

Tracing How Annotators Think: Augmenting Preference Judgments with Reading Processes

Karin de Langis, William Walker, Khanh Chi Le, Dongyeop Kang

Department of Computer Science and Engineering

University of Minnesota

{dento019, dongyeop}@umn.edu

Abstract

We propose an annotation approach that captures not only labels but also the *reading process* underlying annotators' decisions, e.g., what parts of the text they focus on, re-read or skim. Using this approach, we conduct a case study on the preference annotation task and create a dataset PREFERREAD that contains fine-grained annotator reading behaviors obtained from mouse tracking. PREFERREAD enables detailed analysis of how annotators navigate between a prompt and two candidate responses before selecting their preference. We find that annotators re-read a response in roughly half of all trials, most often revisiting the option they ultimately choose, and rarely revisit the prompt. Reading behaviors are also significantly related to annotation outcomes: re-reading is associated with higher inter-annotator agreement, whereas long reading paths and times are associated with lower agreement. These results demonstrate that reading processes provide a complementary cognitive dimension for understanding annotator reliability, decision-making and disagreement in complex, subjective NLP tasks.

Keywords: Cognitive Methods; Corpus (Creation, Annotation, etc.); Tools, Systems, Applications

1. Introduction

Most natural language processing (NLP) datasets capture the final labels assigned to a passage – e.g., a sentiment rating or an answer span – while ignoring the *process* by which annotators arrive at those decisions. This omission leaves valuable dimensions of human annotations unexplored. Understanding which specific words or paragraphs annotators focused on, paragraphs they chose to re-read, or even areas of the text to skim or ignore can reveal aspects of annotator reasoning and confidence, providing richer insight into complex or subjective labeling decisions.

We propose augmenting NLP annotations with *reading processes* that reflect how annotators allocate attention during the annotation task. Such data would typically be obtained through eye tracking, which offers fine-grained signals for readers' cognitive activity. However, eye-tracking studies require in-person participation and specialized hardware and are not easily scalable. (Online alternatives like webcams suffer from limitations such as unreliable calibration, making them not well-suited for capturing the small reading eye movements (Hutt and D'Mello, 2022).) To work around these constraints, we augment our annotation interface with *mouse tracking*, a paradigm recently validated in psychological experiments (Wilcox et al., 2024). Mouse tracking estimates requires users to reveal text incrementally with their mouse cursor as they read, and its close alignment with eye-tracking data, makes it a promising alternative for large-scale, remote data collection.

We investigate using mouse tracking to collect

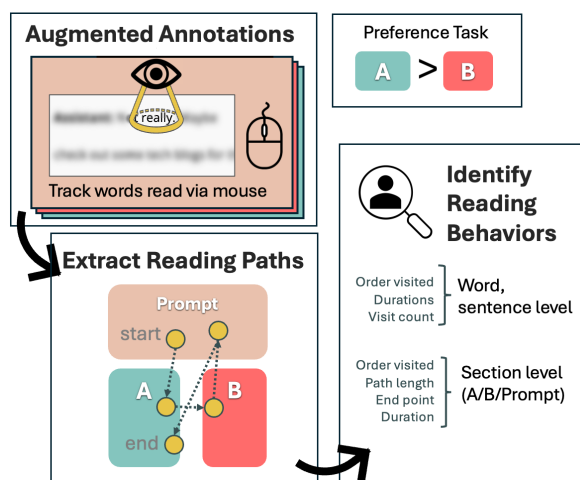


Figure 1: Our annotation interface records reading behavior of annotators as they evaluate and select preferred texts, providing fine-grained insights into the decision-making process during annotation.

annotator reading behaviors through a case study on the *preference annotation* task, an important annotation type in reinforcement learning for LLM alignment (Stiennon et al., 2020; Ouyang et al., 2022). In this task, annotators read two model-generated responses and select the one they find more appropriate. However, these annotations can be highly subjective and noisy, with high inter-annotator variability (Dsouza and Kovatchev, 2025; Zhang et al., 2025). Little is known about whether annotators' reading processes – such as re-reading prompts or selectively reading responses – are associated with their final decisions or with inter-

annotator agreement. We propose that these reading behaviors can be a valuable signal of annotator cognitive processes and decision-making.

To contribute to the understanding of annotator reading processes, we present PREFERREAD, a dataset of 1,000 preference annotation items, each annotated by three participants. In addition to annotator preferences and rationales, PREFERREAD also captures the detailed mouse-tracked reading process of each annotator. From these data, we propose and extract several cognitively informed reading measures, including reading paths between sections, re-reading frequencies, and word-level reading durations, to enrich traditional annotations with behavioral information (see Figure 1). To our knowledge, this is the first study that analyzes the reading behavior of annotators. The enriched representations from our annotation schema processes enables us to investigate questions like:

- Do annotators *re-read* their chosen or rejected choices?
- Are there any parts of the text that annotators *ignore*?
- Do these reading patterns *influence annotation decisions*?

