

SEEING EYE TO AI: HUMAN ALIGNMENT VIA GAZE-BASED RESPONSE REWARDS FOR LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Advancements in Natural Language Processing (NLP), have led to the emergence of Large Language Models (LLMs) such as GPT, Llama, Claude, and Gemini, which excel across a range of tasks but require extensive fine-tuning to align their outputs with human expectations. A widely used method for achieving this alignment is Reinforcement Learning from Human Feedback (RLHF), which, despite its success, faces challenges in accurately modelling human preferences. In this paper, we introduce GazeReward, a novel framework that integrates implicit feedback – and specifically eye-tracking (ET) data – into the Reward Model (RM). In addition, we explore how ET-based features can provide insights into user preferences. Through ablation studies we test our framework with different integration methods, LLMs, and ET generator models, demonstrating that our approach significantly improves the accuracy of the RM on established human preference datasets. This work advances the ongoing discussion on optimizing AI alignment with human values, exploring the potential of cognitive data for shaping future NLP research.

1 INTRODUCTION

Recent advancements in Natural Language Processing (NLP) have led to the emergence of Large Language Models (LLMs) like GPT (OpenAI, 2023), Llama (Dubey et al., 2024), Claude (Anthropic, 2024), and Gemini (Team et al., 2024), which excel across a range of tasks. These models, often consisting of billions of parameters, are trained on massive datasets and typically require extensive fine-tuning to align their outputs with human expectations¹. Several works have focused on refining the way LLMs interpret and respond to user intent (Wang et al., 2023b), which has led to the development of novel alignment techniques. A common approach to achieving human alignment involves leveraging explicit human feedback as preference information. Currently, the most widely adopted method is Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022). RLHF has been implemented in many state-of-the-art LLMs (Cui et al., 2024; OpenAI, 2023; Bai et al., 2022b), and has been shown to help align models to human instructions and mitigate the generation of toxic or harmful content (Kiegeland et al., 2024). However, a persistent challenge with this approach is the difficulty of acquiring sufficient high-quality training data (Casper et al., 2023).

To be able to capture the complexities of real-world user instructions, there is a need for meticulously handcrafted data (Wang et al., 2023b), which are resource-expensive and difficult to scale (Yang et al., 2023). Obtaining high-quality feedback from human annotators, usually provided after examining a model response, suffers from several caveats (Casper et al., 2023). For instance, low inter-annotator agreement can result in inconsistent evaluations of the same model output due to varying interpretations, domain expertise, or biases. Moreover, “scalable oversight” – the ability to supervise models effectively with limited resources (Amodei et al., 2016) – remains an open problem. Inconsistent data quality is another issue, as cost-quality tradeoffs often arise when collecting human feedback.

¹LLMs that are trained only on extensive datasets for language modeling are referred to as “pre-trained” LLMs. Subsequent approaches, such as human alignment, are categorized as “post-training”.

To address these challenges, researchers have increasingly turned to LLMs as a form of AI-driven feedback, referred to as Reinforcement Learning from AI Feedback (RLAIF) Bai et al. (2022b). This method offers improved scalability, easier data collection, and cost-efficiency compared to traditional human-driven approaches (Bai et al., 2022b; Wang et al., 2023a; Madaan et al., 2023). However, it remains unclear what type of feedback signals, or a combination of feedback mechanisms, is optimal to align LLM with human goals (Casper et al., 2023). More research is needed to explore the underlying beliefs and expectations of human users (Casper et al., 2023), and how these can be incorporated into human alignment techniques. Furthermore, the alignment success of a language model is dependent on the quality of the underlying RM (Pace et al., 2024). Various alignment methods, such as RLHF, RLAIF, and Direct Preference Optimization (DPO) (Rafailov et al., 2023), rely on RM to incorporate feedback. Reward modelling is also essential for generating synthetic data for preference alignment and is often used in LLM inference to evaluate model outputs in techniques such as best-of-N sampling (Cui et al., 2024).

In this work, we propose a novel approach that incorporates Eye-tracking (ET) as an additional signal to address the challenge of human alignment. ET measures oculomotor behavior i.e. the movements and fixations of the eyes, which offers insight into visual attention and information processing (Kleinke, 1986; Land & Furneaux, 1997). This allows researchers to correlate observable eye movement patterns with underlying cognitive and perceptual processes during reading and language comprehension tasks (Kleinke, 1986; Krasich et al., 2018). Moreover, ET – unlike other (explicit) forms of feedback (e.g., questionnaire data, data annotation) – does not suffer from human biases, and offers a better temporal and spatial resolution (Zhang & Hollenstein, 2024). Several studies have shown a strong correlation between human eye movements and attention patterns in transformer-based models (Wang et al., 2024a; Bensemam et al., 2022; Sood et al., 2020a). Incorporating ET data into NLP tasks has also proven valuable, as demonstrated by numerous works (Huang et al., 2023; Khurana et al., 2023; Hollenstein et al., 2019; Yang & Hollenstein, 2023; Kiegeland et al., 2024; Deng et al., 2023a; Mathias et al., 2018; McGuire & Tomuro, 2021). Recently, Kiegeland et al. (2024) proposed the integration of ET in controlled sentiment generation to create a dataset that can be used in human alignment methods.

