

Choice Over Control: How Users Write with Large Language Models using Diegetic and Non-Diegetic Prompting

Hai Dang

hai.dang@uni-bayreuth.de
University of Bayreuth
Bayreuth, Bavaria, Germany

Florian Lehmann

florian.lehmann@uni-bayreuth.de
University of Bayreuth
Bayreuth, Bavaria, Germany

Sven Goller

sven.goller@uni-bayreuth.de
University of Bayreuth
Bayreuth, Bavaria, Germany

Daniel Buschek

daniel.buschek@uni-bayreuth.de
University of Bayreuth
Bayreuth, Bavaria, Germany

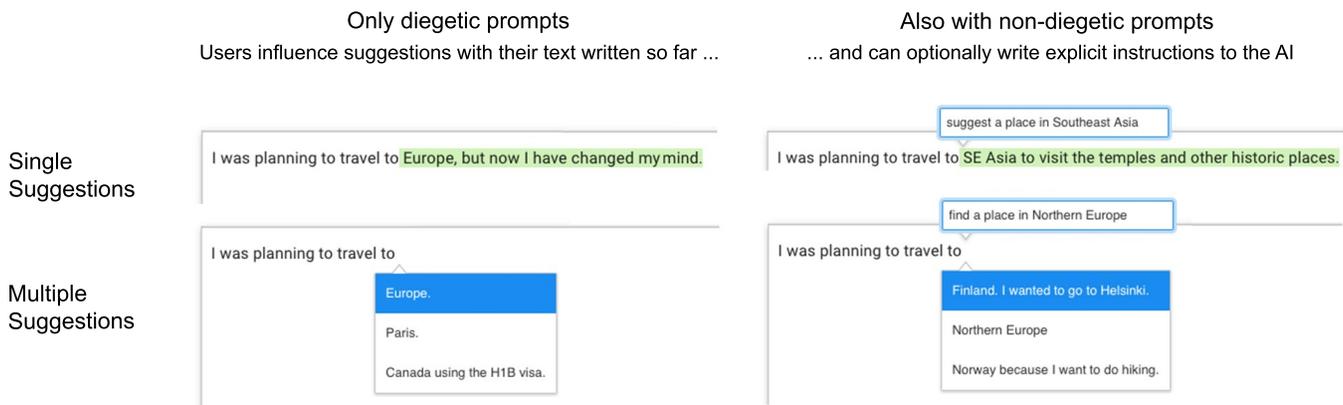


Figure 1: Overview of our four UI variants, showing the user’s written text (black font, i.e. a diegetic prompt), the suggestions (text highlighted in green, and options in the list), and a popup text box that allows users to input an instruction as a zero-shot prompt to the system (i.e. a non-diegetic prompt).

ABSTRACT

We propose a conceptual perspective on prompts for Large Language Models (LLMs) that distinguishes between (1) diegetic prompts (part of the narrative, e.g. “*Once upon a time, I saw a fox ...*”), and (2) non-diegetic prompts (external, e.g. “*Write about the adventures of the fox.*”). With this lens, we study how 129 crowd workers on *Prolific* write short texts with different user interfaces (1 vs 3 suggestions, with/out non-diegetic prompts; implemented with *GPT-3*): When the interface offered multiple suggestions and provided an option for non-diegetic prompting, participants preferred choosing from multiple suggestions over controlling them via non-diegetic prompts. When participants provided non-diegetic prompts it was to ask for inspiration, topics or facts. Single suggestions in particular were guided both with diegetic and non-diegetic information.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
CHI '23, April 23–April 28, 2023, Hamburg, Germany

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-9421-5/23/04...\$15.00
<https://doi.org/10.1145/3544548.3580969>

This work informs human-AI interaction with generative models by revealing that (1) writing non-diegetic prompts requires effort, (2) people combine diegetic and non-diegetic prompting, and (3) they use their draft (i.e. diegetic information) and suggestion timing to strategically guide LLMs.

CCS CONCEPTS

• **Human-centered computing** → **Empirical studies in HCI**; **Text input**; • **Computing methodologies** → **Natural language generation**.

KEYWORDS

Large language models, Co-creative systems, Human-AI collaboration, User-centric natural language generation

ACM Reference Format:

Hai Dang, Sven Goller, Florian Lehmann, and Daniel Buschek. 2023. Choice Over Control: How Users Write with Large Language Models using Diegetic and Non-Diegetic Prompting. In *CHI '23: ACM Conference on Human Factors in Computing Systems, April 23–April 28, 2023, Hamburg, Germany*. ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/3544548.3580969>

1 INTRODUCTION

When writing collaboratively, people coordinate and inspire each other through what they write in the draft itself and through communication beyond it. In this paper, we examine related mechanisms for human-AI co-writing.

Input text provided to a Large Language Model (LLM) as a basis for generating text is referred to as a “prompt”. Providing a few examples of inputs and outputs in such a text prompt can help the model solve a task [5, 50]. This is called few-shot learning. For example, an LLM can be prompted to translate from English to French with a few examples of English sentences and corresponding translations, followed by the English sentence to be translated. By completing this text the LLM then (ideally) translates that sentence. This affords user control: Users can define tasks and delegate them to an LLM ad-hoc. Going further, zero-shot learning prompts the LLM with an instruction *without* examples (e.g. *Translate ‘The weather is nice’ to French*). This is a harder task but from a Human-Computer Interaction (HCI) point of view it frees users from thinking of specific examples when instructing the AI system.

We introduce the terms diegetic prompting and non-diegetic prompting¹ to frame a new perspective on how users influence an LLM in their writing process. A diegetic prompt is part of the users’ narrative. For example, when the user writes about a vacation in South East Asia, the story as written so far forms the diegetic prompt. In contrast, a non-diegetic prompt is an explicit instruction to the LLM (e.g. “suggest activities to do in Singapore”). Crucially, this instruction is not a part of the resulting document (e.g. travel blog); it only serves to guide the LLM’s text generation.

Technically, there may not be a difference between diegetic and non-diegetic prompts for the LLM – both types are received by the model as text input strings. However, from an HCI perspective, this distinction allows us to identify patterns in the perception and interaction of users writing with LLMs. With this new distinction, in this paper we address the research question: *How do users write with Large Language Models using diegetic and non-diegetic prompting?*

Concretely, we propose and compare four UI variants (Figure 1) that allow people to write with these types of prompts, plus a baseline UI without suggestions. We conducted a remote study with 129 crowd workers on Prolific, each writing five stories. We investigate the influence of two independent variables on users’ writing behavior, namely INSTRUCTION with two levels (i_{no} , i_{yes}) and NUMBER of suggestions with three levels (baseline: s_0 , s_1 , s_3).

Users overall prefer choosing from multiple suggestions over controlling them via non-diegetic prompts. They use non-diegetic prompts to ask the LLM for inspiration, topics or facts. Non-diegetic prompts increase effort, for learning how to formulate them and switching between diegetic and non-diegetic writing. Users also prefer the UI with multiple suggestions over seeing single ones, yet allowing them to provide non-diegetic prompts reduces the gap in acceptance rates by boosting it for the UI with single suggestions. Moreover, single suggestions are triggered later in sentences, and less frequently at transition words and to start sentences. Together with people’s comments, this indicates that writers consider diegetic information to guide LLMs. We discuss implications for LLMs and interaction design.

We contribute a new conceptual lens on prompting that distinguishes diegetic and non-diegetic ways in which users can influence LLMs, and a new UI design to combine text continuation suggestions with zero-shot prompt input.

2 RELATED WORK

We relate our work to prompting in Natural Language Processing (NLP) and writing interfaces in Human-Computer Interaction (HCI). Moreover, we present our proposed concept of diegetic and non-diegetic prompting by locating it in existing user interfaces for writing and prompting.

2.1 Prompting in Large Language Models

Language models are trained to predict the next word given the previous words in the text. One primary advantage of Deep Learning-based LLMs is that they can solve several natural language processing tasks without being specifically trained on those. This can be done via text prompts written in natural language [5]. Zhao et al. [50] show that providing a few examples of inputs and outputs can help to steer the model. However, optimizing prompts is not trivial and requires extensive experience [28].

2.1.1 Prompt Engineering. Related work in prompt engineering has proposed several methods to improve prompts: For example, paraphrasing prompts can lead to better model outputs [20, 23, 27, 49]. Another approach involves constructing prompt templates to increase the accuracy for probing knowledge [33], for translation tasks [5], or for text classification tasks [39]. However, optimized prompts constructed in the process of prompt engineering are usually not meant to be consumed by humans; rather, they are designed for LLMs to most effectively perform a task [28]. In contrast, in our study, we explore how non-expert users write and use (zero-shot) prompts when writing with an LLM.

2.1.2 Prompting Interfaces. Several interactive systems have been proposed to enable users to work more effectively with prompts: For example, *AI Chains* by Wu et al. [47] allows users to combine multiple prompt primitives and their outputs to form a chain of prompts that can solve complex language processing tasks. In another study, they introduce an interface for visually programming these chains [46]. Similarly, *PromptMaker* [22] allows users to prototype new AI functionalities using language prompts. Strobel et al. [42] developed a prompt programming environment to allow users to experiment with prompt variations and visualize prompt performance. *Story Centaur* by Swanson et al. [43] supports users in creating few-shot examples for creative writing. Using our terminology, these projects focused on *non-diegetic* prompts as a main output of interaction. In contrast, we integrate non-diegetic prompts into a text editor, with a focus on writing. Concretely, we combine a UI for phrase suggestions with a UI for zero-shot prompt inputs to an LLM, and analyze how users make use of these during their writing.

2.2 Writing Interfaces for LLMs

Here we give a brief overview of key design factors for user interfaces that involve LLMs and text generation.

¹<https://www.merriam-webster.com/dictionary/diegetic>, last accessed March 7, 2023

2.2.1 Scope of Suggestions. Earlier work mainly focused on single word suggestions [11, 12, 18, 34]. This scope favours performance metrics, such as reducing key-strokes, while longer phrase suggestions [6, 26, 37] are perceived more as new ideas for writing [1]. We focus on such phrase suggestions in this paper.

2.2.2 Display of Suggestions. Single text suggestions can be shown inline [4, 8, 17, 48], whereas multiple suggestions are shown as pop-up lists of about three to six entries [6, 26]. Beyond that, Singh et al. [41] evaluated how writers use suggestions displayed as images and sound. Moreover, Bhat et al. [4] used a pop-up text box to show suggestions for insertions in the middle of sentences. We follow these design choices (Figure 1) and show single suggestions inline and multiple ones in a pop-up list. We add a pop-up text field for entering non-diegetic prompts.