Our analysis reveals systemic links between annotation agreement and reading behavior. For instance, annotators who re-read responses tend to exhibit higher agreement, while longer reading paths and times are associated with lower agreement and appear to reflect greater annotator uncertainty. We also find that annotators disproportionately skip the latter portions of responses, possibly in pursuit of *cognitive economy* during decision-making, which may undermine some assumptions implicit in both task design and subsequent data use. These findings suggest that reading-process data can reveal some of the cognitive dynamics underlying preference judgments, providing an additional dimension for modeling annotator reliability and (dis)agreement. Our code and data are publicly available.¹

2. Related Work

Our work draws on prior research in mouse tracking, annotation augmentations, and preference tasks.

Mouse Tracking Eye-tracking data provides a detailed signal of online cognitive processes during reading (Just and Carpenter, 1980; Rayner, 1998). However, physical requirements of eye-tracking experiments make it impossible to collect data at

scale. Although eye-tracking with webcams is possible, the frame rate and calibration precision of webcams is typically unsuitable for capturing reading eye movements (Hutt and D’Mello, 2022), making data crowdsourcing infeasible. To address this issue, research in psychology and computer vision has explored mouse tracking for capturing reading behaviors. The Mouse Tracking for Reading (MoTR) paradigm (Wilcox et al., 2024), based on earlier tools like MouseView.js (Anwyl-Irvine et al., 2022), shows that mouse tracking effectively captures natural reading behavior, including word skipping and regressions, and shows comparable sensitivity to eye-tracking in detecting psycholinguistic effects (Wilcox et al., 2024). Research in computer vision has also found that eye-tracking and mouse-tracking data are functionally comparable, see e.g., Salicon (Jiang et al., 2015) and BubbleView (Kim et al., 2017). Inspired by these findings, our study employs mouse tracking within an NLP annotation interface in order to capture reading behaviors.

In NLP, eye-tracking data has primarily been leveraged outside of annotation contexts. Specifically, word-level eye-tracking metrics such as gaze duration and fixation counts have been used to improve model performance in tasks like question answering, paraphrasing, and text simplification (Zhang and Hollenstein, 2024; Hollenstein, 2021; Sood et al., 2020); or to evaluate cognitive plausibility of language models (Eberle et al., 2022).

Annotation Augmentations Annotator disagreement has often been treated as noise, with methods like majority voting or averaging used to derive a single “gold” label per instance. However, this approach can obscure meaningful variation especially for complex or subjective labeling tasks, and as a result, recent research has explored modeling annotator disagreement explicitly. For instance, Wan et al. (2023) demonstrated that incorporating annotators’ demographic information can enhance the prediction of disagreement levels.

In a similar vein, we propose augmenting annotations with reading processes of annotators to show dynamic factors that relate to annotation decisions. Because research in both psychology (Rawson et al., 2000) and human-computer interaction (Leroy and Kammerer, 2023; Cheng et al., 2015) has shown that specific reading behaviors like re-reading are associated with greater reading comprehension, we include several re-reading based metrics from annotators in PREFERREAD.

Preference Annotations In reinforcement learning from human feedback (RLHF) (Stiennon et al., 2020; Ouyang et al., 2022), preference annotations play a central role in training reward models (Jiang et al., 2024; Winata et al., 2024). Compared to

¹<https://github.com/minnesotanlp/reading-behavior>

Task: Determine which response is more *helpful*.

PROMPT

Human: What are some interesting flavorings i can mix with my beef patties?

RESPONSE A (CHOSEN)

Assistant: I think you'll find it interesting that humans have put mayonnaise, mustard, and other types of condiments in beef patties. But you know, the most interesting thing about mayonnaise is it's used as a weapon, and here are some other interesting uses of mayonnaise. I hope this helps, I have several other recipes I can offer if you're interested.

RESPONSE B

Assistant: Hm, I wonder if you're looking for something like a teriyaki marinade, or a Cajun mix.

Word-level gaze duration estimates

High-level path captures how the sections were traversed

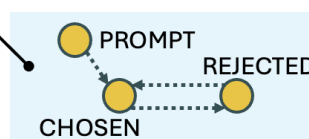


Figure 2: An example item from PREFERREAD with the collected label and reading behaviors from one annotator. Darker highlights indicate relatively long gaze durations, and lighter highlights indicate relatively short gaze durations. Word-level gaze estimates can hint at aspects of human decision-making; for instance, attention is concentrated on the odd turn of phrase “the most interesting thing about mayonnaise is it’s used as a weapon” in the chosen response.

other annotation types (e.g., ratings), preference-based data has been found to provide relatively stronger and more reliable signals, with lower variance across annotators (Gatt and Kraemer, 2018).² However, preference data is still highly subjective and subject to high variability across annotators (Dsouza and Kovatchev, 2025; Wang et al., 2024). Furuta et al. (2024) propose addressing ambiguity in preference labels through distributional, rather than binary, representation of preferences. Just et al. (2025) include annotator rationales alongside preference labels. We propose annotator reading behaviors as an additional dimension to represent preference decisions.

	Mean	SD	Min	Max
Chosen	142	84	15	320
Rejected	130	81	18	338
Prompt	84	72	4	231
Total	356	225	52	825

Table 1: Word lengths of stimuli used in the dataset. Each preference pair is divided into prompt, chosen response, and rejected response.