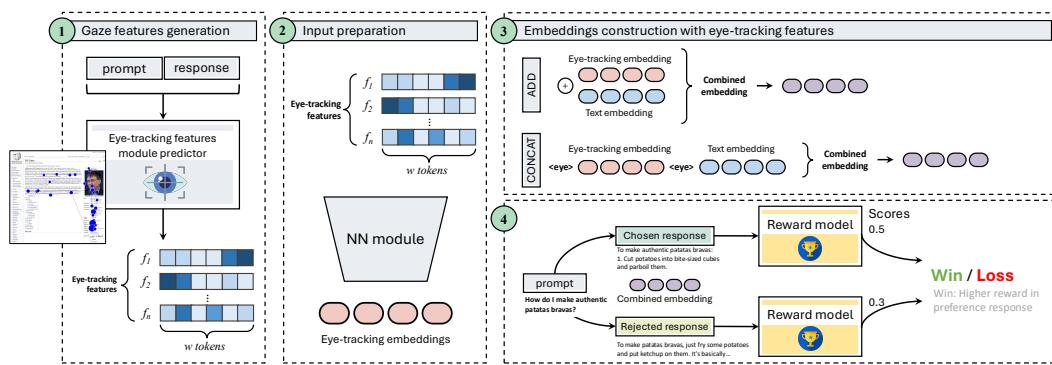


Figure 1: GazeReward Framework for using eye-tracking data for reward modelling. We use a generator model to compute ET features on a preference dataset D and we train the human preference by combining both text and ET embeddings (See section 4 for details.)

Human alignment continues to be one of the biggest challenges in the development of LLM, with RM playing an important role in addressing this issue. This paper investigates how behavioural signals, particularly ET, can be operationalised as implicit feedback to improve human alignment (the proposed approach is shown in Figure 1). Furthermore, we explore the use of ET prediction models that can generate – automatically and with little effort – ET features in response to text input, which makes our solution not only cost-effective but also highly scalable. The main contributions of our work can be summarized as follows:

- We propose **GazeReward**, a novel and scalable framework that integrates implicit feedback in the form of ET data into the RM, a key component in modeling human preferences.

- 108 • We perform for the first time an ablation study that examines several state-of-the-art LLMs,
 109 various ET prediction models, and different methods for incorporating ET features into the
 110 RM.
- 111 • We demonstrate experimentally substantial performance improvements with the GazeRe-
 112 ward framework, showing accuracy gains of over 20% in RM predictions across diverse
 113 human preference datasets.

115 2 PRELIMINARIES

118 2.1 LARGE LANGUAGE MODELS-HUMAN ALIGNMENT

120 LLMs-Human Alignment typically involves training LLMs ² on datasets curated by humans (learning
 121 from human feedback data) (Ouyang et al., 2022). This can be achieved through Supervised
 122 Fine-Tuning (SFT), where the model is trained on pairs of prompts (x) and corresponding human-
 123 generated responses (y) (Liu et al., 2024). Alternatively, alignment can be pursued via preference
 124 optimization, using a human preference dataset that differentiates between a better response (y_w)
 125 and a worse one (y_l) for the same prompt (x): $\mathcal{D} = \{(x^{(i)}, y_w^{(i)}, y_l^{(i)})\}_{i=1}^N$.

126 To this day, RLHF (Ouyang et al., 2022) remains the most popular technique used in state-of-the-art
 127 LLMs like GPT-4 (OpenAI, 2023), Claude (Bai et al., 2022b), Bard (Google, 2023), and Llama
 128 2-Chat (Touvron et al., 2023). Different implementations of RLHF can vary in terms of data col-
 129 lection, training processes, and choice of RL algorithms. Typically, RLHF (Ouyang et al., 2022)
 130 involves three main steps: (1) collecting feedback, (2) training a RM based on that feedback, and (3)
 131 optimising the LLMs using RL techniques, such as Proximal Policy Optimization (PPO) Schulman
 132 et al. (2017). Since RLHF was first introduced, several advancements have been made, including
 133 fine-grained reward systems (Bai et al., 2022b; Wu et al., 2023b; Dong et al., 2023b; Wang et al.,
 134 2023c; 2024c), or replaced the original PPO algorithm with other RL techniques (Wu et al., 2023a).

135 An alternative to RLHF is DPO (Rafailov et al., 2023), which employs an offline RL approach to
 136 optimize language models based on preference data, without the need for a separate reward model.
 137 While DPO can be used independently, it is often complementary to other training methods such
 138 as SFT or statistical rejection sampling, to further improve human alignment based on a RM (Zhao
 139 et al., 2023; Liu et al., 2024; Dubey et al., 2024). Statistical rejection sampling, also called best-of-N
 140 or top-k-over-N (Bai et al., 2022b; Touvron et al., 2023; Dubey et al., 2024) is another widely used
 141 technique. Moreover, certain methods perform human alignment without RL to avoid instabilities,
 142 and fine-tune the model on filtered samples by a RM, or other sources (Dong et al., 2023a; Yuan
 143 et al., 2023).

