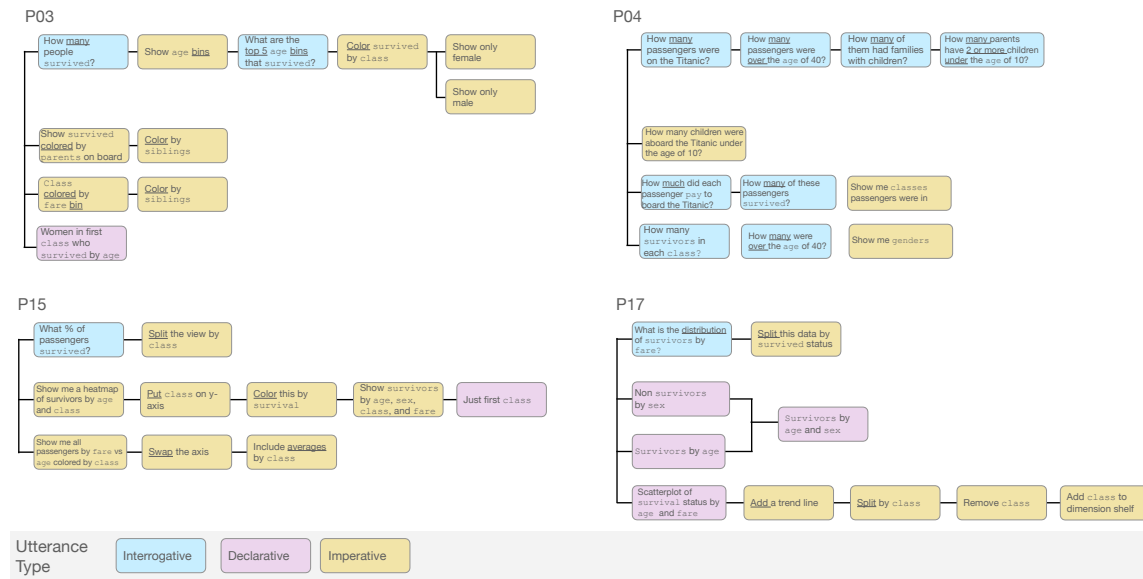


Fig. 15. Visualizations of discord chunks derived from participants in [64].

*A.1.3  The evolution of analytic conversations.*  We combined the branching structures together with the classification of individual utterances to examine the evolution of analytic conversations from this dataset. The individual approaches of participants varied, with some tending to pose more interrogative utterances throughout their sessions. However, participants typically began with interrogative utterances and then proceed imperative utterances within a branch (discourse chunk). Participants also tended to use interrogative utterances at the beginning of sessions but eventually transition to more imperative, and sometimes declarative, utterances to pose new questions (start a new branch). Generally, all questions made reference to some data attributes (often explicitly or in a semantically similar way) and some data operation. It was less common was the specify a specify chart type and more common to refine existing charts wizard-generated by binding additional attributes to available encoding channels. Within a branch, the majority of participants used unspecified imperative utterances. These utterances only make sense when taking into account the conversational context of the specific branch. In only two instances did we observe the explicit merging (P17, Figure 15) or splitting (P03) of two branches; the majority of branches went from the more general to a specific in a linear manner.

## A.2  Replication in ChatGPT Advanced Data Analysis

We compare the responses from the WoZ study to those produced by the ChatGPT Advanced Data Analysis (ADA) plugin. We identified the following key differences.

**Text vs. visualization default responses.** The WoZ study defaulted all responses to visualizations, whereas ADA tended to default to textual responses unless prompted to produce a visualization. However, ADA's behavior could be

modified by providing instructions at the beginning of the session, for example instructing it to *"provide all responses as data visualizations by default and where appropriate"*[7]. Alternatively, participants' utterances could make explicit references to creating a visualization in the imperative (e.g., *"show me a visualization of [...]"* or simply *"show me"* or interrogative forms (e.g., *"can you show me a visualization..."*) When ADA provided a textual response it was also possible to ask it to provide a visualization as a follow-up utterance. This interplay between text and visualizations is absent in the WoZ study but is something ADA is capable of supporting. While some prior research [13, 21] examined this interplay, the data from these studies is limited.