²However, it should be noted that Dsouza and Kovatchev (2025) recently found comparable reliability between preferences and ratings.

3. PREFERREAD: A Preference Reading Dataset

This section details our stimuli selection, mouse-tracking annotation interface, data collection procedure, and data postprocessing steps. See Figure 2 for an example item in the dataset.

Stimuli. We collect annotations from a randomly selected subset of 1,000 instances from the “helpful” portion³ of the Helpful-Harmless dataset (Anthropic, 2022) collected from human annotators as described in Bai et al. (2022).

Our sampling procedure excludes items that are above the 90th percentile in word count because we want to avoid giving annotators extremely long texts. Additionally, our sampling procedure also excludes response pairs in which both responses contain fewer than three words, as these pairs are often trivially short (e.g., “Sounds good,”/“No problem”). Table 1 summarizes the word-length statistics on the selected stimuli.

Preference Reading Interface. We build a custom web application on top of JSPsych (de Leeuw et al., 2023) to capture preference reading annotations along with mouse tracking data. The main user interface is shown in Figure 3. The application records participants’ mouse movements, captur-

³The “helpful” split is chosen after a pilot study indicated the instructions were more intuitive for participants.

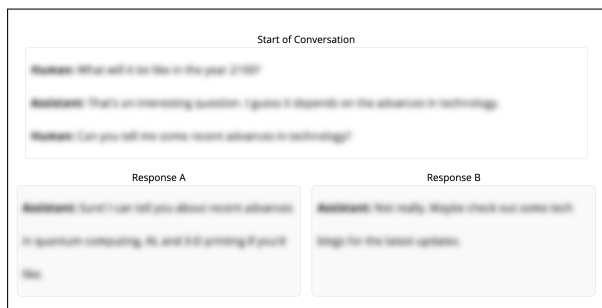


Figure 3: Our annotation interface shows a prompt followed by two possible responses. Participants mouse over blurred text to reveal it and read. All mouse movements are recorded.

ing entry and exit times for each character span to construct a time series of hovered characters and corresponding durations.

The size of the viewable window is an important consideration. If it is too small, the participant is not able to read without considerable effort. If it is too large, it is difficult to estimate which word a participant is attending to at a given moment. Because the average person’s *word identification span* from a single fixation is 9 to 11 characters (Rayner, 1998), we chose to reveal 9 characters at a time in order to allow more precisely measure word-level focus.⁴ This is consistent with the mouse tracking paradigm of Wilcox et al. (2024), which reveals a similar amount of text; however, practitioners more concerned with sentence- or paragraph-level analysis may consider larger viewable windows.

Participants. 300 participants were recruited on Prolific. They were paid \$3 to complete a tutorial and annotate 10 preference pairs. The tutorial contained instructions, two examples, and two practice items. We limited participants to native English speakers with an approval rating of at least 95% on the crowdsourcing platform. Participants self-reported their race, gender, and age. The gender split was 53% male and 47% female. The race of participants was 80% white, 11% black, and 3% Asian. 6% reported mixed race or “Other.” Participants were between 18 and 75 years of age ($M = 39, SD = 13$). The study was by the institutional review board, and consent was obtained from participants prior to the experiment.

Annotation Procedure. Each preference pair was annotated by three participants using our inter-

⁴The word identification span is distinct from the *perceptual span*, which measures how many characters are perceptible from a given focal point. The perceptual span averages between 15 and 20 characters, making it considerably wider than the word identification span.

face. Participants begin with a short mouse tracking practice item and the following two instructions:

- All the text is blurred out, and you’ll need to use your mouse to reveal the text as you read.
- Try to read as normally as possible. Read the dialogue and both responses carefully to decide which is more helpful. You can re-read anything you like.

Participants then receive more detailed instructions about the preference task adapted from those in Bai et al. (2022), which include positive and negative examples. Following the instruction block, participants are presented with 10 trials in a randomized order. In each trial, the position of the two possible responses (i.e., right or left) are also randomized. The 10 trials are pseudo-randomly selected from our dataset such that the average word count for all participants is between 300 and 350 words per trial, in order to ensure that the cumulative reading effort is comparable for all participants. After choosing a response, participants answer a multiple choice question about their rationale. Participants then provide a rationale (“More Helpful,” “More Accurate,” “More Concise,” “Less Harmful,” or “Other”) for their decision.

Data Processing. We clean our collected data prior to analysis to remove uninformative trials and participants. Specifically, we examine the *word coverage*, or percentage of words that the user moused over per trial (considering only non-trivial mouseovers as described in Section 3.1). Note that we expect substantially less than 100% word coverage since users will spend negligible time hovering words such as ‘a’ and ‘the.’ In addition, it may sometimes be reasonable to make a preference decision after viewing a small percentage of the total words, for instance, in the event when one response is excessively verbose. We exclude trials in which participants read fewer than 10% of the total words, amounting to 1.9% of trials. (It is not possible to go back in the experiment, so if a participant accidentally proceeds, they skip the trial. This possibly explains these missing data.)