144 A significant challenge in all human alignment techniques is data acquisition (Casper et al., 2023).
 145 This includes problems such as evaluator misalignment, supervision difficulties, and feedback qual-
 146 ity issues (Casper et al., 2023). However, as AI systems continue to improve, LLMs are increasingly
 147 employed for tasks that were traditionally handled by humans, such as data annotation and gen-
 148 eration. Unlike human feedback, AI-generated feedback offers better scalability, enabling faster
 149 and more cost-effective data collection. For example, RLAIF, introduced by Bai et al. (2022b),
 150 is a promising approach that trains reward models based on preferences generated by off-the-shelf
 151 LLMs. Variations of RLAIF have been explored in several studies (Lee et al., 2023; Jiao, 2023;
 152 Cui et al., 2024; Li et al., 2024; Yang et al., 2024). In the context of self-generating instructions,
 153 approaches like Self-Instruct (Wang et al., 2023a), Self-Refine (Madaan et al., 2023), and Self-
 154 Alignment (Li et al., 2023) demonstrate how models can autonomously generate datasets based on
 their learned human preferences.

155 Different alignment methods like RLHF and RLAIF rely on the RM to incorporate the human feed-
 156 back. The RM learns to predict human preference based on labeled examples, serving as a proxy
 157 for human judgment later. Therefore, the success of language model alignment relies heavily on the
 158 quality of the underlying reward model (Pace et al., 2024), which in turn dictates the behaviour of
 159 the resultant chatbot (Shen et al., 2023). Even in LLM inference, methods like best-of-N sampling
 160 use the RM to evaluate model outputs (Cui et al., 2024). RM has also become crucial for generating

²Before the process of human alignment, these models are referred to as “pre-trained” LLMs.

162

163

Table 1: Eye-tracking (ET) features computed per word.

164

165

Feature	Definition
First Fixation Duration (FFD)	Time spent on the initial fixation
Go-Past Time (GPT)	Cumulative fixation time before moving to the right
Total Reading Time (TRT)	Overall time spent fixating on a word
Number of Fixations (nFix)	Number of fixations on each word
Proportion of participant (fixProp)	Proportion of participants that fixated on the word

166

167

171 synthetic data for preference alignment. In recent RLAIF methods, reward modeling has expanded
 172 beyond its traditional role and is now used to generate artificial feedback.

173

174

2.1.1 REWARD MODELING

175

176

In the original implementation (Ouyang et al., 2022), the goal of RM training is to train a classifier
 that predicts the probability of human preference p^* between two completions (Equation 1), mod-
 elled by a Bradley-Terry model (Bradley & Terry, 1952). The typical setup involves showing two
 completions, with preferences being measured using win-loss-tie outcomes or a Likert scale to cap-
 ture the strength of preference (Bai et al., 2022a). The data is processed into prompt-chosen-rejected
 trios, where the chosen completion, y_w , is preferred over the rejected one, y_l , forming the basis for
 training (Ouyang et al., 2022).

182

183

$$p^*(y_w \succ y_l | x) = \frac{\exp(r^*(x, y_w))}{\exp(r^*(x, y_w)) + \exp(r^*(x, y_l))}. \quad (1)$$

184

185

186

187

2.2 EYE-TRACKING

188

189

Eye-tracking (ET) systems monitor oculomotor behavior, such as eye movements and fixations, of-
 fering valuable insights into visual attention, information processing, and expands our understanding
 of reading and language comprehension. (Zhang & Hollenstein, 2024). Specifically, ET data often
 include fixations – pauses in eye movement to focus on specific areas (Mathias et al., 2022); sac-
 cades – rapid movements between two points (McGuire & Tomuro, 2021); scanpaths – sequences of
 fixations that reveal saccades and regressions (Yang & Hollenstein, 2023); and other temporal and
 spatial gaze behavior features (Zhang & Hollenstein, 2024). Incorporating ET data into NLP tasks
 often involves the use of several features listed in Table 1.

190

191

192

193

194

195

While several publicly available datasets such as ZUCO (Hollenstein et al., 2020b), ZUCO2 (Hol-
 lenstein et al., 2018), PROVO (Luke & Christianson, 2018), ETSA-I (Mishra et al., 2016), ETSA-II
 (Mishra et al., 2018), GECO (Cop et al., 2017), GECO-MT (Colman et al., 2022) are widely used in
 ET research, obtaining real ET data for NLP tasks remains a challenge. This is primarily due to the
 cost and precision requirements of ET equipment, the unavailability of gaze data during inference, as
 well as privacy concerns (Khurana et al., 2023). To address these challenges, two main approaches
 have been proposed. The first involves integrating ET data into the model during training through
 methods like Multi-task learning (MTL), which eliminates the need for ET data during inference
 (Mishra et al., 2018; Klerke et al., 2016; Ren & Xiong, 2023; Yu et al., 2024; Deng et al., 2024).
 The second approach involves techniques that directly predict users’ gaze behaviour (Deng et al.,
 2024; 2023a; Zhang & Hollenstein, 2024; Wang et al., 2024a), creating synthetic ET data for any
 given text or stimulus (Deng et al., 2023b; Bolliger et al., 2023; Khurana et al., 2023; Li & Rudzicz,
 2021; Hollenstein et al., 2021; 2022).