**Phrasal classification of utterances does not impact ADA's ability to respond.** Generally, we found that it did not matter whether an utterance was posed to ADA was in an interrogative form or not (e.g., *"how many people survived?"* vs. *"total number of people that survived"*). Even constructing non-interrogative sentences that make little sense in the context of language dialogue still yielded a correct response(e.g. inputting just *"many people survived"* without any context but the dataset). While out of the scope of this study, it would be interesting to explore this observation further since we see evidence that auto-regressive agents like ADA do not require utterances resemble human-to-human dialogue. That is, humans may converse with an AI agent in a manner similar to conversing with other humans, but, this is not strictly necessary for an AI Agent to respond appropriately.

What did matter more were references to attributes, data operations, and encodings in the utterance. Even the terse *"many people survived"* arguably contains references to a count operation (semantically similar to many) and a data attribute ('survived') compared to a more general and imperative utterance such as *"summarize survival"* where summarize can refer to many data operations, such as counting, computing a proportion, computing a rate, etc. On this dataset, we found that, in general, ADA correctly mapped explicit, semantic, and value-based references to attributes. It also did so for the majority of data operations requested by participants (binning, counting, averages). When it was unsure of the mapping it would prompt for further clarification, but this was not consistent behavior sometimes it also produced an error. Despite the non-deterministic behavior of GPT model outputs, it was also possible to provide a non-ambiguous imperative response (*"Scatter plot of survival status by age and fare"*). While there was still variability in ADA's overall it produced the expected output.

**Difficulties incorporating context appropriately.** ADA was not able to consistently incorporate the contextual dependency of utterances. For example, P03 in [64] has a chain of utterances that elaborate from the overall survival to modifying the chart by including age than gender (Figure 15). ADA treats these utterances as independent and does not further amend the visualization to add age and gender to the survival. The result is that instead of incremental refinements, it produces a single chart for each utterance. We can improve the performance if we add an anaphoric reference such as 'this' or 'the previous visualization' to the utterances, for example, the utterance *"how many survived"* followed by *"add age bins to this"* produces a comparable visualization to the WoZ results. We observed that after several visualizations related to one chart, it is not always necessary to include a reference term. This behavior is important as the majority of analysis refinements are ambiguous imperative statements that are dependent on the context of the discourse block.

When ADA learned to refine a chart without needing to have an anaphoric reference it became necessary to provide explicit discourse markers to begin new analysis or visual encoding. For example, if refining a chart, it was possible to provide a 'start over' essentially shifting the agent's focus to something new. However, the overall behavior of ADA

---

[7]In August a custom instructions option was provided as well that could prompt ADA to response using visualizations by default and across all conversations.

was unpredictable because it was not clear how an utterance would be processed until the output was displayed. In circumstances when the dependency of utterances was interpreted incorrectly, it was necessary to either begin the analysis a new or provide an independent imperative prompt (e.g. *"show me a visualization of survival status according to the top five age bins in male passengers"*)

Lastly, ADA can switch context to answer questions outside of the dataset. For example, there are facts about the Titanic, such as the name of its captain, the number of rescue boats, or the name of the rescue ship, that might be of interest but are not directly in the dataset. This was not tested in the WoZ study.

**Similar visualizations between WoZ and ADA.** The visualizations that ADA produced were primarily bar charts and the same was true for the WoZ. One notable difference was that the Tableau system used in the WoZ study tended to use facets, whereas ADA did not. By default ADA either created compounded variable names (e.g., women-1, indicated females that survived) or would chose to filter the data and only display the subset. Finally, ADA also provided textual descriptions and interpretations of charts that it produced, which was absent in the WoZ study.

**Errors interrupt the conversational flow.** On occasion, the WoZ study reached error states resulting from either unclear on incomplete utterances or from errors with the Tableau Software used to generate the visualization. However, these were less common. In contrast, ADA frequently produced errors when generating a visualization, once again related to utterance clarity but also because the code it generated was not executable. In this latter case, ADA would automatically attempt to regenerate an appropriate response.