2.2.3 Implicit vs. Explicit Trigger. In writing interfaces, suggestions can be triggered explicitly or implicitly. Related work showed suggestions automatically after short inactivity [4, 6] or gated by a utility function [24]. Alternatively, recent work has also explored designs in which users explicitly request suggestions with a hotkey [7, 17, 26, 41, 48]. We also use this design with an explicit request key to better understand how and when users request suggestions.

2.3 Diegetic and Non-Diegetic Prompting in Existing Writing Interfaces

Here we apply the proposed lens to analyse how existing systems use diegetic and non-diegetic information in their writing interfaces. Traditionally, systems mainly use diegetic information, that is, they predict text based (only) on the preceding text [11, 12, 18, 34]. Some also added other information (e.g. hand posture, body movement [15, 16]). These show early examples of non-diegetic input to the language model. In this work, we focus on textual diegetic and non-diegetic information.

From a technical perspective, for recent systems that use LLMs to generate text suggestions, there might be no difference between the user’s text draft (i.e. diegetic text) and other text inputs to the language model (e.g. instructions to the model, i.e. non-diegetic text). Therefore, systems in which the UI did not afford text prompts explicitly made the implicit choice of only using diegetic information as their input to the LLM [6, 7, 26, 41].

In contrast, writing interfaces that indeed allow users to explicitly enter prompts often use a mix of diegetic and non-diegetic information. Gero et al. [14] propose “sparks”, i.e. sentences generated from LLMs to inspire new ideas for scientific writing. The user-provided prompts to generate these sparks are not part of the final outcome text, thus they are non-diegetic. Similarly, other systems (e.g. *Wordcraft* [48], *LaMPPost* [17]) allowed users to select a part of the written text and modify it via predefined functionality (internally these functionalities also use prompting: e.g. a button for “rewrite selection” + text entry field for prompt). The selected text in this example is diegetic information while the prompt template and user-provided prompts are non-diegetic information. Related, we include the entire user written text draft as diegetic information and allow users to provide non-diegetic custom text prompts to further guide the LLM.

3 INTERACTION CONCEPT

Here we describe our UI and interaction concept (also see Figure 1 and Figure 8): It closely integrates diegetic and non-diegetic prompting in the same UI; users can use both types without having to take the hands off the keyboard.

3.1 Inline (Single) Suggestions (Figure 1 top row)

When a user requests a new suggestion (**TAB**) a preview of the suggestion appears after the current caret position in the text editor. Users can press **TAB** repeatedly to get new suggestions. The suggestion preview is visually highlighted in green to indicate that it is not part of the text yet. We decided for this design instead of e.g. a greyed out suggestion text (as e.g. used in Google’s Smart Compose [8]) because pilot tests showed that grey text can be difficult to read for some people and makes readability more dependant on screen brightness settings, which we cannot control in an online study. If the suggestion is accepted (**ENTER**) the preview style (green background) is removed and the suggested text becomes part of the text document. Alternatively, the user can cancel the current suggestion preview by pressing **ESC** or by continuing to type without confirming the suggestion. When the suggestion is cancelled in one of these ways, the previewed suggestion is removed from the text editor.

3.2 Multiple Suggestions (Figure 1 bottom row)

Our system follows current practices for multiple suggestions (see Section 2.2) and shows each phrase suggestion as a separate item in a list of three. Again, users can press **TAB** to get suggestions (and repeatedly to get new ones). Users can use the **↑ UP / ↓ DOWN** keys to navigate this list and confirm a suggestion with **ENTER**. Selection via mouse is also possible.

3.3 Pop-up Textbox for Non-Diegetic Prompts (Figure 1 right column)

In the study, we described non-diegetic prompts as “instructions to the AI”. Users can request suggestions with **TAB** as before. Additionally, they can enter an instruction by typing in a popup box that appears above the caret position. Thus, users have the option to input an instruction but are not forced to do so to request suggestions. Input focus is automatically switched from the text editor to the pop-up textbox when requesting suggestions so users can type instructions directly after pressing **TAB**. Users can submit the instruction with **TAB** or **Enter**. They can then press **Enter** again to accept the selected suggestion. Alternatively, they can revise their instruction to update the suggestions.

4 PROTOTYPE IMPLEMENTATION

Here we provide details about the web prototype used in the study. For screenshots, see Figure 1 and Appendix A.

4.1 Web System

The prototype was implemented with ReactJS² and CKEditor³. Each suggestion request from the client was passed to and parsed by a backend server which used FastApi⁴ as a lightweight webserver. The server forwarded these requests to OpenAI's *text-davinci-edit-001* model along with the entire written text as well as an instruction for the suggestion model (see Section 3). We chose this model and API because it is reportedly trained specifically to take in a given text as well as a (separate) instruction relating to the text.

4.2 Language Model Prompts

We used two default prompt prefixes to retrieve sentence completions from GPT-3 (*text-davinci-edit-001*): (1) *Complete the sentence*. (2) *Complete the sentence and <user_instruction>*. The system automatically used (1) when there was no option for the participants to provide explicit instructions to the AI, or when users did not provide an instruction. When they did write instructions, these were appended to (2). For instance, if the user wrote the instruction: 'suggest colors', the resulting full instruction sent to the model was: 'Complete the sentence and suggest colors'. During the pre-study we experimented with other default instructions such as: 'Continue' or 'Continue the text', as well as more complex ones, but found them to be less suitable (e.g. produced longer text or less consistent). We applied a post processing step to trim the model's output and display only the generated continuation.

4.3 Information Box

For the user study, we implemented an information box (Figure 8 in Appendix A) which explains the different features of the current text editor setup. Concretely, it showed an image that demonstrates the usage of the UI as well as an explanation of the available action keys.

5 METHOD

We used the following methods, in line with related studies on human-AI writing (e.g. cf. [6, 26]).

5.1 Questionnaires

To assess participants' backgrounds, an initial questionnaire asked about demographics and experience with writing features and language models. Participants also filled in one questionnaire after each UI variant (see Figure 3) to give subjective feedback per UI. To extend on this with overall feedback, a final questionnaire asked for (optional) open comments on changes to the system and experiences with suggestions and instructions.

5.2 Interaction Logging

To analyze interaction behaviour in detail, we logged interaction events, i.e. key and mouse events, during the writing tasks (see Appendix A). Each event included a *timestamp*, *task id*, and the *current text* in the editor. Depending on the event it included information

such as the suggestion trigger position in the text or the instruction to the AI.

5.3 Coding of Open Questions

We analysed the open comments from the final questionnaire in an approach adopting coding steps from Grounded Theory [10, 29], in order to identify and report on the emerging aspects: First, two researchers inductively proposed codes for the data of 20 people. They then compared and clustered these codes to develop a common codebook. Then, they coded the first 20 plus 32 more participants and checked each other's codings, with slight adjustments to the codebook. Finally, one researcher coded the remaining data and another one checked this coding. Throughout the process, disagreements were resolved via discussion.

5.4 Evaluation of User Written Text

We used LanguageTool⁵, a multilingual grammar and spell-checker, to count the number of grammar and spelling mistakes. To evaluate the degree to which participants engaged with the selected writing prompts during the user study (cf. Section 6), three researchers independently reviewed the stories and provided comments on their connection to the prompts. Finally, one researcher reviewed all comments to ensure consistency.

6 USER STUDY

6.1 Study Design

Our study uses a within-subject design with two independent variables: The NUMBER of (parallel) suggestions with two levels: one and three suggestions (s_1, s_3); and the opportunity for INSTRUCTION with two levels (i_{no}, i_{yes}). This results in four UI variants with suggestions. In addition, we included a baseline UI without any suggestions (s_0). The order of these five UIs was fully counterbalanced. As dependent variables we included interaction measures as well as questionnaire data.

6.2 Participants

We conducted a pre-study with 6 participants with direct discussions for rich feedback, followed by our main study with 129 participants ($M=71, F=57, NB=1$). We recruited on *Prolific*⁶ and screened participants for written and spoken fluency in English, as well as access to a computer with a keyboard. Participants reported ages ranged from 18 to 70 with a median age of 32. Following the platform recommendations, participants were compensated with 8 £/h.

6.3 Procedure

The study started with a description page, including information about the collected data and GDPR, in line with our institute's regulations. After giving their consent, participants were directed to a page with an overview of the procedure and involved UI variants. Following this, people were guided through the five writing tasks in a counterbalanced order, and then to the final questionnaire. The study had an estimated duration of 45 minutes (actual mean was 45 minutes and 41 seconds).

²<https://reactjs.org/>, last accessed March 7, 2023

³<https://ckeditor.com>, last accessed March 7, 2023

⁴<https://fastapi.tiangolo.com>, last accessed March 7, 2023

⁵<https://languagetool.org>, last accessed March 7, 2023

⁶<https://www.prolific.co/>, last accessed March 7, 2023

6.3.1 Topic Selection. For each task, participants first selected a writing topic. Repeated selections were allowed, as in related work [26], yet we asked them to choose at least two different topics overall. The topic order was also randomized and shown one at a time to encourage variety in the topic choices overall.

Gero et al. [13] suggested three tasks for writing support tools, including story writing and argumentative essay writing. We thus selected five topics for creative writing⁷ and five topics for argumentative writing from the same source⁸ as Lee et al. [26].

6.3.2 Writing Task. Participants were told to write about the previously selected topic for five minutes and finish their text with a clear ending. The description encouraged to try out all features but also to write their own text. A timer was shown below the text editor. It was mentioned that the timer was not a hard cut-off, but served as a reminder of when to finish the task. We set a minimum time of 15 seconds before participants could submit their story but people stayed close to the five minutes anyway (see Section 7.2). People filled in a questionnaire after each task (Section 5.1).

7 RESULTS

Here we present our study results. For statistical testing we use R [35], concretely, (generalised) linear mixed-effects models (LMMs with the packages *lme4* [2], *lmerTest* [25]). The models account for participants' individual differences, as well as for the type of their chosen topics (creative story writing, argumentative writing), via random intercepts. As fixed effects, the models have `INSTRUCTION` and `NUMBER`. Moreover, we use the R package *multgee* [44] to analyse the Likert results (i.e. ordinal data) with Generalized Estimating Equations (GEEs). We report significance at $p < 0.05$.

We define a *suggestion session* as continuous interaction with suggestions, from requesting them until cancellation or acceptance (e.g. a session might involve three subsequent “tab” presses to browse suggestions). Participants triggered 3097 suggestion sessions. The mean in tasks with suggestions enabled was 6.47 (SD 4.00), comparable to related work [26].