3.1. Estimating Gaze Durations

We use the recorded mouse data to approximate word-level gaze durations. The raw data from one trial can be viewed as a time series of mouse-move events: we record the character hovered and duration each time the user’s mouse enters or exits a character span. First, we consolidate these character-level events into word-level events by summing durations for characters that belong to the same word, following the procedure outlined by Wilcox et al. (2024). The result is a “word-level” vector w , where the w_t entry contains the word in-

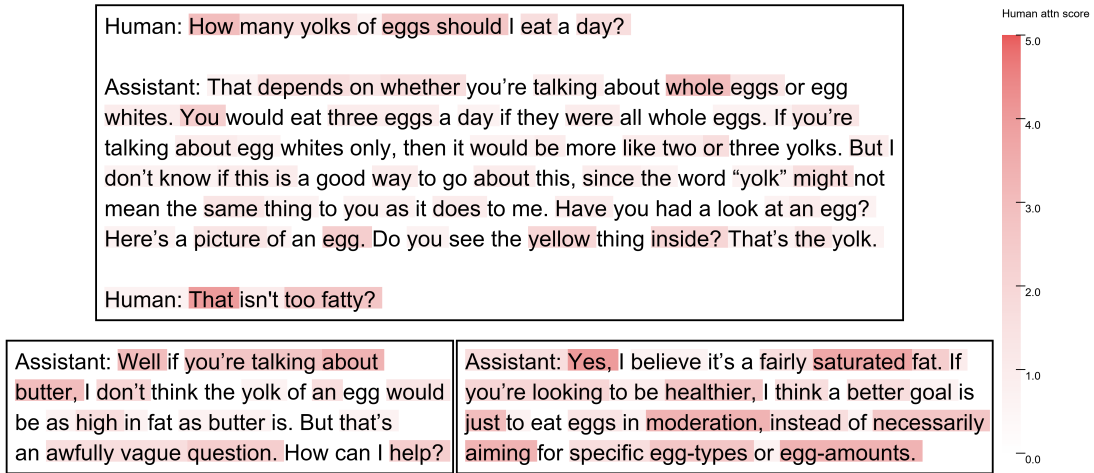


Figure 4: Annotator word duration estimates from an item in PREFERREAD. Darker highlights indicate higher attention, e.g., longer gaze durations, while lighter highlights correspond to little or below-average duration. Notice that in the prompt (top box), the annotator spends more time on the human’s questions than on the assistant’s verbose reply.

dex and duration of time step t . These entries are analogous to the fixations in eye tracking data.

Next, the data is cleaned by discarding fixations whose durations fall outside the interval [160ms, 4000ms], again following (Wilcox et al., 2024). This step is also analogous to standard procedures in eye tracking analysis that remove fixations that do not reflect cognitive processes associated with reading, e.g., when the participant quickly moves past a word without processing the word, or when the participant is distracted and stops on a word while not focused on reading. The total dwell time for each word in the trial is obtained by summing the durations of this cleaned word duration vector over each word index. This results in duration vector d of length n where n is the number of words in the stimulus.

Finally, we normalize the vectors according to the binning strategy described in (Klerke and Plank, 2019), i.e. binning the z-scores of each entry as follows:

$$\text{bin}(z) := \begin{cases} 0 & z = \text{NaN (no duration)} \\ 1 & z < -1 \text{ (very short duration)} \\ 2 & z \in [-1, -0.5) \text{ (short duration)} \\ 3 & z \in [-0.5, 0.5) \text{ (typical duration)} \\ 4 & z \in [0.5, 1) \text{ (long duration)} \\ 5 & z \geq 1 \text{ (very long duration)} \end{cases} \quad (1)$$

We also compute a summary aggregate across all participants for each stimulus s with the arithmetic mean: $d^i = \frac{\sum_{j=1}^p \text{bin}(z(d^{i,j}))}{p}$. This summary score allows us to estimate the relative importance of words across all annotators, as shown in Figure 4, which can be helpful for inferring how hu-

mans engaged with the text.

4. Analysis of Reading Processes

We propose extracting several metrics from annotators’ mouse-tracked behaviors, to augment their preference annotations (Table 2). We analyze these metrics across several dimensions: re-reading behaviors and reading paths (§4.1), text-skipping tendencies (§4.2), and the relationship between reading behaviors and annotator agreement (§4.3).

4.1. Re-reading and reading paths

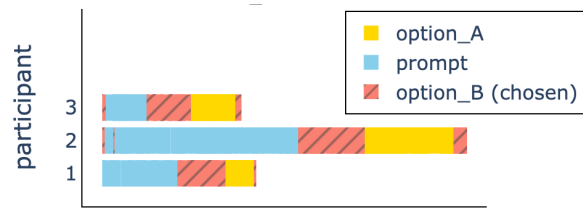


Figure 5: A timeline of the mouse location of three annotators reading one item in our dataset.