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

3 RELATED WORK

212

213

214

215

Reward Modelling. The most popular approach to reward modeling follows the framework intro-
 duced by Ouyang et al. (2022). Several studies have examined alternative versions for refining RMs.
 For instance, Bai et al. (2022b) proposed more fine-grained reward structures, evaluating helpfulness
 and harmlessness separately. Other approaches have explored different reward modelling strategies
 (Wu et al., 2023b; Dong et al., 2023b; Wang et al., 2023c). Another line of research has focused

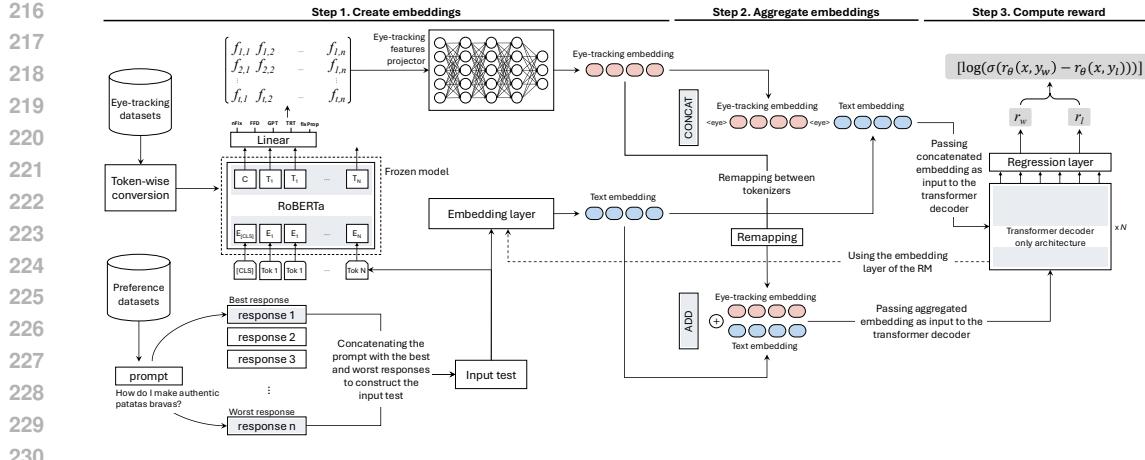


Figure 2: Overview of the **GazeReward** framework, incorporating eye-tracking features into the reward model. The architecture is illustrated in the figure using the second ET prediction model, but it would be identical if the first one were used instead (see subsection 5.1)

on Process Based Reward Models (PRMs) (Lightman et al., 2024; Uesato et al., 2022) which differ from conventional RMs by predicting the correctness of intermediate steps, rather than solely evaluating final outputs. Other studies implement data augmentation techniques (Shen et al., 2023), or cross-attention mechanisms between encoded input text and candidate pairs (Jiang et al., 2023b). Moreover, some works have leveraged synthetic preference data for reward modelling (Cui et al., 2024; Jiao, 2023). Wu et al. (2024b) built upon the LLM-as-a-Judge framework Zheng et al. (2023) by introducing LLM-as-a-Meta-Judge, which evaluates the model’s judgments to generate preference pairs that enhance its decision-making capabilities. Finally, Pace et al. (2024) incorporated a self-training approach to improve reward model training. However, to date, no research has explored the integration of ET or other implicit feedback signals into RM.

Eye-tracking in Natural Language Processing. Several studies have investigated the use of ET data for a variety of NLP tasks, such as named entity recognition (Hollenstein & Zhang, 2019; Ren & Xiong, 2023; Yu et al., 2024; Hollenstein et al., 2019), text comprehension (Ahn et al., 2020; Reich et al., 2022; Sood et al., 2020b), language modelling (Huang et al., 2023; Huang & Hollenstein, 2023; Deng et al., 2023b), and question answering (Zhang & Hollenstein, 2024; Wang et al., 2024a). Other applications include code comprehension (Alakmeh et al., 2024), code summarization (Zhang et al., 2024) and hallucination detection (Maharaj et al., 2023). Eye-tracking has also been applied to sentiment analysis and sarcasm detection tasks (Mishra et al., 2016; 2017; 2018; Barrett et al., 2018; Huang et al., 2023; Khurana et al., 2023; Hollenstein et al., 2019; Yang & Hollenstein, 2023; Kiegeland et al., 2024; Deng et al., 2023a; Mathias et al., 2018; McGuire & Tomuro, 2021). The most relevant work to our approach is by Kiegeland et al. (2024), which introduced a dataset generation method using ET signals for DPO, building on the controlled sentiment generation framework proposed by Deng et al. (2023a); Yang & Hollenstein (2023). While this study has contributed to the first steps towards integrating ET for human alignment in LLMs, it is task- and dataset-specific, often relying on ranking criteria that underutilize the potential of ET feedback. In contrast, our approach presents a more general framework by directly incorporating implicit feedback into the RM, rather than limiting its application to dataset creation.