## B  AI THREADS ADDITIONAL DETAILS

### B.1  Editing Actions & Utterance Intent Examples

AI Threads can modify a visual encoding in a variety of ways, in Figure 16 we provide a concrete set of examples from the edit operations described in GraphScape [29]. Additionally, we explored how AI Threads responds utterances posed with different types of analytic intents. We draw on examples utterances from the SNOWY system [61]. Figure 17 shows an example of the utterance, the resulting visualization, and a text output generated by the system to describe the data visualization that was generated. *Note that due to the size of the figures, they follow the subsequent Section B.2.*

### B.2  Few-Shot Examples

The following are the few shot examples provided to the GPT-3.5 ChatCompletion API via AI Threads as conversational context following the system instructional prompt (Section 6.2).

During pilot studies, we noticed that GPT-3.5 tended to hallucinate semantically relevant yet non-existent colors or color palettes when generating or modifying visualizations from color-related user utterances, unless it learned how to specify colors in the demonstrated approach. Defining every variable in the examples, even the chart's `width` and `height`, before they are used also allowed the model to learn how to correctly modify charts. The model is also taught how to incorporate computations to its captions.

**Example 1:**

[ROLE: USER]

```
How many rows are in the dataset?
```

[ROLE: ASSISTANT]

```python
```

```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import base64
import io

tab10 = ["#4e79a7","#f28e2c","#e15759","#76b7b2","#59a14f","#edc949","#af7aa1","#ff9da7","#9c755f",
"#bab0ab"]
sns.set_palette(tab10)
sns.set(style="whitegrid")

# Load the dataset
df = pd.read_csv('./workspace/{filename}')

# Get the number of rows in the DataFrame
num_rows = df.shape[0]

# Create the bar chart using seaborn, adjusting the size of the chart as needed
width = 8
height = 2
plt.figure(figsize=(width, height))
sns.barplot(x=[num_rows])

# Add labels, annotations, and/or title
plt.xlabel('Number_of_Rows')

# Display the chart as a Bytes64 PNG representation with a tight layout
plt.tight_layout()
buffer = io.BytesIO()
plt.savefig(buffer, format='png')
buffer.seek(0)
png_data = base64.b64encode(buffer.getvalue()).decode('utf-8')
buffer.close()

plt.clf()
plt.close()

# Print the chart image
print(png_data)

# Print a caption
print("This_plot_shows_a_single_bar_representing_the_total_number_of_rows_in_the_dataset,_which_is_
" + str(num_rows) + ".")
```
```

**Example 2:**

[ROLE: USER]

```
Plot the distributions of every continuous measure.
```

[ROLE: ASSISTANT]

```python
import pandas as pd
import seaborn as sns
```

```python
import matplotlib.pyplot as plt
import base64
import io

tab10 = ["#4e79a7","#f28e2c","#e15759","#76b7b2","#59a14f","#edc949","#af7aa1","#ff9da7","#9c755f",
"#bab0ab"]
sns.set_palette(tab10)
sns.set(style="white")

# Load the dataset
df = pd.read_csv('./workspace/{filename}')

# Selecting the continuous variables from the DataFrame
continuous_vars = []

for col_name, dtype in df.dtypes.iteritems():
    if dtype==float:
        continuous_vars.append(col_name)

# Creating a new DataFrame with only the continuous variables
continuous_df = df[continuous_vars]

# Melt the DataFrame to create a long-form representation
melted_df = pd.melt(continuous_df)

# Plotting a faceted histogram
g = sns.FacetGrid(melted_df, col="variable", col_wrap=2, sharex=False, sharey=False)
num_bins = 20
g.map(sns.histplot, "value", bins=num_bins)
g.set_titles("{col_name}")
g.set_axis_labels("", "Count")

# Display the chart as a Bytes64 PNG representation with a tight layout
plt.tight_layout()
buffer = io.BytesIO()
plt.savefig(buffer, format='png')
buffer.seek(0)
png_data = base64.b64encode(buffer.getvalue()).decode('utf-8')
buffer.close()

plt.clf()
plt.close()