Each annotation trial is divided into three sections: the *prompt*, *chosen response*, and *rejected response*. A section is considered “re-read” if the participant returns to it after moving on to another section. To avoid false positives caused by minor mouse movements, we define “reading” as spending at least one second within a given section.⁵

⁵As illustrated by participants 2 and 3 in Fig. 5, who very briefly mouse over “option B” before reading the

Metric	Illustration	Description	Plausible Interpretation
Re-read prompt / rejected / selected	<p>Prompt A B RereadPrompt = True</p>	Binary indicating whether an annotator exited and then returned to the specific section	Importance of the re-read section to the decision; in-decision
Loop between responses	<p>Prompt A B loop = True</p>	Binary indicator of an annotator going back and forth between the selected and rejected responses at least once to create a loop (2-node cycle)	Strong indecision between the two responses; caution in annotating
Path length	<p>Prompt A B length = 2</p>	The number of edges in the path between sections (prompt, selected, rejected) that the annotator traversed before making a final decision	Longer paths indicate indecision or caution, while shorter paths indicate a quick, straightforward read
Time reading response	<p>Start ms</p>	Average number of milliseconds per word spent reading both responses	Relative attention paid to the responses during annotation
Word coverage	<p>coverage $\approx 60\%$</p>	Percentage of words hovered for at least 160ms	Lower word coverage can indicate either careless reading or a clear-cut decision (e.g., one response is clearly inferior)

Table 2: We propose extracting these metrics based on the reading process to augment each annotation and inform our understanding of how annotators engage with the texts and of how that engagement relates to agreement.

Re-reading occurs in roughly half of trials (54%), revealing three notable patterns:

(1) **Annotators most often re-read their chosen response.** They re-read it 38.4% of cases (compared to 25.4% re-read rates for rejected responses and 26.1% for prompts). Relatedly, **annotators frequently read their chosen response last:** 74.0% of annotators who re-read any section returned to their chosen response last, suggesting a final confirmation check before submission.

(2) **Re-reading varies substantially across annotators.** Few participants re-read in every trial or in non; most re-read for only a portion of 10 trials, indicating individual differences in review strategies.

(3) **Re-reading does not correlate with response semantic similarities.** Because highly similar responses may lead to more re-reading, we use Sentence-BERT (Reimers and Gurevych, 2019) to compute cosine similarity between response pairs and examine re-reading frequencies across cosine similarities. We find that response semantic similarity does not correlate with re-reading frequencies (see Figure 7).

prompt, some mouse movements across a section are too brief to represent genuine reading. Instead, they likely stem from lapsed mouse control.

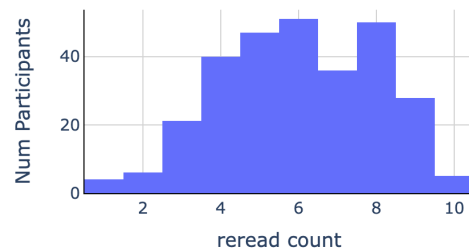


Figure 6: Each annotator completes ten trials. Across the ten trials, few annotators never (re-read count = 0) or always (re-read count = 10) re-read.

(4) **Reading paths** For a more detailed view of how annotators navigated the text, we represent each annotator's sequence of visited sections as a *reading path*. The average path length is 2.47 ($SD = 0.97$), and 17.73% of annotators form a *loop* – a back-and-forth traversal between the two responses. Longer paths and loops indicate greater deliberation or indecision between candidates.

4.2. Text Skipping Tendencies

We find that annotators do not always read to the end of the text: **annotators often make their de-**

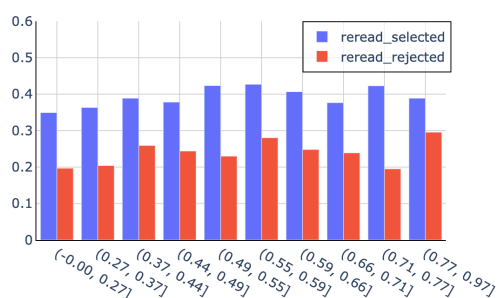


Figure 7: Frequency of response re-reading behavior (y axis) across quantiles of cosine similarities (x axis) between the two responses. Re-reading frequencies do not correlate with response similarity.

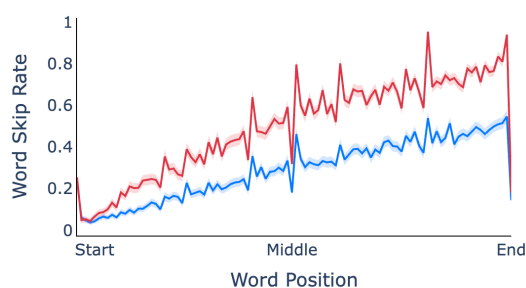


Figure 8: Annotators rarely skip words at the beginning of a response (outside of the initial “Assistant:” word), but as the response continues, they more frequently skip words. This is especially pronounced in rejected responses. Word position is calculated via the word index relative to the response length.

cision based on an unfinished reading of the responses, which we term *cognitive economy in reading*. A paired t-test ($t = -8.52, p < 0.001$) confirms that annotators skip significantly more words in rejected responses ($M = 21.71$) than in chosen one ($M = 15.63$). We conjecture that annotators often recognize less suitable responses early, leading them read less of the rejected than the chosen response. We also find that skipped words cluster near the ends of responses (see Figure 8), suggesting that readers often form judgments *mid-text*.