4 GAZERWARD: REWARD MODELING WITH ET FEEDBACK

In this section, we discuss the proposed framework for augmenting the RM using implicit feedback derived from ET signals (Figure 2). Initially, we generate the ET features (subsection 4.1). For this step, we consider two state-of-the-art ET prediction models. Next, we combine the ET features with the text (subsection 4.2), producing different types of combined embeddings, and finally pass them as input into the RM to obtain the reward for the prompt and its corresponding response (subsection 4.3).

270
271
272
273
274
275
276
277

Table 2: Overview of different corpora used in the study to train the reward model.

Corpus	Train set	Val. set	Test set	Lang.	Reference
OASST1	6567	1160	416	EN*	Köpf et al. (2023)
HelpSteer2	5938	1049	364	EN	Wang et al. (2024c)

278 4.1 EYE-TRACKING FEATURES GENERATION

280 As discussed in subsection 2.2, obtaining organic ET features for NLP applications presents several
 281 challenges. In this work, we consider an approach inspired by RLAIF research, where feedback is
 282 artificially generated from pre-trained LLMs and, in particular, from ET prediction models. Specifi-
 283 cally, we incorporate the output of two different ET prediction models (Li & Rudzicz, 2021; Huang
 284 & Hollenstein, 2023) and evaluate the impact of different set of features. As input to these mod-
 285 els, we pass the same text as we do in the RM: a combination of prompt x and response y . The
 286 output is a set of ET features, denoted as f_{et} , for each token $f_{et} = \{f_1, f_2, \dots, f_w\} \in \mathbb{R}^{w \times f}$,
 287 where w represents the number of tokens in the tokenizer used by the ET prediction model, and f
 288 is the number of features. Depending on the specific model, between one and five synthetic features
 $f_{et} = \{f_1, f_2, \dots, f_w\} \in \mathbb{R}^{w \times f}$ are generated per token for the input text.

290 4.2 RM AUGMENTATION USING EYE-TRACKING FEATURES

292 We implement two different approaches for incorporating ET features into the RM, as shown in
 293 Figure 2. In the first approach, **GazeConcat**, we concatenate the ET embeddings with the text
 294 embeddings. In the second approach, **GazeAdd**, we add the ET embeddings to the text embeddings.
 295 Furthermore, we concatenate the prompt and the response to be evaluated and pass them through the
 296 pre-trained embedding layer of , to generate the embeddings $H = \{h_1, h_2, \dots, h_n\} \in \mathbb{R}^{n \times d}$, where
 297 n is the number of tokens in the tokenizer used by the RM and d is the model embedding size.

298 To project these features to the model embedding size (d), we use a Multilayer Perceptron (MLP)
 299 ET feature projector $fp()$. The $fp()$ consists of two linear layers, two dropout layers, two Layer
 300 Normalization layers, and ReLU activation, designed for stable, non-linear ET feature representation
 301 and overfitting prevention. The model’s input dimension dynamically adjusts to accommodate the
 302 number of features used during training. The ET features projector can be formulated as $emb_{ETF} =$
 303 $fp(f_{et}) \in \mathbb{R}^{w \times d}$ (Figure 2). This formula describes the projection of ETF features (f_{et}) through a
 304 function fp , resulting in an embedding matrix emb_{ETF} with dimensions $w \times d$, where w represents
 305 the number of time steps and d the embedding dimension.

306 **GazeConcat:** The ET embedding, denoted as emb_{ETF} , is concatenated with the text embedding
 307 H to form the input for the RM. To distinguish between the two modalities, we introduce two
 308 special tokens: $\langle eye \rangle$ and $\langle /eye \rangle$, which flag the start and end of the ET embedding, respectively
 309 (Figure 2). These special tokens are randomly initialized as one-dimensional vectors and added to
 310 the embedding layer or the RM model for training. The final input is structured as: $(emb(\langle eye \rangle) \circ$
 $emb_{ETF} \circ emb(\langle /eye \rangle) \circ H)$. The same process is applied to the attentions masks.

311 **GazeAdd:** The input to the RM consists of the ET embedding emb_{ETF} and the text embedding H ,
 312 which are added in an elementwise fashion: $(emb_{ETF} + H)$. The two ET prediction models use
 313 different tokenizers, which also differ from those used by the base models in the RM. As a result,
 314 the number of tokens n in the input for the RM and the number of tokens w generated by the ET
 315 prediction model may not match. To address this embedding alignment issue, and have the same
 316 dimension, we remap the ET features from the w -token space to the n -token space used by each
 317 base model in the RM. Further implementation details can be found in Appendix A.1.3.