7.1 Suggestion Acceptance

We define the acceptance rate as the number of accepted suggestions divided by the number of triggered suggestion sessions. We found considerable differences between the UIs (Means: $s_1=0.55$, $s_3=0.74$, $i_{no}=0.59$, $i_{yes}=0.69$). The grand mean acceptance rate was 0.64 (SD 0.29). The mean for suggestion requests with a written instruction was in line with this (0.64, SD: 0.33). We fitted a generalised LMM on the acceptances as binomial data (i.e. for each shown suggestion we logged if it was accepted or not), summarized in Table 1 (row 1). Figure 2 (top left) shows the descriptive data. In summary, showing one suggestion (instead of three) significantly decreased the chance of acceptance, yet enabling users to write instructions significantly reduced this gap by increasing their acceptance (Mean rate of 0.45 for s_1 without instructions vs 0.65 with them).

⁷<https://www.reddit.com/r/WritingPrompts>, last accessed March 7, 2023

⁸<https://www.nytimes.com/2021/02/01/learning/300-questions-and-images-to-inspire-argument-writing.html>, last accessed March 7, 2023

7.2 Task Completion Time

We measured task time from starting the task to submitting it (Means: $s_0=281$, $s_1=305$, $s_3=306$, $i_{no}=303$, $i_{yes}=309$). As a fixed writing time was given, we do not expect large differences here. Indeed, an LMM fitted on this data for the suggestion UIs did not reveal significant effects (Table 1, row 2). Another such model compared the suggestion UIs against the baseline (Table 1, row 3): Here we found that writing with three suggestions took significantly longer than without suggestions. This is in line with the descriptive picture in Figure 2 (top center): Participants followed the task description of writing for five minutes, and writing with the suggestion UIs took slightly longer.

7.3 Text Length

In total, submitted texts contained 87,640 words, including text from accepted suggestions. The grand mean number of words per text was 134 words (SD 54). We fitted a generalised (Poisson) LMM on the word count data to compare the four tasks with suggestions (Table 1, row 4), and another such model to compare the suggestion UIs against the baseline without suggestions (Table 1, row 5). The results match the descriptive pattern visible in Figure 2 (top right): In summary, texts are significantly shorter when writing with single suggestions or with a UI allowing for instructions. However, writing with multiple suggestions leads to significantly longer texts. These differences are rather small, about 6-10 words (Means: $s_0=136$, $s_1=130$, $s_3=140$, $i_{no}=138$, $i_{yes}=131$).

7.4 Moments of Suggestion Requests

We analysed at which moments participants requested suggestions.

7.4.1 After Sentence vs Mid-sentence. We analysed how often suggestions started a sentence (e.g. “Hello, world! [tab]”) vs in the middle (e.g. “Hello world, how [tab]”). We fitted a generalised LMM on the requests as binomial data (i.e. for each request we logged if it was at the beginning of a new sentence or not), summarized in Table 1 (row 6). Figure 2 (bottom left) shows the descriptive data. In summary, showing one suggestion (instead of three) significantly decreased the chance of requesting suggestions at the beginning of a new sentence (Means: $s_1=21.70\%$, $s_3=31.02\%$).

7.4.2 Number of Words in Sentence. For the suggestion requests in the middle of sentences we further analysed after how many words in that sentence they were requested. We fitted an LMM on the mean numbers of words in sentences with suggestion requests per text, summarized in Table 1 (row 7). Figure 2 (bottom center) shows the descriptive data. In summary, showing one suggestion (instead of three) significantly increased the number of words in a sentence after which suggestions were requested – by about 1.5 words (Means: $s_1=10.93$, $s_3=9.48$; i.e. a relative increase of 15.3%), while `INSTRUCTION` seemed to make no difference (Means: $i_{no}=10.18$, $i_{yes}=10.34$). Note that 1-2 words later in a sentence is considerable because it may lead to very different constraints that users give to the system for possible continuations (e.g. “The...” vs “The man said...”).

7.4.3 Words at the Suggestion Requests. We further analysed the type of words after which suggestions were requested. Concretely,

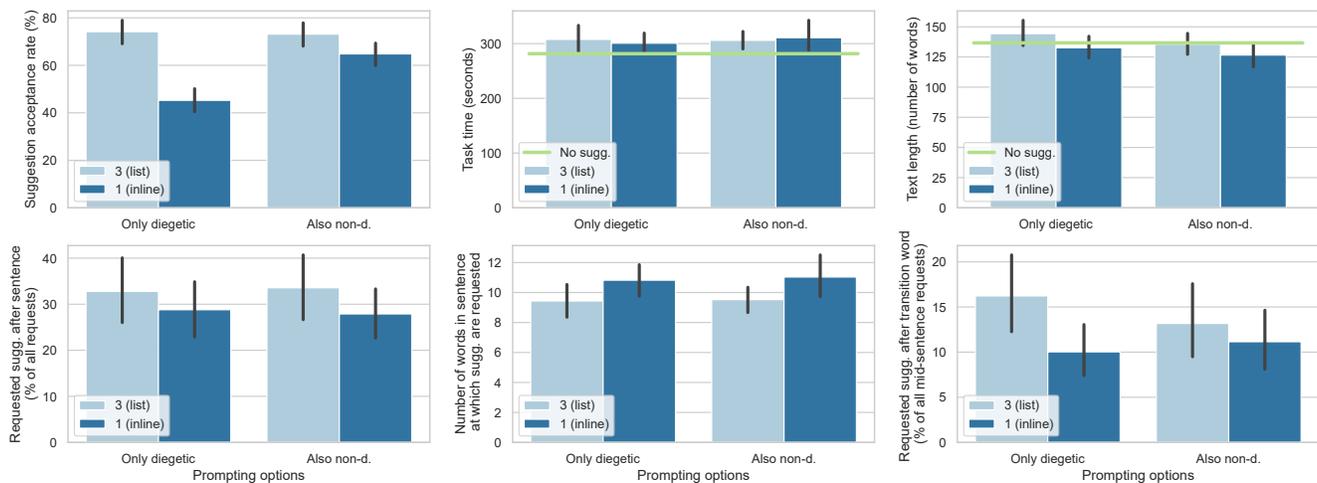


Figure 2: Overview of the interaction metrics in our study. In summary, we observe: (1) Giving users the option to write instructions (i.e. non-diegetic prompts) increases the acceptance rate of suggestions for single suggestions, but not beyond that of multiple suggestions (top left). (2) Writing time did not vary much and texts were slightly shorter with single suggestions and instructions (top center/right). (3) Single suggestions were requested less often at the start of sentences (bottom left), about 1.5 words later in a sentence (bottom center), and less often after transition words (bottom right). See text for details.

we categorised these “trigger words” into *transition* words and other words, using online lists of English transition words⁹. For example, transition words mark causes (e.g. “because”, “since”), opposites (e.g. “while”, “despite”), effects (e.g. “therefore”, “then”), and other aspects. We provide the full list we used in the project repository. We fitted a generalised LMM on the requests as binomial data (i.e. for each request we logged if it was after a transition word or not), summarized in Table 1 (row 8). Figure 2 (bottom right) shows the descriptive data. In summary, showing one suggestion (instead of three) significantly decreased the chance of requesting suggestions after a transition word (Means: $s_1=11.80\%$, $s_3=14.03\%$), while INSTRUCTION seemed to make no (sig.) difference (Means: $i_{no}=12.03\%$, $i_{yes}=13.32\%$).

7.5 Perception of the Tasks and UIs

We used Likert items to assess participants’ perception after each writing task (Figure 3). Descriptively, suggestions received favourable ratings by the majority and had almost no perceived grammatical or factual errors. However, no UI was clearly “best” for everyone: Across questions and UI variants, there is a spread of opinions, including for the perceived usefulness of being able to write instructions (i.e. non-diegetic prompting). This spread fits to the different pros and cons and preferences that participants commented on (see Section 8). Here, we report on the results from our GEE analysis. Since we have 16 questions, we summarize this analysis according to the emerging bigger picture.

7.5.1 Perceived Differences for NUMBER of Suggestions (s_1 vs s_3). Showing a single suggestion was rated worse than having a list of three suggestions. This was significant for several questions. The

GEE model estimates that the odds of giving a higher rating with a single suggestion were “ x ” times the odds of that with the list of three suggestions, with x as follows: Single suggestions were rated as significantly more distracting ($x=1.95$, $p<0.005$), less helpful ($x=0.60$, $p=0.02$), leading to more manual editing ($x=1.91$, $p<0.005$), feeling less in control ($x=0.67$, $p=0.03$), and providing less diverse suggestions ($x=0.67$, $p=0.04$).

7.5.2 Perceived Differences for INSTRUCTION (i_{no} vs i_{yes}). We found a tradeoff in the perception of instructions: On the negative side, the UIs that allowed users to enter instructions to the AI received ratings of manually editing suggestions significantly more ($x=1.54$, $p=0.02$) and being significantly more distracting ($x=2.03$, $p<0.0005$).

On the positive side, giving instructions was rated significantly better on being able to influence the suggested text ($x=1.92$, $p<0.001$). Descriptively, it was also rated better on feeling in control of the suggested text (see Q_{10} in Figure 3), although this was not significant ($x=1.42$, $p=0.058$).

7.5.3 Interactions of NUMBER and INSTRUCTION. As mentioned in the previous two parts, both single suggestions and the ability to give instructions were perceived as significantly more distracting. However, there was also a significant negative interaction effect of INSTRUCTION and NUMBER on distraction. The increase in distraction between s_1 compared to s_3 was lower for i_{yes} than i_{no} (which also matches the picture for Q_2 in Figure 3). This seems to be in line with the earlier finding for acceptance rates (Section 7.1): Possibly, finding more useful single suggestions with instructions reduced the otherwise perceived distraction of single suggestions and/or instructions. That said, note that for all suggestion UIs, the majority did not find them distracting. We return to the aspect of distraction in more detail when analysing the open feedback (Section 8).