# Print the chart image
print(png_data)
# Print a caption
print("This plot shows a grid of histograms, each representing a different continuous variable. The
 x-axis represents the value range of each variable, and the y-axis represents the count of
occurrences.")
```

| Category | Edit Operation | Example Initial | Revised |
|---|---|---|---|

**Mark** — Mark A -> Mark B

Change the mark type from A to B, or vice versa

**Utterance** : Show me a scatter plot to two highly correlated and numeric variables in the dataset

**Utterance** : Make this a box plot



**Transform** — Scale

Add, change or remove an axis scale transform

**Utterance** : Make the y-axis a log scale

Sort

Change the sort order of a visualized data field.

**Utterance** : sort the x-axis from the lowest average fare paid to the highest average fare paid

Bin

Discretize values of a visualized field into bins.

**Utterance** : Show me a scatter plot between fare paid and the passenger age

**Utterance** : can you create density bins to summarize this?

Aggregate

Aggregate values according to an aggregation function such as sum, mean or median.

**Utterance** : Can you add a line to the y-axis that shows the median fare paid?

Modify / Filter/ Add / Remove

Filter data records according to a filter predicate function.

**Utterance**: can you filter out passengers that paid more than $150

**Encoding** — Transpose

Swaps the (x, y) or (row, column) channels.

**Utterance**: can you swap the x and y axes?

Add / Remove

Add or remove a field from an encoding channel

**Utterance** : Show me a scatter plot between fare paid and the passenger age

**Utterance**: Can you add passenger class?

Move

Move a field on one channel to another channel.

**Utterance**: **Utterance**: can you facet by passenger class instead of using color?

Modify

Replace a field in a channel with another field.

**Utterance**: can you color points by survival status?

Fig. 16. **Example editing operations.** The arrows indicate that the editing operation was continues in a thread. The above figure show three prompts producing a encoding that is the further modified via editing operations in a thread. Arrows indicate a continuation from the previously refined encoding. Edit operations were taken from [29].

| Intent | Example | | |
|--------|---------|---|---|
| | **Utterance** | Visualization | Output Text |
| Correlate | How are fare and passenger's age correlated? |  | This scatter plot shows the relationship between Fare and Age of the passengers. The x-axis represents the Fare amount paid by the passengers, while the y-axis represents their age. The dots represent individual passengers. The correlation coefficient between Fare and Age is 0.09606669176903888, indicating the strength and direction of the linear relationship between the two variables. |