319 4.3 REWARD MODEL

321 The RM’s architecture and hyperparameters are identical to those of the pretrained LLM, except
 322 that the classification head used for next-token prediction is replaced with a regression head that
 323 outputs a scalar reward (Touvron et al., 2023). This scalar reward indicates the quality of the model
 generation, corresponding to the predicted score for the final reply in a conversation. Differences in

these rewards represent the log-odds that one response is preferred over another. The loss function is defined in Equation 2, where y_w refers to the preferred response in a pair of completions y_w and y_l . The dataset D consists of human comparisons, where $r_\theta(x, y_w), r_\theta(x, y_l)$ represents the RM θ scalar outputs for the preferred and less preferred completions, respectively Ouyang et al. (2022).

$$\text{loss}(\theta) = -E_{(x, y_w, y_l) \sim D} [\log (\sigma(r_\theta(x, y_w) - r_\theta(x, y_l)))] \quad (2)$$

In the proposed method, we augment the traditional RM, which uses text input (a combination of the prompt x and response y), by incorporating (artificial) implicit feedback through ET features generated from the same text. These ET features provide valuable information for capturing human preferences, thereby improving the RM’s performance.

5 EXPERIMENTS

5.1 EXPERIMENTAL SETUP

Datasets. For our experiments, we use the OpenAssistant Conversations dataset’s (OASST1) (Köpf et al., 2023) and HelpSteer2 (Wang et al., 2024c) (Table 2). OASST1 is a human-generated, human-annotated, assistant-style conversation, created through global crowdsourcing and widely used for human alignment tasks (Köpf et al., 2023; Dettmers et al., 2023; Wu et al., 2024b). We filtered all non-English text, as the ET prediction models were exclusively trained on English data. Among the different responses in the dataset, we selected the two most distinct responses to compare the chosen and the rejected responses (Wang et al., 2024b). HelpSteer2 is a more recent, English-only dataset that has been used in studies such as Wang et al. (2024b;c). The dataset provides annotations for five response attributes: helpfulness, correctness, coherence, complexity, and verbosity. To transform it into a preference dataset, we designate the response with the higher helpfulness score as the chosen response and the other as the rejected response, following a method similar to that used in DPO training (Wang et al., 2024c) (see Appendix A.1.1 for more details about the datasets).

Dataset Preparation. To tune LLMs for human-AI interaction, we need to define a chat dialogue protocol that allows the model to understand human instructions and rate them. To this end, we adopt a chat protocol that utilizes special header and termination tokens, similar to the format used in Llama 3. For example, in the case of the Llama 3 8B model, the concatenation of a prompt and its corresponding response would follow this structure: <*im_start*>user Example Prompt <*im_end*> <*im_start*>assistant Example Response <*im_end*> (see Appendix A.1.1 for more details).

Models. As RM base models we use the pretrained checkpoint of Hugging Face ([Appendix A.1.4](#)) for Llama 3 8B, Llama 3 8B-instruct (Dubey et al., 2024) and Mistral 7B (Jiang et al., 2023a).

ET prediction models. In our analyses, we utilise two state-of-the-art ET prediction models, both pre-trained to predict ET features and kept frozen in our implementation. The input to these models is the same text used for the RM, with minimal modifications (Appendix A.1.2). The first model (Huang & Hollenstein, 2023), consists of a T5 embedding layer (Raffel et al., 2020), a two-layer BiLSTM (Hochreiter & Schmidhuber, 1997), and a one-hidden-layer MLP. This model was trained on the Dundee, GECO (Cop et al., 2017), ZuCo1 (Hollenstein et al., 2018), and ZuCo2 (Hollenstein et al., 2020a) datasets, and predicts total reading time (TRT) per token (Figure 3). The second model (Li & Rudzicz, 2021), is based on RoBERTa (Liu et al., 2019) with a regression head on each token. This head is a linear layer that outputs five features: FFD, fixProp, GPT, TRT, and nFix (Table 1). The model is initialized with pre-trained weights and fine-tuned on the ZuCo1 (Hollenstein et al., 2018), ZuCo2 (Hollenstein et al., 2020a) and PROVO (Luke & Christianson, 2018) datasets. Since RoBERTa’s maximum sequence length is 512 tokens and our input sequences are longer, we employ a sliding window approach. The input is split into 512-token segments with a 50-token overlap, and the results are combined using a linear weighted approach. Further details on these models and their integration into our framework are provided in Appendix A.1.2.

Baseline models. To evaluate the improvement in accuracy for a RM that incorporates implicit feedback, and specifically ET signals, we compare the same RM with and without the ET embeddings. For each dataset and model, we train and evaluate all combinations of integrating and combining ET features, and then compare them against a RM trained on the same base model and dataset but without implicit feedback.

378
 379 Table 3: Reward modeling accuracy (%) for OASST1 dataset. The highest results are in bold and
 380 the second highest are underlined.