⁹e.g.: <https://www.grammarly.com/blog/transition-words-phrases/>, <https://writingcenter.unc.edu/tips-and-tools/transitions/>, last accessed March 7, 2023

Section	Aspect	Sig. pos. predictors	Sig. neg. predictors	Sig. interaction	Takeaway in words
1	7.1 Suggestion acceptance		s_1 ($\beta=-1.26$, SE=0.11, CI _{95%} =[-1.48, -1.03], $p<.0001$)	NUMBER * INSTRUCTION ($\beta=0.72$, SE=0.17, CI _{95%} =[0.39, 1.05], $p<.0001$)	Showing one suggestion (instead of three) decreases chance of acceptance; more so without instructions than with them.
2	7.2 Task time, comparing sugg. UIs				No sig. differences in task completion times were found between the four UIs with suggestions.
3	7.2 Task time, comparing sugg. UIs against baseline (no suggestions)	s_3 ($\beta=24.59$, SE=9.5, CI _{95%} =[5.96, 43.22], $p<.01$)			Writing with three suggestions took longer than without suggestions.
4	7.3 Text length, comparing sugg. UIs		s_1 ($\beta=-0.08$, SE=0.01, CI _{95%} =[-0.10, -0.06], $p<.0001$); i_{yes} ($\beta=-0.06$, SE=0.01, CI _{95%} =[-0.09, -0.04], $p<.0001$)		Texts are slightly shorter when writing with single suggestions or with a UI allowing for instructions...
5	7.3 Text length, comparing sugg. UIs against baseline (no suggestions)	s_3 ($\beta=0.05$, SE=0.01, CI _{95%} =[0.03, 0.07], $p<.0001$)	s_1 ($\beta=-0.03$, SE=0.01, CI _{95%} =[-0.05, -0.01], $p<.005$); i_{yes} ($\beta=-0.05$, SE=0.01, CI _{95%} =[-0.07, -0.03], $p<.0001$)		..., also compared to the baseline. However, writing with multiple suggestions leads to slightly longer texts.
6	7.4.1 Requesting suggestions after sentence vs mid-sentence		s_1 ($\beta=-0.83$, SE=0.13, CI _{95%} =[-1.08, -0.57], $p<.0001$); i_{yes} ($\beta=-0.32$, SE=0.14, CI _{95%} =[-0.59, -0.05], $p=0.018$)		Showing one suggestion (instead of three) decreased the chance of requesting suggestions at the beginning of a new sentence.
7	7.4.2 Number of words in sentence at suggestion request	s_1 ($\beta=1.55$, SE=0.78, CI _{95%} =[0.01, 3.08], $p=0.049$)			Showing one suggestion (instead of three) increased the number of words in a sentence after which suggestions were requested.
8	7.4.3 Type of words at suggestion request		s_1 ($\beta=-0.37$, SE=0.14, CI _{95%} =[-0.65, -0.09], $p=0.010$)		Showing one suggestion (instead of three) decreased the chance of requesting suggestions after a transition word.

Table 1: Overview of the (generalised) LMM results and takeaways of the significant results. Empty cells indicate no significant results. See Sections 7.1 - 7.4 for details and Figure 2 for a descriptive overview of the data.

Finally, we also asked two questions that focused on the instructions directly (Q_{15} and Q_{16}) and thus could only be asked for those UIs with instructions (i.e. there’s only a non-diegetic row in Figure 3 for Q_{15} and Q_{16}). For these two questions, we found no significant differences between s_1 and s_3 .

7.6 Instruction Usage and Content

In total, participants used the non-diegetic prompting option to send 397 instructions to the system, with an average of 3.08 instructions per person (SD: 3.40). The mean instruction length was 14.20 characters (SD: 9.01) and 2.52 words (SD: 1.74). On average, participants had a ratio of 0.19 (SD: 0.17) of entering an instruction text when requesting suggestions, for those tasks that offered to do so. That is, about every fifth suggestion request used instructions.

We identified three main instruction “styles”: The most common one (171 usages) was to use single *keywords* (or comma-separated lists of keywords). We also found an *imperative* style with 59 occurrences (e.g. starting the prompt text with “suggest”, “give”, “find”, “describe”). In 12 cases, participants formulated a *question* (e.g. starting with a w-word like “what”, “who” and so on, and/or ending with a “?”). Other cases included instructions consisting of multiple

words to describe something (e.g. “somewhere in Italy”). Qualitatively, we found a range of approaches (Table 2).

7.7 Evaluation of Text Quality

The mean number of spelling and grammar mistakes per word was 0.0025, which is comparable to values reported in previous research [26]. Approximately 3.5% of all texts (23 out of 645) did not align with the selected topics. Of these, the majority (13 out of 23) were written for the category of “shapeshifter”, which may have been misunderstood as a metaphor for a specific set of desired traits in a partner. Despite this potential misunderstanding, the majority of participants demonstrated attentiveness to the task and provided thoughtful reflections on the topic.

8 OPEN FEEDBACK

We analyzed the final feedback as described in Section 5. We structure this report by the emerging aspects.

8.1 Comments on Suggestions

The majority preferred multiple suggestions (75 people stated this preference vs 23 for single suggestions). Main reasons were higher chances of finding fitting suggestions (coded 36 times) and more

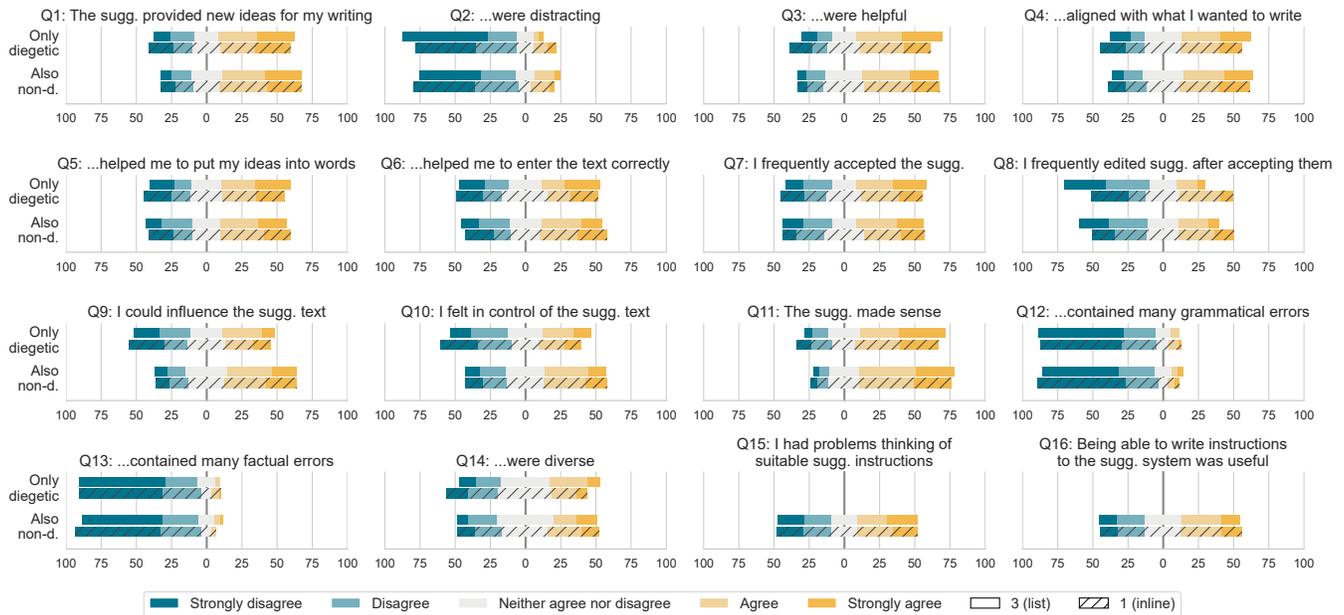


Figure 3: Overview of the Likert results. These questions were asked after each writing task. Note that Q15 and Q16 relate to the instructions (i.e. non-diegetic prompts) and thus were only asked for the corresponding tasks.

Approach	Examples
providing a topic	“school”, “book”, “retirement”, “zoo”, “event”
providing adjectives	“good”, “horrendous”, “bad”, “long, too much, insane”, “friendly”, “funny”, “scared”
request for inspiration	“give me a horror story”, “suggest a place”, “suggest the next step”, “things we do in the morning”, “suggest an activity for a middle aged man”
make idea more concrete	“suggest something disgusting”, “what is wrong with dad”, “suggest a cocktail”, “suggest a type of pistol”
request for variation	“another phrase”, “another action outside”, “suggest a different approach”, “anything”
request for writing help	“are words for stereotypical”, “find a synonym for valued”, “suggest a word for young people”, “another word for talent”
ask for opinion/advice	“are books good”, “what do i do next”
retrieve facts	“closest galaxy”, “side effect of anti ageing”, “a place on the Danube”

Table 2: An overview of the different approaches for writing non-diegetic prompts during the user study. Participants used single keywords to suggest topics and adjectives. Multiple-word prompts were often written in the imperative style, phrased as questions or phrased as incomplete sentences without a verb.

inspiration (coded 8 times). As P_{38} wrote: “I found the multiple suggestions much more user friendly and also much more inspiring due to the multiple options.”

Those preferring single suggestions found them more intuitive (coded 7 times), faster to work with (coded 3 times) or less distracting (coded 3 times): “I strongly preferred inline due to how intuitive

they were to use.” (P_{113}) Or: “I like seeing how the sentences actually looks in its actual place, and the inline suggestions allowed this.” (P_{109}). Others liked not having to decide (coded 2 times) but noted that this might lead to choosing a less than optimal suggestion.

Fourteen participants reflected on benefits for both, such as: “On the one hand, the inline suggestions felt less cluttered and I could just press tab again if the first suggestion wasn’t suitable. On the other hand, displaying multiple suggestions at once could lead me to a better suggestion when I might just have settled for the first one.” (P_{13})

Participants commented on why and how to use suggestions, mentioning inspiration (coded 39 times), overcoming writer’s block (coded 3 times), or finishing sentences (coded 7 times). For example: “I used the suggestions if they aligned with what I was writing or if I felt a little stuck with what to say next.” (P_{117}) Or: “I tended to start with a vague idea of my own and see what ideas it had.” (P_{14}).

Eighteen participants explicitly commented on suggestion quality: Eight were negative (P_{91} : “[...] had to edit most of it.”). Nine felt suggestions were hit or miss (P_{130} : “Sometimes [the suggestions] helped, sometimes it didn’t.”). Two left positive comments (P_{27} : “I was sceptical about whether the AI would align with my ideas or suggest phrasing that I would actually use but most often it did so and I was pleasantly surprised by the results.”). P_{109} further noted that “[instructions] helped the AI write more detailed and interesting sentences” when the direction of the sentence was known beforehand. Using “one or two words in the instructions to get better sentences” that participant continued that “[s]ometimes [it] worked, but quite often I just ended up writing my own sentences, or changing the suggested sentences substantially.” (see Figure 5).

“I liked it better without [the instructions], just let the AI do its thing. That seems more human, that’s how I share story telling with my grand children, we just take turns.” (P₃₄)

It was a normal Thursday morning when Matt Damon was kidnapped. It was like he just disappeared off the face of the earth. This caused a huge worldwide search. People from all over went to extreme lengths to try and find the beautiful actor and no one was willing to give up until he was safe. An old couple who were a huge fan of Matt, spent hours walking around places they'd never been to before in hopes they'd find him. Carrying around weapons just in case, they were putting their lives on the line for him. After a few hours when it was getting dark, they saw a sighting - they weren't sure what it was, it didn't look human. As they got closer, they realised it was Matt Damon in the flesh but he wasn't.. him. It looked like he had just morphed into a completely different person. Different features, different voice - the couple weren't sure whether to believe it was actually him. At the end of the day they decided they should bring it to the police, they were their only hope in finding out what happened, or to potentially get the old Matt back. To be continued..