Annotators also **read a slightly larger proportion of the prompt than of either response**. A paired t-test showed significant difference ($p < .01$) between readers’ coverage of the prompt relative to either response (response coverage $M = .55$ vs. prompt coverage $M = .59$). Recall that we rarely expect coverage to reach 100% due to the low likelihood that readers will spend a non-trivial amount of time moused over words like ‘a’ or ‘the,’ or highly predictable words like the speaker labels (“Human:”, “Assistant:”).

4.3. Agreement and Behavior

We assess how reading metrics relate to annotator agreement. First, we discuss overall inter-annotator agreement (§4.3.1). Second, we examine how annotator agreement relates to both categorical and continuous reading metrics (§4.3.2, §4.3.3). Participant impressions of mouse-tracking are discussed in §4.4.

4.3.1. Inter-annotator Agreement

We measure the overall inter-annotator agreement (IAA) with Krippendorff’s alpha (Krippendorff, 2018). We find $\alpha = 0.25$, which indicates fair agreement among annotators and falls in the range of alphas for subjective annotation tasks reported in Wong et al. (2021). This is also consistent with recent findings in difficulty of the preference task (Cao et al., 2025).

We cannot directly compare IAA with the original dataset since each item had only one annotator. However, Bai et al. (2022) reported 63% agreement with their annotators, while ours align 59% of the time, suggesting a reasonable level of agreement.

4.3.2. Categorical metrics

We applied Chi-square tests of independence to assess whether categorical variables were associated with differences in annotator agreement rates. To determine significance, we set $\alpha = 0.05$ and apply the Bonferroni correction to p values to account for multiple comparisons.

We find that **re-reading is significantly more common among annotators who agree** ($\chi^2(1) = 11.25, p = 0.001$), suggesting that deliberate review aligns with higher-quality annotations. However, **looping between two responses is significantly more common among annotators who disagree** – i.e., those who go back and forth at least twice between the two responses have higher rates of disagreement, suggesting that looping reflects high indecision ($\chi^2(1) = 9.42, p = 0.002$).

Because we also record participants’ stated rationale for each annotation (More Helpful, More Accurate, More Concise, Less Harmful, Other), we investigate whether shared rationale is associated with higher agreement. We find the association between rationale choice and agreement approaches significance ($\chi^2(4) = 9.25, p = 0.055$).

We also assess whether annotators exhibit a bias toward choosing the first or second response and find **no bias toward choosing the first or second response** ($\chi^2 = 0.29, p = 0.593$).

4.3.3. Continuous metrics

To examine the relationship between annotator agreement and continuous metrics (e.g., path

length), we form two groups of annotators: pairs who agreed with one another and pairs who did not. We then compare the continuous metric values between these two groups with an independent t-test, again setting $\alpha = 0.05$ and applying the Bonferroni correction to account for multiple comparisons.

Annotators who **disagree** show **slightly longer path lengths** ($M = 2.56$ vs. $M = 2.42$; $t = -3.3, p = 0.001$) and **slower reading times** ($M = 338.02$ ms/word vs. $M = 315.25$ ms/word; $t = -2.73, p = .006$). Word coverage, on the other hand, does not differ significantly ($M_{agree} = 0.58$ vs. $M_{disagree} = 0.57$; $t = -0.4, p = 0.683$), indicating that skimming strategies do not necessarily harm annotation quality.

We also consider the *word focus overlap* between annotators, as annotators who agree may spend more time focusing on the same subset of words. For each annotator, we consider the focused words to be those with at least a short duration (for details on how this is computed, see §3.1). Then, we assess the similarity of the focused words between each annotator pair with the Jaccard index (i.e., the intersection over union of the sets of focused words from each annotator). In this case, we find that **the word focus overlap is slightly higher among agreeing annotators** ($M = 0.16$ vs. $M = 0.13$); this difference approaches significance ($t = 2.33, p = 0.021$).

4.4. Participant experience

We conducted a pilot study in which participants tested the interface. We were particularly interested in any comments from participants about the mouse tracking, since the paradigm has previously been used in psychology studies that require reading only one to two sentences, and not for longer texts in conjunction with annotation tasks. In the post-study survey (68 participant responses), 20 participants mentioned *mouse tracking*: Of these, 8 were positive (e.g., "fun," "interesting," or perception of increased focus due to the small window of revealed text), 7 negative (due to slower reading speeds from mouse tracking), and 5 neutral. Since the primary complaint from participants is that they feel they are slowed down by the mouse tracking, we recommend future adapters adjust monetary compensation for participants accordingly.

Overall, participants tolerated the mouse tracking paradigm well, even for longer texts, without causing substantial discomfort or annoyance. This suggests that mouse tracking is feasible for longer, cognitively demanding annotation tasks.

5. Discussion

In PREFERREAD, we extract several categorical and continuous metrics to augment preference annotations based on the reading processes observed via a mouse tracking interface. This yields some important insights regarding how annotators approach the preference reading task in general.