		Llama-3-8B-Instruct	Llama-3-8B	Mistral-7B			
	baseline	64.6	diff(%)	58.6	diff(%)	62.0	diff(%)
GazeConcat	f_{comb1}	68.5	6.1%	68.5	17.0%	63.9	3.1%
	$f_{comb2.5}$	<u>70.7</u>	<u>9.4%</u>	70.4	20.3%	71.2	<u>14.7%</u>
	$f_{comb2.2}$	<u>68.5</u>	<u>6.1%</u>	73.1	24.8%	<u>72.6</u>	17.1%
GazeAdd	f_{comb1}	70.9	9.8%	71.2	<u>21.5%</u>	-	-
	$f_{comb2.5}$	51.9	-19.6%	<u>54.3</u>	<u>-7.2%</u>	-	-
	$f_{comb2.2}$	70.0	8.3%	69.0	17.8%	-	-

389
 390
 391 Table 4: Reward modeling accuracy (%) for Helpsteer2 dataset. The highest results are in bold and
 392 the second highest are underlined.

		Llama-3-8B-Instruct	Llama-3-8B	Mistral-7B			
	baseline	54.3	diff(%)	53.2	diff(%)	51.25	diff(%)
GazeConcat	f_{comb1}	62.8	15.7%	59.6	12.0%	63.3	23.6%
	$f_{comb2.5}$	<u>59.9</u>	<u>10.3%</u>	53.3	0.2%	53.0	3.5%
	$f_{comb2.2}$	59.3	9.2%	53.3	0.2%	<u>57.1</u>	<u>11.5%</u>
GazeAdd	f_{comb1}	63.2	16.3%	64.6	21.3%	-	-
	$f_{comb2.5}$	59.9	10.3%	58.8	10.5%	-	-
	$f_{comb2.2}$	59.1	8.7%	<u>59.6</u>	<u>12.1%</u>	-	-

402
 403
 404 **Evaluation metrics.** Performance is determined by measuring the model’s ability to predict the
 405 better response from pairs of replies with different ranks. Accuracy is calculated as the percentage of
 406 cases where the reward score for the preferred response is higher than that of the less preferred re-
 407 sponse, based on a held-out dataset. This method follows similar approaches found in Touvron et al.
 408 (2023); Yuan et al. (2023); Köpf et al. (2023); Cui et al. (2024). We use the test split proposed by the
 409 authors for each dataset (Table 2). We also conduct a complementary evaluation on RewardBench
 410 (Lambert et al., 2024), a benchmark dataset created for evaluating performance and safety features
 411 of RM’s (see Appendix B for more details).

412 **Training procedure.** In our implementation, the ET prediction model remains frozen while the
 413 ET features projector is trained (Figure 2). For the RM, we fine-tune the open-source models pre-
 414 viously introduced with QLoRA (Dettmers et al., 2023) a Parameter-Efficient Fine-Tuning (PEFT)
 415 method based on Low-Rank Adaptation (LoRA) (Hu et al., 2021), with other memory optimization
 416 techniques. We follow the training process for the RM as outlined in Touvron et al. (2023); Ouyang
 417 et al. (2022). We independently train each model on its respective dataset for two epochs, as detailed
 418 in Wang et al. (2024c). For hyperparameter tuning, we reserve 15% of each dataset for validation
 419 (shown in Table 2). The best-performing checkpoints are selected based on the lowest validation
 420 loss and used for performance evaluation. We perform a grid search to determine the optimal batch
 421 size and testing values of {8, 16, 32}. The AdamW optimizer (Loshchilov & Hutter, 2019) is used,
 422 with the Learning Rate (LR) is tuned over the range {1, 5, 10, 50}e-6, following the values reported
 423 in Touvron et al. (2023); Cui et al. (2024); Wang et al. (2024c). Additionally, we evaluate different
 424 LR schedulers: constant, linear, and cosine with a minimum LR. Further hyperparameter values and
 425 implementation details for both the RM and the ET projector can be found in Appendix C.

426 5.2 RESULTS

427
 428 The results of our experiments on the OASST and HelpSteer datasets, covering all possible combi-
 429 nations of ET features, models, and inclusion methods, as shown in Table 3 and Table 4 respectively.
 430 For the Mistral model, results for the **GazeAdd** method are unavailable due to the inability to map
 431 features between the ET prediction model’s tokenizer and the reward model’s tokenizer (details in
 432 Appendix A.1.3). In what follows, we discuss the main findings.

Effect of Model Initialization. We evaluate the impact of model initialization on performance. Open-access LLMs typically come in two forms: a pre-trained version without human alignment and a final version that has undergone alignment with human feedback. Since we lack access to intermediate checkpoints, we experiment with both pre-trained models (Mistral 7B and Llama 3) and models that are already human-aligned (Llama 3 Instruct). Our goal is to confirm that our method is effective for RM initialized with both pre-trained and human-aligned checkpoints. When comparing accuracy improvements relative to the baseline (without ET features), all models show considerable gains from incorporating implicit feedback. Notably, the Llama 3 8B and Mistral 7B models, which had no prior alignment, demonstrate over 20% performance improvement from the incorporation of ET features, indicating that unaligned models benefit more from implicit feedback.