Accepted text suggestions

Figure 4: Text sample of P₃₄ who took turns with the AI to write about the kidnapping of Matt Damon. The suggestions were taken verbatim and mostly requested at the start or in the middle of a sentence.

funny

restaurant

Fun way

train

Dating is a funny thing. It can be like a rest[au]rant, where you pay tons of money and you expect a great meal but you might get the worst meal in the world. Nowadays, this has been complicated by the rise of social media, and apps, such as Facebook and tinder. Personally, I'm still a fan of the "old school" way of dating. It is more of a fun way of dating, and it is a great way to meet new people. It also is a little more exciting, and I think rewarding whether or not t[he] date is successful. On one date, I remember a going to the train station with the girl I was dating and we ate a nice meal together and had a great time. going to the zoo. It was great.

Accepted text suggestions

User provided non-diegetic prompt

“I think the suggestions helped the AI write more detailed and interesting sentences. I usually had an idea of where I wanted the sentence to go, and I used a word or two in the instructions to get better sentences. Sometimes this worked, but quite often I just ended up writing my own sentences, or changing the suggested sentences substantially.” (P₁₀₉)

Figure 5: Text sample of P₁₀₉ who provided non-diegetic prompts to guide the LLM. For the last instruction (“train”), P₁₀₉ decided to then modify the topic to “zoo”.

8.2 Comments on Instructions

Opinions diverged on instructions: 21 participants explicitly stated they preferred the UIs allowing for instructions, while 24 preferred those without them. 22 participants reflected on both pros and cons. The main reasons for using instructions were getting more suitable suggestions (coded 12 times) (e.g. P₂₇: “I found the instructions more helpful as I could guide the AI when needed.”), inspiration for words (coded 29 times) (e.g. P₆₅: “[...] it gave me inspiration when i was stuck for words.”) and delegating tasks like coming up with places, names or synonyms. (coded 5 times). One person used the AI “[...] to get suggestions for and against the point I was trying to make.” (P₈₅).

In contrast, some found it hard to write instructions (see Section 8.4 for details). Six participants described a trial-and-error approach to find out how to best write instructions.

It was also reported that coming up with instructions can disrupt the writing flow (coded 3 times) and thus reduces efficiency, or is not worth the effort. For example: “It made no difference, as i never felt the need to give it specific instructions. I felt it did a pretty good job of knowing what sort of suggestions I wanted.” (P₃₈).

Some said writing with instructions felt less natural (coded 3 times): “I mostly enjoyed writing without the instructions. I felt more like I was ‘one’ with the AI and it felt like it was more of a team member with me than a piece of software. I think because it removed that feeling of using a computer to help me write I felt like the suggested writing was an extension of myself.” (P₁₃₂). And P₃₄ wrote: “I liked it better without [the instructions], just let the AI do its thing. That seems more human, that’s how I share story telling with my grand children, we just take turns.” (see Figure 4).

8.3 Control and Influence

Eleven participants commented on control and influencing suggestions. For example: “I prefer[r]red multiple because - literally - there were multiple to choose from and that gave me a better feeling of control over the story.” (P₅₉). Another commented: “I like the suggestion systems especially when I was able to provide guidance.” (P₈₂). Overall, multiple suggestions and instructions were mentioned here as contributing to feeling in control, matching the Likert results on control and influence (Q₉ and Q₁₀ in Figure 3).

Moreover, participants commented on strategies around what we now call diegetic prompting in this paper. For example, some preferred influencing the suggestion with the diegetic approach: “I didn’t have much success providing instructions, was having trouble thinking of suggestions quickly and instead focused on directing the topic towards a place where viable suggestions would be made without interactive input.” (P₉₉). Similarly, P₁₁₁ said: “Often it was just as difficult to think of the instruction as it would be to actually write something. It seemed just as easy to start writing what I wanted in order to push the AI in the direction I wanted it to go.”

In contrast, some disliked diegetic prompting: “Without instructions was highly annoying, had to shape your lead-in sentences to get it to say something relevant. The instructions were intuitive and usually got it right.” (P₇₈, also see Figure 6).

Finally, others noticed influences on their own writing processes related to diegetic prompting: “[W]hen I was on my own I just rambled on but while working with the AI I was mentally setting up what

I wrote to be able to ask for a suggestion at a point where the ideas could go in different directions, depending on what was suggested.” (P₉₉, also see Figure 7). And similarly, P₉ wrote: “[I] noticed that the more time I spent the more my tendency was to find a way to write that would facilitate the suggestion to be meaningful and at the same time interesting to add to give more in-depth to my story.”

8.4 Learnability

Several participants (33) touched on challenges of learnability and writing instructions: “I found coming up with suggestions [to the AI] difficult, really. Having to type the start of a sentence and then type what I wanted in a smaller box felt quite clunky and not worth the effort for what was generated. It felt much more fluid when the AI recognized what I wanted and completed the writing without needing suggestions.” (P₂₉). P₆ said: “I almost felt stressed trying to think of some instructions to give to the AI; it felt really hard to me. I’m glad that the option was there, but I guess I wasn’t taking full advantage of it.” Fittingly, 25 participants said they did not use instructions much because, for example, “[...] I wasn’t very good of thinking of them.” (P₂). Some of the previous comments (Section 8.3) fit this aspect as well.

8.5 Distraction

Nine participants explicitly reflected on distraction. For example P₂₀ wrote: “I actually found the suggestions fairly distracting and not helpful. I tended to already know what I wanted to say so the chances of the suggestions aligning with my thoughts were fairly slim.” P₄₃ perceived instructions in particular as distracting: “I feel like writing without instructions help me focus more and [I] am less distracted which allows my sentences to flow and be more natural. Instructions are good if [I] am stuck and need help.”

8.6 Perception of the AI and Expectations

The comments indicate two fundamental views on the role of the AI: Some expected the system to serve efficiency. For example: “I think there is a lot of value in this system, but inputting instructions make it quite long winded and onerous, negating any benefits there may be. I preferred the multi suggestions without instruction.” (P₂₆). Also see the first quote on “alignment” above (Section 8.5). In contrast, others saw the system as serving inspiration (also see Section 8.2 and Section 7.6). They asked the AI for content suggestions or were curious to see in which direction the AI would take the story. This included feeling inspired by suggestions even without accepting them: “I was reading the suggestions either to use them or to just get ideas of what I was writing about” (P₆₇).

9 DISCUSSION

9.1 Choice vs. Control

Our findings contribute to the literature on prompt-based interaction with generative systems for writing: Participants overall preferred choosing from multiple text suggestions presented to them, over actively writing instructions, in short creative and argumentative writing tasks. This is evident from highest acceptance rates with the multiple suggestions UI (Section 7.1), which were not improved through instructions, and from the qualitative feedback,

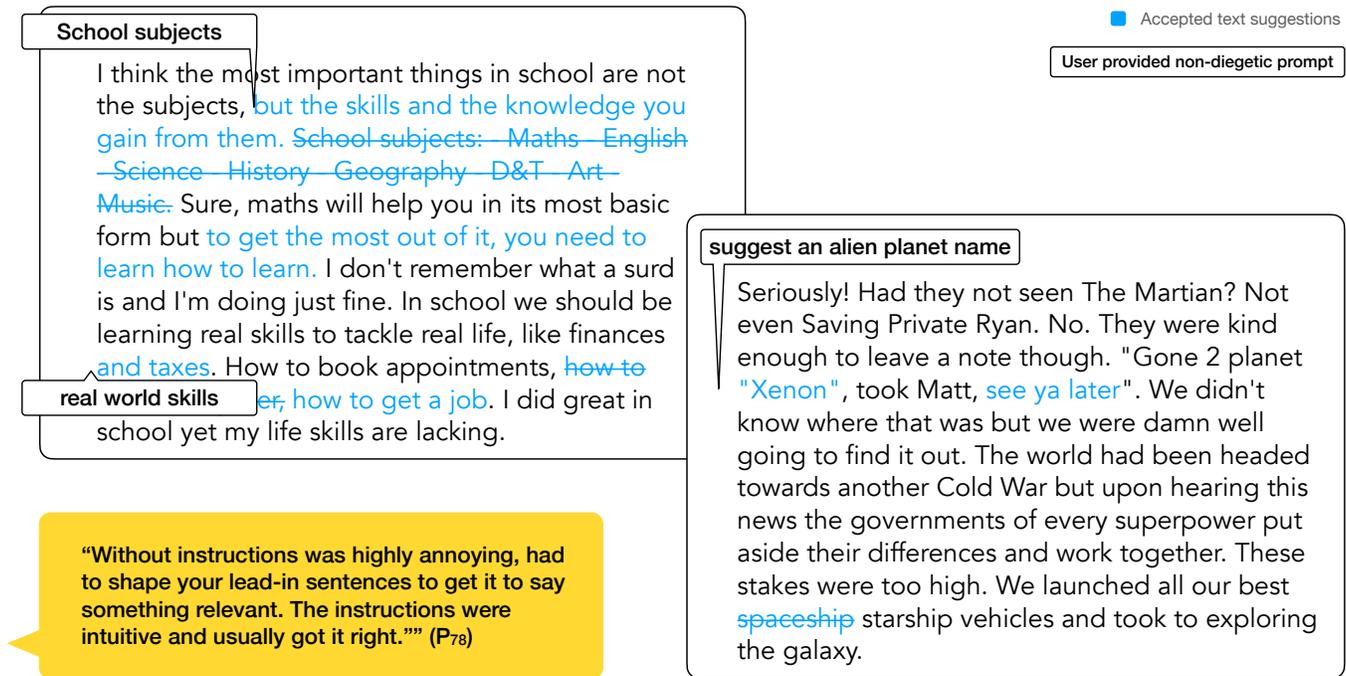


Figure 6: Text sample of P_{78} who used non-diegetic prompting to retrieve a list of “school subjects”. The accepted suggestion is highlighted in blue. Part of the accepted suggestions was later on deleted.