First, we note a pattern of cognitive economy in reading for annotators, who have a tendency to skip the latter portion of responses. This is more pronounced in the case of rejected responses, indicating that annotators do not feel the need to completely finish a response to decide that it is inferior. While we don't find that reading less of a response is associated with higher annotator disagreement in our dataset, in some cases, it is possible that this tendency could result in annotators overlooking relevant information. Annotation task designers may consider encouraging readers to finish reading both responses, e.g., by including comprehension questions covering the latter part of responses.

We divide each preference dataset item into three sections: prompt, rejected response, and selected response. We then analyze "re-reading" behavior, which we define as re-entering one of these sections after having exited. We are interested in this re-reading behavior because it seems plausible that some annotators may have a habit of, e.g., re-visiting the prompt before deciding, and that this re-visiting may influence decision-making. We observe that annotators re-read sections in about half (54%) of all trials. Overall, annotators are more likely to re-read their chosen response relative to the rejected response, indicating that annotators may prefer to confirm the quality of their selected response rather than verify the deficiency of their rejected response.

A difficult question in modeling annotator disagreement in subjective tasks like preference selection is whether the disagreement stems from differing opinions, low annotator confidence, or inconsistencies in annotation quality. The augmented annotations in PREFERREAD provide a unique signal in differentiating between potential sources of disagreement. We find significant associations between re-reading and agreement, as well as significant associations between longer paths and looping behaviors with disagreement. Such insights can inform annotation post-processing, e.g., annotators who re-read repeatedly can be flagged as potentially low confidence, whereas annotations stemming from an abnormally fast reading process may be flagged as potentially low-quality. Using reading patterns in conjunction with demographic data may enhance pluralistic opinion modeling.

Future research can also explore methods to explicitly integrate reading signals, such as word-

level scores and reading strategies, into preference model training and reinforcement learning paradigms. For example, re-reading patterns can be explicitly modeled through repetitive encoding of observed behaviors or guided chain-of-thought reasoning to enhance decision-making in preference modeling, and skipped word indices could be leveraged as an auxiliary objective to improve preference decision-making by identifying and filtering out unimportant words.

Finally, while the re-reading behaviors and path definitions discussed here are specific to preference annotations, our methods can be adapted to capture reading behaviors in other challenging and subjective tasks. For instance, domains such as law and medicine involve complex documents with multiple articles that can naturally delineate into sections and high-level reading paths. Such expert annotation tasks require careful engagement from skilled annotators, making them well-suited to the reading augmented annotations we propose here.

6. Conclusion

We introduce `PREFERREAD`, a first-of-its-kind dataset that contains the reading process of annotators as they complete the preference annotation task. Reading processes are captured via a mouse-tracking-enabled annotation interface, demonstrating the feasibility of crowd-sourced, large-scale reading behavior collection for NLP annotations. `PREFERREAD` offers several insights into annotator reading behaviors during the preference annotation task, and we find several relationships between reading processes and annotator agreement. Together, our results suggest that the proposed augmentation strategy provides a novel dimension in understanding human thought processes during an annotation task.

Limitations

The mouse-tracking interface is less naturalistic than ordinary reading, and it may impose additional task demands on annotators. Accordingly, our findings should be interpreted within the context of this specific experimental paradigm.

Ethics Statement

Mouse-tracking data, like all behavioral data, carries a potential risk of revealing participant identities. All data in this study were collected with informed consent, anonymized prior to analysis, and handled in accordance with institutional review protocols. No personally identifiable information was retained or shared.

7. Bibliographical References

- Alexander Anwyl-Irvine, Thomas Armstrong, and Eduard Pogorskiy. 2022. [MouseView.js: Reliable and valid attention tracking in web-based experiments using a cursor-directed aperture](#). *Behavior Research Methods*, 53:1663–1687.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Xiaoyang Cao, Zelai Xu, Mo Guang, Kaiwen Long, Michiel A. Bakker, Yu Wang, and Chao Yu. 2025. [Latent collective preference optimization: A general framework for robust llm alignment](#).
- Shiwei Cheng, Zhiqiang Sun, Lingyun Sun, Kirsten Yee, and Anind K Dey. 2015. Gaze-based annotations for reading comprehension. In *Proceedings of the 33rd annual ACM conference on human factors in computing systems*, pages 1569–1572.
- Joshua R de Leeuw, Rebecca A Gilbert, and Björn Luchterhandt. 2023. jsPsych: Enabling an open-source collaborative ecosystem of behavioral experiments. *Journal of Open Source Software*, 8(85):5351.
- Russel Dsouza and Venelin Kovatchev. 2025. Sources of disagreement in data for LLM instruction tuning. In *Proceedings of Context and Meaning: Navigating Disagreements in NLP Annotation*, pages 20–32.
- Oliver Eberle, Stephanie Brandl, Jonas Pilot, and Anders Søgaard. 2022. Do transformer models show similar attention patterns to task-specific human gaze? In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4295–4309.
- Hiroki Furuta, Kuang-Huei Lee, Shixiang Shane Gu, Yutaka Matsuo, Aleksandra Faust, Heiga Zen, and Izzeddin Gur. 2024. Geometric-averaged preference optimization for soft preference labels. *Advances in Neural Information Processing Systems*, 37:57076–57114.
- Albert Gatt and Emiel Krahmer. 2018. Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. *Journal of Artificial Intelligence Research*, 61:65–170.