| Group | On average what was the total fare paid by class? |  | This bar plot shows the average fare paid by passenger class. The x-axis represents the passenger class (1st, 2nd, 3rd), and the y-axis represents the average fare. The height of each bar represents the average fare paid by passengers in that class. It can be observed that on average, passengers in 1st class paid the highest fare, followed by 2nd class and then 3rd class. |
| Aggregate | Show median fare paid instead |  | This bar plot shows the median fare paid by passenger class. The x-axis represents the passenger class (1st, 2nd, 3rd), and the y-axis represents the median fare. The height of each bar represents the median fare paid by passengers in that class. It can be observed that the median fare paid increases from 3rd class to 1st class, indicating a higher fare for higher passenger classes. |
| Trend | Show the survival by 10 year age groups |  | This bar plot shows the survival rate for each ten-year age group. The x-axis represents the ten-year age groups, and the y-axis represents the survival rate. The height of each bar represents the survival rate of passengers in that age group. It can be observed that the survival rate is highest for children (0 - 9 years) and decreases as the age group increases until it reaches the lowest for the elderly (70 - 79 years). |
| Distribution | What was the age age distribution by survival status? |  | This histogram shows the distribution of age by survival status. The x-axis represents the age of the passengers, and the y-axis represents the count of passengers. The histogram is stacked, with the different survival statuses (0 = Not Survived, 1 = Survived) represented by different colors. The vertical axis also includes a smoothed density estimate (Kernel Density Estimate) to visualize the underlying distribution. It can be observed that there is a higher count of young passengers who did not survive compared to those who survived. |
| Filter | Show the top three age groups by survival status |  | This bar plot shows the count of survivors and non-survivors in the top three age groups. The x-axis represents the age group, and the y-axis represents the count of passengers. The bars are stacked, with different colors representing the survival status (0 = Not Survived, 1 = Survived) within each age group. It can be observed that the highest count of survivors is in the age group 20-29, followed by the age group 30-39, while the highest count of non-survivors is in the age group 20-29, followed by the age group 30-39. |

Fig. 17. **Examples of AI Threads responds to utterances with different intents**. Utterances are inspired by those suggested in [61].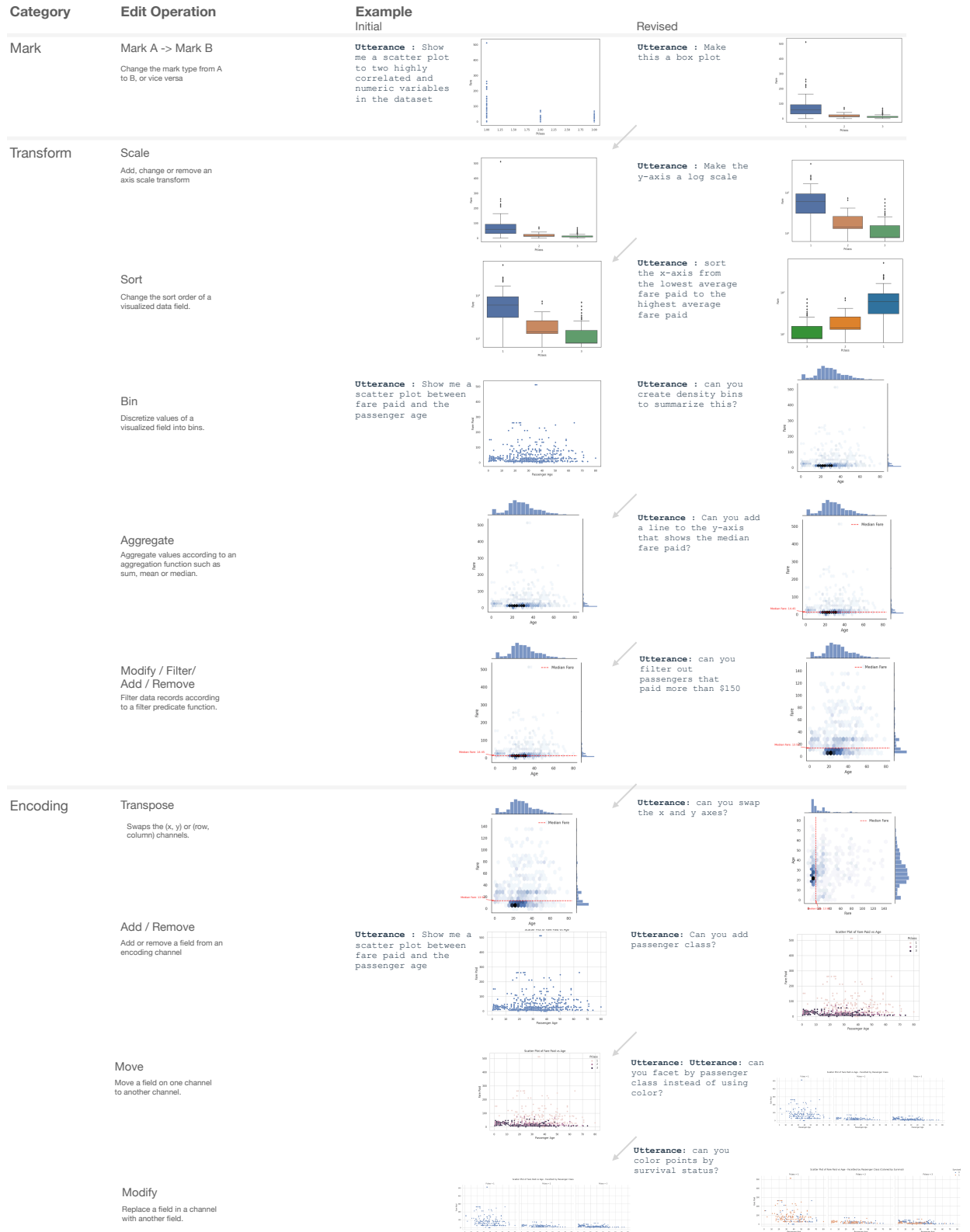