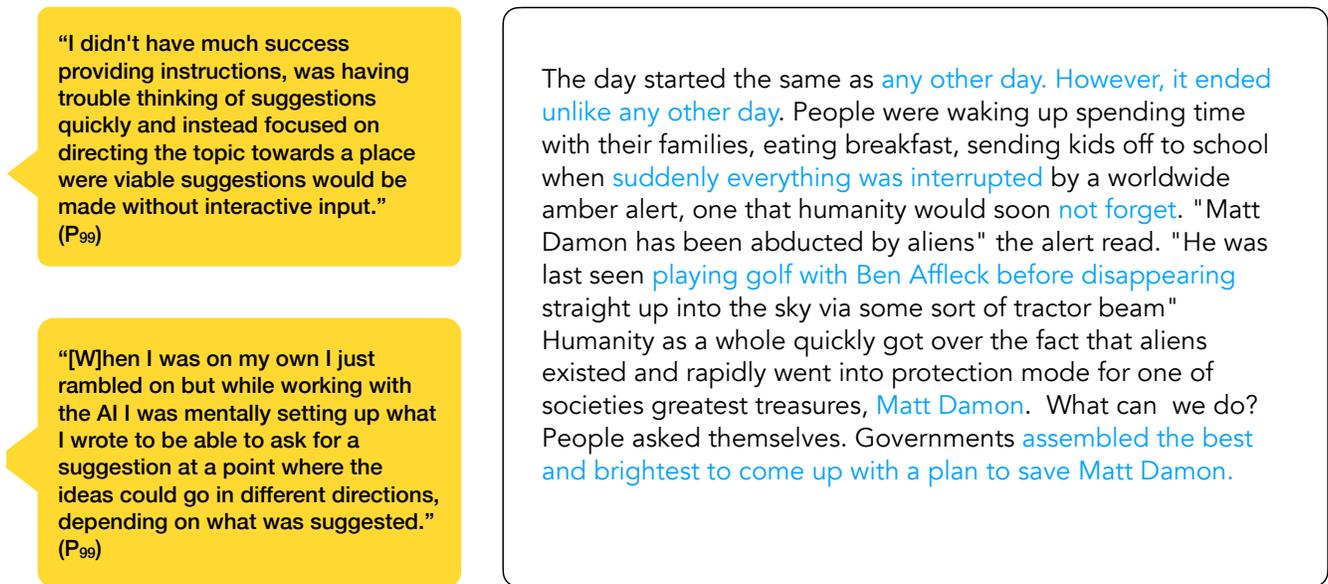


Figure 7: Text sample of P_{99} who found it difficult to provide non-diegetic prompts and instead focused on guiding the suggestion through diegetic content, e.g. requesting suggestions after “...sending kids off to school when” (line 3). By setting the sentence up in this way before requesting a suggestion this participant guided the LLM to suggestions that “go in different directions” (cf. comment in the second yellow box).

where a clear majority favored multiple suggestions, while opinions were divided on instructions (Section 8).

However, giving users more control options in the UI by adding non-diegetic prompting partially mitigated the drawback of a lack of suggestion choice: Instructions increased acceptance rates for single suggestions – although these still did not reach the rate for multiple suggestions (Section 7.1). This indicates that the control offered by instructions was useful to guide single suggestions but not better than having a choice of three suggestions to begin with.

We discuss possible reasons: First, participants might *satisfice* [40], that is, accept a “good enough” suggestion rather than trying to “optimize” it via instructions. Suggestions might also already be good enough so that there is no need for instructions, as supported by some comments (Section 8.2). Second, a known usability principle is *recognition over recall* [31]: Users might find it easier to recognize a presented suggestion as suitable (or not), compared to coming up with an instruction and typing it in. Third, *convenience* might lead participants in the study to accept suggestions without instructions to get through the tasks quickly. However, participants accepted suggestions at a rate comparable with related work (with multiple suggestions and explicit request via tab key: 74 % here vs 72 % in [26]). For suggestions based on user instructions, our rate (64 %) is higher than in a related study design where users could enter requests in a sidebar (17.6 % in [48]): This suggests that potential influences of the study setup do not necessarily work against instructions, or are less dominant than the effects of the UI design (e.g. sidebar vs integration at text cursor). Moreover, times and texts, in combination with the comments, further support the conclusion that participants took the tasks seriously (see Section 7.7).

At the same time, instructions were indeed (situationally) useful: Participants commented on their benefits (Section 8.2), used them in every fifth suggestion request, and experimented with different styles (Section 7.6). Together, these findings motivate the HCI community to further explore the *integration* of choice and control via prompting. For example, future work could build on our conceptual lens to envision further UI designs that combine diegetic and non-diegetic prompting, and use our data as a benchmark in their evaluation.

9.2 Guiding Suggestions with Diegetic Prompts

Our results add to the literature on writing with AI by revealing that people specify more diegetic information to offset the lack of suggestion choice in UIs that display only a single suggestion. This is based on the first large-scale analysis of where in the text users request suggestions: Users wrote about 1.5 more words in the sentence before requesting single suggestions, compared to multiple ones. Moreover, single suggestions were requested less frequently to start a new sentence and to continue after a transition word. Possibly, receiving a single suggestion is less useful here, given that new sentences and transition words signal “openness” for potential changes to the direction of the narrative.

Currently, there is one other (small-scale) analysis of trigger moments (N=4 in [7]). Thus, we encourage the community to analyze trigger moments whenever studying UIs with explicit suggestion triggers.

Fittingly, we indeed recently see high interest in interaction designs where users explicitly request suggestions (e.g. [7, 26, 41]). Our study explores this design space further by looking at how it interacts with the number of suggestions: Here, we contribute evidence that people consider when to request suggestions, and in particular for single suggestions they request them at points in their text that are expected to give clearer guidance to the text continuation system. Future work could examine whether this holds in other writing contexts and to what extent users actively think about when to request suggestions while writing. Based on people’s comments, at least some strategically thought about what we term diegetic prompting (see Section 8.3).

As a related aspect, prior work focused on how people *react* to suggestions (e.g. evaluation fatigue [4], integrative leaps [41]). Complementary, the above results indicate that there is also a *proactive* direction: Writers think about suggestions *before* seeing them. Future work could investigate this in more detail, in particular for UIs in which users explicitly request suggestions.

9.3 Challenges of Integrating Non-Diegetic Prompts

We extract two concrete challenges of interacting via non-diegetic prompts to guide future research and design.

9.3.1 Non-Diegetic Prompts Interrupt the Writing Process. Writing involves multiple cognitive processes, such as coming up with a thought, turning it into words, and entering it [21]. Recently, Bhat et al. [4] studied (without non-diegetic prompts) how this is impacted by text suggestions. For example, writers need to evaluate displayed suggestions. Here, our study adds insights into the relative impact of diegetic vs non-diegetic prompts: Crucially, switching from diegetic writing to non-diegetic instructing forces writers to shift from thinking about their narrative or argument to thinking about instructions to the system. This is reflected in people’s comments (Section 8.2, 8.3, 8.4) and the Likert results on distraction and problems with thinking of instructions (Q_2 and Q_{15} in Figure 3). In contrast, diegetic prompts do not require such shifts, although they still require engagement with displayed suggestions [4, 6].

9.3.2 Non-Diegetic Prompts can be Hard to Write. Even after making that shift, then writing effective non-diegetic prompts is difficult, adding to related findings in the literature [48]: Many participants struggled with this and recognised that they did so in self-reflection (Section 8.2, 8.3, 8.4). More positively, the non-diegetic prompts collected in our study show how users experimented with different styles. These might evolve further with longer use. At the moment, none of these styles go beyond what would also be a meaningful comment to a human co-author.

9.4 Perceived Role of the AI

Here we discuss how users perceived the AI and support this discussion by reflecting on three writing processes as in the framework for analyzing writer-suggestion interactions by Bhat et al. [4]: (1) *proposing* new topics or ideas, (2) *translating* abstract thoughts or keywords into sentences, (3) *transcribing* (i.e. entering) words.

9.4.1 Two Perspectives on the Main Role: Proposer vs Transcriber. Some people clearly saw the system as something that serves input

efficiency (i.e. *transcriber*), whereas others saw it as providing inspiration (i.e. *proposer*). The former are more critical about the system since it would only be good if it is fast and predicts exactly what they want. Based on the qualitative feedback we think that the chosen topic as well as participants' familiarity with the topic might have an influence on their writing mindset. For argumentative writing and, more generally, when people already had an opinion about a topic, they felt that the AI was distracting if it proposed something other than what participants had in mind. Future work may have a closer look at the influence of topic genre and prior knowledge about a topic on the perception of the role of the AI. Study designs should take this difference into account when choosing writing topics to calibrate metrics for performance or exploration.

9.4.2 Non-diegetic Prompts Reflect Users' Perception of the AI. We can further discuss how the content of non-diegetic prompts reflects varying perceptions of the role of the AI: Considering the writing processes [4], non-diegetic prompts from our dataset show that users requested the AI to *propose* inspirational ideas. Sometimes users also only provided partial phrases or keywords, or asked for word choices, which puts the AI into the role of *translating* these abstract ideas into full sentences. At other times, they perceived the AI as a *transcriber* for input efficiency (Section 8).

Other non-diegetic prompts indicate influences on the perceived role beyond these writing processes: For example, people asked the AI for an opinion or advice, or to lookup information. Thus, non-diegetic prompts may shift perception of the AI's role towards a writing collaborator.

9.5 Limitations and Reflections on Methodology

People wrote for five minutes with each UI. Hence, they spent ten minutes in total with each individual UI feature across the writing tasks (single and multiple suggestions, with and without instructions). This is comparable to related work (e.g. 11 min [26], 4 min [6], 10-12 min [48]). Future studies should investigate long-term use, in particular to observe how non-diegetic prompts evolve as writers gain experience with a system.

We prototyped our system with GPT-3 via an API. We did not have access to the model directly and we do not claim to have identified the "best" settings for our specific usage of the model. We noticed two limitations: Sometimes, suggestions were repetitive (e.g. similar ones in one list) or repeated the instruction text (which seems unhelpful). Nevertheless, suggestions were rated highly overall (Section 7.5).

Potential changes to the model over time are beyond our control. This limits exact replicability for studies like this. We see a trend of limited direct access to state-of-the-art LLMs for parts of the academic community, which is not easy to resolve. On the positive side, our work shows that it is possible to construct and study in detail interactive applications built on existing models.

We chose an online setup in line with recent related work (e.g. [6, 26]) to collect logging data from interactions of many people. However, we could not observe people directly or ask questions at interesting moments in the interaction, except for in our pre-study, which we used to refine our design. A small-N study with direct

observation and think-aloud could complement our work, for example, to understand decision-making around triggering suggestions and writing non-diegetic prompts in more detail. Nevertheless, we received rich qualitative feedback as well (Section 8).

It is possible that the instruction styles (Section 7.6) are biased by the provided examples (Figure 8 in Appendix A). Our pre-study showed that such examples are needed to help people get started with this new feature. Nevertheless, people experimented beyond these examples (e.g. questions, writing help, advice, etc.; see Table 2).

With the pop-up box, we tested one way of integrating instructions. This UI element is motivated as a simple way of integrating instructions with the established design of a suggestion list (or inline suggestion). A similar pop-up is used in recent related work (not for instructions but for suggestions in the middle of sentences; cf. [4]). Other designs should be explored in the future.