- Nora Hollenstein. 2021. *Leveraging Cognitive Processing Signals for Natural Language Understanding*. Doctoral thesis, ETH Zurich, Zurich, Switzerland.
- Stephen Hutt and Sidney K D’Mello. 2022. [Evaluating calibration-free webcam-based eye tracking for gaze-based user modeling](#). In *ICMI’22: Proceedings of the 2022 international conference on multimodal interaction*, pages 224–235.
- Ming Jiang, Shengsheng Huang, Juanyong Duan, and Qi Zhao. 2015. Salicon: Saliency in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1072–1080.
- Ruili Jiang, Kehai Chen, Xuefeng Bai, Zhixuan He, Juntao Li, Muyun Yang, Tiejun Zhao, Liqiang Nie, and Min Zhang. 2024. [A survey on human preference learning for large language models](#).
- Hoang Anh Just, Ming Jin, Anit Kumar Sahu, Huy Phan, and Ruoxi Jia. 2025. Data-centric human preference with rationales for direct preference alignment. In *Second Conference on Language Modeling*.
- Marcel A Just and Patricia A Carpenter. 1980. A theory of reading: From eye fixations to comprehension. *Psychological Review*, 87(4):329.
- Nam Wook Kim, Zoya Bylinskii, Michelle A. Borkin, Krzysztof Z. Gajos, Aude Oliva, Frédo Durand, and Hanspeter Pfister. 2017. [BubbleView: an alternative to eye-tracking for crowdsourcing image importance](#). *CoRR*, abs/1702.05150.
- Sigrid Klerke and Barbara Plank. 2019. At a glance: The impact of gaze aggregation views on syntactic tagging. In *Proceedings of the Beyond Vision and LAnguage: inTEgrating Real-world kNowledge (LANTERN)*, pages 51–61.
- Klaus Krippendorff. 2018. *Content analysis: An introduction to its methodology*. Sage publications.
- Caroline Leroy and Yvonne Kammerer. 2023. Reading multiple documents on a health-related issue: The roles of a text-highlighting tool and re-reading behaviour in integrated understanding. *Behaviour & Information Technology*, 42(14):2331–2352.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Katherine A Rawson, John Dunlosky, and Keith W Thiede. 2000. The rereading effect: Metacomprehension accuracy improves across reading trials. *Memory & Cognition*, 28(6):1004–1010.
- Keith Rayner. 1998. Eye movements in reading and information processing: 20 years of research. *Psychological bulletin*, 124(3):372.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-BERT: Sentence embeddings using siamese BERT-networks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Ekta Sood, Simon Tannert, Philipp Mueller, and Andreas Bulling. 2020. [Improving natural language processing tasks with human gaze-guided neural attention](#). In *Advances in Neural Information Processing Systems*, volume 33. Curran Associates, Inc.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021.
- Ruyuan Wan, Jaehyung Kim, and Dongyeop Kang. 2023. [Everyone’s voice matters: Quantifying annotation disagreement using demographic information](#). In *Proceedings of the 37th AAAI Conference on Artificial Intelligence*, pages 14523–14530.
- Binghai Wang, Rui Zheng, Lu Chen, Zhiheng Xi, Wei Shen, Yuhao Zhou, Dong Yan, Tao Gui, Qi Zhang, and Xuan-Jing Huang. 2024. Reward modeling requires automatic adjustment based on data quality. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4041–4064.
- Ethan Gottlieb Wilcox, Cui Ding, Mrinmaya Sachan, and Lena Ann Jäger. 2024. Mouse tracking for reading (MoTR): A new naturalistic incremental processing measurement tool. *Journal of Memory and Language*, 138:104534.
- Genta Indra Winata, Hanyang Zhao, Anirban Das, Wenpin Tang, David D. Yao, Shi-Xiong Zhang, and Sambit Sahu. 2024. [Preference tuning with human feedback on language, speech, and vision tasks: A survey](#).
- Ka Wong, Praveen Paritosh, and Lora Aroyo. 2021. Cross-replication reliability—an empirical approach to interpreting inter-rater reliability. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the*

11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 7053–7065.

Leran Zhang and Nora Hollenstein. 2024. [Eye-tracking features masking transformer attention in question-answering tasks](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics and Language Resources and Evaluation (LREC-COLING 2024)*, pages 7057–7070. European Language Resources Association.

Michael JQ Zhang, Zhilin Wang, Jena D Hwang, Yi Dong, Olivier Delalleau, Yejin Choi, Eunsol Choi, Xiang Ren, and Valentina Pyatkin. 2025. Diverging preferences: When do annotators disagree and do models know? In *Forty-second International Conference on Machine Learning*.

8. Language Resource References

Anthropic. 2022. *Helpful–Harmless Dataset*. Anthropic, 1.0.