Finally, we emphasize the importance of open writing tasks in HCI research. Historically, transcription tasks have dominated text entry research (cf. [45]). With the rising interest in human-AI co-creation, research on writing tools needs new tasks. These might not necessarily focus on measuring input speed but rather cover a range of topics, text types, and other aspects. Pragmatically, writing tasks from writer communities and custom tasks have been used in recent studies (e.g. [6, 26, 41, 48]), including ours. As a community, we should systematically evaluate and curate such writing tasks if they are to become a lasting key methodological component.

9.6 Beyond Writing: Diegetic and Non-diegetic Interaction in Generative Systems

We have studied diegetic and non-diegetic prompts to draft text (i.e. *text to text*). Here we reflect on this new perspective by discussing concrete examples of how other interactive generative systems use diegetic and non-diegetic prompting.

- *Visual to Text* Chung et al. [9] proposed a new story ideation tool that uses visual sketching to guide a LLM. Here the sketch is translated to a text prompt. This interaction is non-diegetic.
- *Text to Visual* Recent text to image models allow users to generate images from text descriptions [30, 32, 36]. These are non-diegetic prompts, because they are not part of the visuals.
- *Visual to Visual* Bau et al. [3] show an example of "painting shapes" to guide image models: Users draw simple shapes such as a triangle to symbolize a mountain. The image model then translates these shapes into a high-fidelity rendering. Since the abstract shape is usually not part of the outcome we consider this interaction non-diegetic. On the other hand, Ha and Eck [19] enable users to start painting a part of an image (i.e. providing diegetic information) and let the system continue or finish the painting.

Differentiating these two perspectives therefore allows researchers to analyse users' intention and behavior when interacting or designing systems with generative AI. As shown in the following discussion we can use this understanding to derive implications on the design of interactions for generative models.

9.7 Implications for LLMs and User Interfaces

In recent work by Schick et al. [38], their LLM “PEER” is explicitly trained to follow non-diegetic prompts related to text revision. Effectively, our study contributes the HCI counterpart – an investigation of a UI and interaction design to integrate an LLM in such a role into the writing process. Our results guide future work at this intersection of HCI and NLP in two concrete ways:

First, based on our collected non-diegetic prompts these LLMs should be trained to understand a broader range of inputs. For instance, PEER is trained on the *imperative*-style but we found the *keyword*-style to be more common. Alternatively, users need to be guided towards the supported style via the UI.

Second, while LLMs are rapidly improving, even the best model cannot eliminate cognitive costs and interaction costs of switching between diegetic and non-diegetic writing. This motivates further studies on interaction designs that require such switches and potential pathways to making them easier and more efficient.

10 CONCLUSION

Our new understanding highlights that people use two types of prompting to guide LLMs for text generation. While related work has presented systems that focused on non-diegetic prompts, our findings reveal that users additionally think about and shape their text to guide LLMs through diegetic information. With our UI design that allows for both types, using GPT-3, participants preferred choosing from multiple suggestions over writing instructions. We conclude by highlighting three key takeaways based on our results:

First, writing instructions to the AI requires effort, including switching between diegetic and non-diegetic writing. Second, people combine diegetic and non-diegetic prompting, as single suggestions benefitted from both. Third, writers use their draft (i.e. diegetic information) and suggestion timing to strategically guide LLMs, based on our analysis of when people request suggestions, as well as their self-reflection in comments.

We encourage future work to further analyze these prompt types to develop better writing tools and generalize to other domains (e.g. interaction with generative models for images). To facilitate this, we release our prototype and material on the study and analysis here:

<https://osf.io/qwakj>

ACKNOWLEDGMENTS

We thank Lukas Mecke for feedback on the manuscript. This project is funded by the Bavarian State Ministry of Science and the Arts and coordinated by the Bavarian Research Institute for Digital Transformation (bidt).

REFERENCES

- [1] Kenneth C. Arnold, Krysta Chauncey, and Krzysztof Z. Gajos. 2018. Sentiment Bias in Predictive Text Recommendations Results in Biased Writing. In *Proceedings of the 44th Graphics Interface Conference* (Toronto, Canada) (*GI '18*). Canadian Human-Computer Communications Society, Waterloo, CAN, 42–49. <https://doi.org/10.20380/GI2018.07>
- [2] Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. 2015. Fitting Linear Mixed-Effects Models Using lme4. *Journal of Statistical Software* 67, 1 (2015), 1–48. <https://doi.org/10.18637/jss.v067.i01>
- [3] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Bolei Zhou, Joshua B. Tenenbaum, William T. Freeman, and Antonio Torralba. 2018. *GAN Dissection: Visualizing and Understanding Generative Adversarial Networks*. Technical Report arXiv:1811.10597. arXiv. <https://doi.org/10.48550/arXiv.1811.10597> [cs] type: article.
- [4] Advait Bhat, Saaket Agashe, Niharika Mohile, Parth Oberoi, Ravi Jangir, and Anirudha Joshi. 2022. Studying writer-suggestion interaction: A qualitative study to understand writer interaction with aligned/misaligned next-phrase suggestion. <https://doi.org/10.48550/ARXIV.2208.00636>
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 1877–1901. <https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf>
- [6] Daniel Buschek, Martin Zürn, and Malin Eiband. 2021. The Impact of Multiple Parallel Phrase Suggestions on Email Input and Composition Behaviour of Native and Non-Native English Writers. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (*CHI '21*). Association for Computing Machinery, New York, NY, USA, Article 732, 13 pages. <https://doi.org/10.1145/3411764.3445372>
- [7] Alex Calderwood, Vivian Qiu, Katy Ilonka Gero, and Lydia B. Chilton. 2020. How Novelists Use Generative Language Models: An Exploratory User Study.. In *HAI-GEN+ user2agent@ IUI*.
- [8] Mia Xu Chen, Benjamin N. Lee, Gagan Bansal, Yuan Cao, Shuyuan Zhang, Justin Lu, Jackie Tsay, Yinan Wang, Andrew M. Dai, Zhifeng Chen, Timothy Sohn, and Yonghui Wu. 2019. Gmail Smart Compose: Real-Time Assisted Writing. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Anchorage, AK, USA) (*KDD '19*). Association for Computing Machinery, New York, NY, USA, 2287–2295. <https://doi.org/10.1145/3292500.3330723>
- [9] John Joon Young Chung, Wooseok Kim, Kang Min Yoo, Hwaran Lee, Eytan Adar, and Minsuk Chang. 2022. TaleBrush: Sketching Stories with Generative Pretrained Language Models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (*CHI '22*). Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3491102.3501819> event-place: New Orleans, LA, USA.
- [10] Juliet M Corbin. 1990. *Basics of qualitative research: Grounded theory procedures and techniques*. Sage.
- [11] Mark Dunlop and John Levine. 2012. Multidimensional Pareto Optimization of Touchscreen Keyboards for Speed, Familiarity and Improved Spell Checking. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Austin, Texas, USA) (*CHI '12*). Association for Computing Machinery, New York, NY, USA, 2669–2678. <https://doi.org/10.1145/2207676.2208659>
- [12] Andrew Fowler, Kurt Partridge, Ciprian Chelba, Xiaojun Bi, Tom Ouyang, and Shumin Zhai. 2015. Effects of Language Modeling and Its Personalization on Touchscreen Typing Performance. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems* (Seoul, Republic of Korea) (*CHI '15*). Association for Computing Machinery, New York, NY, USA, 649–658. <https://doi.org/10.1145/2702123.2702503>
- [13] Katy Gero, Alex Calderwood, Charlotte Li, and Lydia Chilton. 2022. A Design Space for Writing Support Tools Using a Cognitive Process Model of Writing. In *Proceedings of the First Workshop on Intelligent and Interactive Writing Assistants (In2Writing 2022)*. Association for Computational Linguistics, Dublin, Ireland, 11–24. <https://aclanthology.org/2022.in2writing-1.2>
- [14] Katy Ilonka Gero, Vivian Liu, and Lydia Chilton. 2022. Sparks: Inspiration for Science Writing Using Language Models. In *Designing Interactive Systems Conference* (Virtual Event, Australia) (*DIS '22*). Association for Computing Machinery, New York, NY, USA, 1002–1019. <https://doi.org/10.1145/3532106.3533533>
- [15] Mayank Goel, Leah Findlater, and Jacob Wobbrock. 2012. WalkType: Using Accelerometer Data to Accommodate Situational Impairments in Mobile Touch Screen Text Entry. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Austin, Texas, USA) (*CHI '12*). Association for Computing Machinery, New York, NY, USA, 2687–2696. <https://doi.org/10.1145/2207676.2208662>
- [16] Mayank Goel, Alex Jansen, Travis Mandel, Shwetak N. Patel, and Jacob O. Wobbrock. 2013. ContextType: Using Hand Posture Information to Improve Mobile Touch Screen Text Entry. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Paris, France) (*CHI '13*). Association for Computing Machinery, New York, NY, USA, 2795–2798. <https://doi.org/10.1145/2470654.2481386>
- [17] Steven Goodman, Erin Buehler, Patrick Clary, Andy Coenen, Aaron Michael Donsbach, Tiffanie Horne, Michal Lahav, Bob MacDonald, Rain Breaw Michaels, Ajit Narayanan, Mahima Pushkarna, Joel Christopher Riley, Alex Santana, Lei Shi, Rachel Sweeney, Phil Weaver, Ann Yuan, and Meredith Ringel Morris. 2022. LaMPost: Evaluation of an AI-assisted Writing Email Editor Prototype for Adults with Dyslexia. <https://arxiv.org/abs/2207.02308>

- [18] Mitchell Gordon, Tom Ouyang, and Shumin Zhai. 2016. WatchWriter: Tap and Gesture Typing on a Smartwatch Miniature Keyboard with Statistical Decoding. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (San Jose, California, USA) (CHI '16). Association for Computing Machinery, New York, NY, USA, 3817–3821. <https://doi.org/10.1145/2858036.2858242>
- [19] David Ha and Douglas Eck. 2017. A Neural Representation of Sketch Drawings. <http://arxiv.org/abs/1704.03477> arXiv:1704.03477 [cs, stat].
- [20] Adi Haviv, Jonathan Berant, and Amir Globerson. 2021. BERTese: Learning to Speak to BERT. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Association for Computational Linguistics, Online, 3618–3623. <https://doi.org/10.18653/v1/2021.eacl-main.316>
- [21] John R. Hayes. 2012. Modeling and Remodeling Writing. *Written Communication* 29, 3 (July 2012), 369–388. <https://doi.org/10.1177/0741088312451260>
- [22] Ellen Jiang, Kristen Olson, Edwin Toh, Alejandra Molina, Aaron Donsbach, Michael Terry, and Carrie J. Cai. 2022. PromptMaker: Prompt-based Prototyping with Large Language Models. In *CHI Conference on Human Factors in Computing Systems Extended Abstracts*. 1–8.
- [23] Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. 2020. How Can We Know What Language Models Know? *Transactions of the Association for Computational Linguistics* 8 (Dec. 2020), 423–438. https://doi.org/10.1162/tacl_a_00324
- [24] Anjali Kannan, Karol Kurach, Sujith Ravi, Tobias Kaufmann, Andrew Tomkins, Balint Miklos, Greg Corrado, Laszlo Lukacs, Marina Ganea, Peter Young, and Vivek Ramavajjala. 2016. Smart Reply: Automated Response Suggestion for Email. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (San Francisco, California, USA) (KDD '16). Association for Computing Machinery, New York, NY, USA, 955–964. <https://doi.org/10.1145/2939672.2939801>
- [25] Alexandra Kuznetsova, Per B. Brockhoff, and Rune H. B. Christensen. 2017. lmerTest Package: Tests in Linear Mixed Effects Models. *Journal of Statistical Software* 82, 13 (2017), 1–26. <https://doi.org/10.18637/jss.v082.i13>
- [26] Mina Lee, Percy Liang, and Qian Yang. 2022. CoAuthor: Designing a Human-AI Collaborative Writing Dataset for Exploring Language Model Capabilities. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (CHI '22). Association for Computing Machinery, New York, NY, USA, Article 388, 19 pages. <https://doi.org/10.1145/3491102.3502030>
- [27] Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What Makes Good In-Context Examples for GPT-3? <http://arxiv.org/abs/2101.06804> arXiv:2101.06804 [cs].
- [28] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing. *arXiv:2107.13586 [cs]* (July 2021). <http://arxiv.org/abs/2107.13586> arXiv: 2107.13586.
- [29] Michael J. Muller and Sandra Kogan. 2012. Grounded Theory Method in Human-Computer Interaction and Computer-Supported Cooperative Work. In *The Human-Computer Interaction Handbook* (3 ed.). CRC Press. Num Pages: 21.
- [30] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2021. GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models. (2021). <https://doi.org/10.48550/ARXIV.2112.10741> Publisher: arXiv Version Number: 3.
- [31] Jakob Nielsen. 1994. Enhancing the Explanatory Power of Usability Heuristics. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Boston, Massachusetts, USA) (CHI '94). Association for Computing Machinery, New York, NY, USA, 152–158. <https://doi.org/10.1145/191666.191729>
- [32] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. 2021. StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery. (2021). <https://doi.org/10.48550/ARXIV.2103.17249> Publisher: arXiv Version Number: 1.
- [33] Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. 2019. Language Models as Knowledge Bases? <http://arxiv.org/abs/1909.01066> arXiv:1909.01066 [cs].
- [34] Philip Quinn and Shumin Zhai. 2016. A Cost-Benefit Study of Text Entry Suggestion Interaction. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (San Jose, California, USA) (CHI '16). Association for Computing Machinery, New York, NY, USA, 83–88. <https://doi.org/10.1145/2858036.2858305>
- [35] R Core Team. 2020. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org>
- [36] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical Text-Conditional Image Generation with CLIP Latents. (2022). <https://doi.org/10.48550/ARXIV.2204.06125> Publisher: arXiv Version Number: 1.
- [37] Melissa Roemmele and Andrew S. Gordon. 2015. Creative Help: A Story Writing Assistant. In *Interactive Storytelling (Lecture Notes in Computer Science)*, Henrik Schoenau-Fog, Luis Emilio Bruni, Sandy Louchart, and Sarune Baceviciute (Eds.). Springer International Publishing, Cham, 81–92. https://doi.org/10.1007/978-3-319-27036-4_8
- [38] Timo Schick, Jane Dwivedi-Yu, Zhengbao Jiang, Fabio Petroni, Patrick Lewis, Gautier Izacard, Qingfei You, Christoforos Nalmpantis, Edouard Grave, and Sebastian Riedel. 2022. PEER: A Collaborative Language Model. <https://doi.org/10.48550/ARXIV.2208.11663>
- [39] Timo Schick and Hinrich Schütze. 2021. Exploiting Cloze Questions for Few Shot Text Classification and Natural Language Inference. <http://arxiv.org/abs/2001.07676> arXiv:2001.07676 [cs].
- [40] Herbert Alexander Simon. 1996. *The sciences of the artificial* (3. ed. ed.). MIT Press, Cambridge, Mass.
- [41] Nikhil Singh, Guillermo Bernal, Daria Savchenko, and Elena L. Glassman. 2022. Where to Hide a Stolen Elephant: Leaps in Creative Writing with Multimodal Machine Intelligence. *ACM Trans. Comput.-Hum. Interact.* (jan 2022). <https://doi.org/10.1145/3511599> Just Accepted.
- [42] Hendrik Strobelt, Albert Webson, Victor Sanh, Benjamin Hoover, Johanna Beyer, Hanspeter Pfister, and Alexander M. Rush. 2022. Interactive and Visual Prompt Engineering for Ad-hoc Task Adaptation with Large Language Models. <http://arxiv.org/abs/2208.07852> arXiv:2208.07852 [cs].
- [43] Ben Swanson, Kory Mathewson, Ben Pietrzak, Sherol Chen, and Monica Dinulescu. 2021. Story Centaur: Large Language Model Few Shot Learning as a Creative Writing Tool. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*. Association for Computational Linguistics, Online, 244–256. <https://doi.org/10.18653/v1/2021.eacl-demos.29>
- [44] Anestis Touloumis. 2015. R Package multgee: A Generalized Estimating Equations Solver for Multinomial Responses. *Journal of Statistical Software* 64, 8 (2015), 1–14. <http://www.jstatsoft.org/v64/i08/>
- [45] Keith Vertanen and Per Ola Kristensson. 2014. Complementing Text Entry Evaluations with a Composition Task. *ACM Trans. Comput.-Hum. Interact.* 21, 2, Article 8 (feb 2014), 33 pages. <https://doi.org/10.1145/2555691>
- [46] Tongshuang Wu, Ellen Jiang, Aaron Donsbach, Jeff Gray, Alejandra Molina, Michael Terry, and Carrie J. Cai. 2022. PromptChainer: Chaining Large Language Model Prompts through Visual Programming. <http://arxiv.org/abs/2203.06566> Number: arXiv:2203.06566 arXiv:2203.06566 [cs].
- [47] Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022. AI Chains: Transparent and Controllable Human-AI Interaction by Chaining Large Language Model Prompts. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (CHI '22). Association for Computing Machinery, New York, NY, USA, Article 385, 22 pages. <https://doi.org/10.1145/3491102.3517582>
- [48] Ann Yuan, Andy Coenen, Emily Reif, and Daphne Ippolito. 2022. Wordcraft: Story Writing With Large Language Models. In *27th International Conference on Intelligent User Interfaces* (Helsinki, Finland) (IUI '22). Association for Computing Machinery, New York, NY, USA, 841–852. <https://doi.org/10.1145/3490099.3511105>
- [49] Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. BARTScore: Evaluating Generated Text as Text Generation. <http://arxiv.org/abs/2106.11520> arXiv:2106.11520 [cs].
- [50] Tony Z. Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate Before Use: Improving Few-Shot Performance of Language Models. <http://arxiv.org/abs/2102.09690> arXiv:2102.09690 [cs].

A APPENDIX

In this appendix, we provide a table of logged events and additional screenshots.

No.	Interaction Event	Description
1	EVENT_CONFIRM_INSTRUCTION	User has confirmed instruction (Enter Key)
2	EVENT_CANCEL_INSTRUCTION	User has cancelled the instruction (ESC key or clicking outside the instruction box)
3	EVENT_OPEN_INSTRUCTION_BOX	User has triggered new suggestions in the “with instructions” writing setting (Tab Key)
4	EVENT_SELECT_NEXT_SUGGESTION	User has selected next suggestion (Down Arrow Key)
5	EVENT_SELECT_PREV_SUGGESTION	User has selected previous suggestions (Up Arrow Key)
6	EVENT_REQUEST_SUGGESTIONS	User has requested new suggestions (Tab Key)
7	EVENT_SUGGESTIONS_RESPONSE	System returned suggestions
8	EVENT_CONFIRM_SUGGESTION	User has selected and confirmed one suggestion (Enter Key or Mouse Selection)
9	EVENT_CANCEL_SUGGESTION	User has cancelled the suggestions (ESC key or clicking outside the suggestion box)
10	EVENT_TASK_STATUS	Can be either “task started” or “task finished”
11	EVENT_KEYDOWN	User has pressed a key, e.g “A” or “TAB”

Table 3: An overview of the interaction events logged in the user study.

Inline Suggestion + Instructions

suggest a place in Southeast Asia

I was planning to travel to SE Asia to visit the temples and other historic places.

Inline suggestion with an instruction.

- Suggestions are shown inline
- You can control suggestion by instructing the AI.
- Please try out the suggestion and instruction features, but also write yourself

You can give keywords: e.g.

travel, train, summer

You can tell the AI what to do: e.g.

suggest a warm destination.

You can request information: e.g.

somewhere in Asia

Anything else you want to try:

anything

Tab Request new suggestions

Enter Request new suggestions given the (optional) instructions.

Enter Press again to confirm the instruction to the AI and the current selected suggestion.

ESC Cancel instruction and suggestions.

Suggestion Text: Confirm the suggestion with **ENTER**. Typing other keys will cancel the suggestion.

Spend at most **5 minutes** to write about the topic. Once you are done writing, click on the **I'm done writing** and then the **Next** on the bottom right corner to proceed.

Note: The **I'm done writing** will be enabled after 15 seconds.

Selected Topic

Is Listening to a Book Just as Good as Reading It? Do you listen to audiobooks? What are the benefits, in your opinion, of listening instead of reading? Are there advantages to reading that cannot be gained by listening? Which method do you prefer? Why?

Listening to a book is like going on an adventure. Optional: Type here are free, your mind has time to wander and imagine. You can listen to the narration of another person and enjoy their voice. A good narration voice can make the story come to life. **You can become the character.**

Time: 1:11

I'M DONE WRITING

Figure 8: Screenshot of the writing interface. (Left Side) The info box describes the available functionalities in the current setting, (Top Middle) the selected topic, (Bottom Middle) the text editor with the current written text and an inline suggestion.

Choose one topic that you wish to write about.

Selected Topic

How Worried Should We Be About Screen Time During the Pandemic? The coronavirus pandemic ended the screen time debate: Screens won. We all now find ourselves on our screens for school, work, and connecting with family and friends during this time of social distancing and increased isolation. But should we be worried about this excessive screen use right now? Or should we finally get over it and embrace the benefits of our digital devices?

[SHOW PREVIOUS TOPIC](#) [SHOW NEXT TOPIC](#)

I'M READY TO WRITE ABOUT THIS TOPIC

Figure 9: The topic selection panel. Users can browse through the topics and indicate that they are ready to write about the depicted topic.