Inclusion method. The results shown Table 3 and Table 4 indicate that both **GazeConcat** methods and **GazeAdd** introduce a substantial performance improvements to the RM. Across both datasets, concatenating embeddings (**GazeConcat**) delivers more consistent results. Incorporating ET information through specialized separator embeddings allows the model to process both text and ET features more robustly. However, in the HelpSteer dataset (Table 4), directly adding ET information to the text embeddings (**GazeAdd**) results in the greatest improvement over the baseline.

Eye-tracking (ET) feature importance. Different ET features capture distinct aspects of reading behaviour and information processing, influencing model performance uniquely (Zhang & Hollenstein, 2024). Here, we examine how model performance varies when incorporating three different feature combinations generated by two different ET prediction models: f_{comb_1} – TRT generated by the first ET prediction model; $f_{comb_{2.5}}$ – five features (nFix, FFD, GPT, TRT, fixProp) generated by the second ET prediction model; and $f_{comb_{2.2}}$ – TRT and FFD generated by the second ET prediction model. TRT and FFD are widely used in ET research (Huang et al., 2023; Huang & Hollenstein, 2023; Zhang & Hollenstein, 2024; Maharaj et al., 2023; Wang et al., 2022), and they have been shown to correlate with attention scores from pre-trained transformer models (Wang et al., 2024a; Bensemann et al., 2022; Sood et al., 2020a). Other studies have used gradient-based saliency methods to explore similar correlations (Hollenstein & Beinborn, 2021; Wu et al., 2024a). When comparing results, we observe that the RM benefits from implicit feedback regardless of the ET feature combination or ET prediction model used. Specifically, in most cases, f_{comb_1} yields the best results, particularly with the **GazeAdd** method. For **GazeConcat**, $f_{comb_{2.2}}$ performs best for the OASST1 dataset, while f_{comb_1} excels in HelpSteer2. We attribute the superior performance of f_{comb_1} to how the ET prediction model generating the fixations was trained, including the data and preprocessing methods used (see Appendix A.1.2). Moreover, when comparing $f_{comb_{2.2}}$ and $f_{comb_{2.5}}$ – both generated by the same model – only in one case does integrating nFix, GPT, and fixProp improves performance. In some instances, using $f_{comb_{2.5}}$ results in worse performance than the baseline, confirming findings provided by prior studies, which suggest that features related to reading time, such as FFD and TRT, contribute most to performance gains.

RewardBench. As a side contribution, we evaluate our best performing models (trained on the OAAST1 dataset) on RewardBench. This evaluation is not intended to directly compare our method with larger, more resource-intensive RM, but rather to show that through the integration of multi-modal signals like ET features we can significantly enhance the performance of RM models. The results shown in Table 5 demonstrate consistent improvements as previously observed, with gains exceeding 40% for the Mistral model – a notable gain considering that the base RM is the same. We note that the performance of the baseline models is impacted by RM trained on base models with less than 9B parameters and on relatively small datasets (see details in Appendix B).

6 DISCUSSION

In this work, we introduced a novel framework for integrating implicit feedback into the Reward Model, a key component for aligning LLMs and generating synthetic data for further alignment. We validated our approach using widely-adopted, open-source models such as Llama 3 and Mistral, for initializing the RM. By employing two different models to generate ET features, our results show that incorporating implicit feedback consistently improves the RM’s ability to predict user preferences, regardless of the model used and without the need to reach large parameter counts or train on massive datasets. Additionally, our method leverages ET features generated by models, making it fully scalable and applicable to various human alignment methods, including those that

486
 487 Table 5: Reward modeling accuracy (%) evaluating on RewardBench dataset. All models are trained
 488 on OASST1 dataset. The highest results are in bold and the second highest are underlined.

		Llama-3-8B-Instruct	Llama-3-8B	Mistral-7B			
	baseline	46.9	diff(%)	50.9	diff(%)	41.2	diff(%)
GazeConcat	f_{comb1}	57.8	23.1%	<u>58.4</u>	<u>14.5%</u>	59.9	45.4%
	$f_{comb2.5}$	58.4	24.4%	<u>58.1</u>	<u>14.1%</u>	60.3	46.4%
	$f_{comb2.2}$	<u>58.1</u>	<u>23.8%</u>	58.5	14.8%	60.5	46.9%
GazeAdd	f_{comb1}	56.5	20.3%	56.6	11.2%	-	-
	$f_{comb2.5}$	54.9	16.9%	53.8	5.6%	-	-
	$f_{comb2.2}$	55.4	17.9%	52.5	3.1%	-	-

497
 498
 499 involve artificially generated datasets. This work advances the ongoing discussion on optimizing
 500 AI alignment with human values and shows the potential of multimodal signals for NLP research
 501 **enhancing current methods.**

503 6.1 LIMITATIONS & FUTURE WORK

505
Data. A limitation of our study is that both ET prediction models were trained on a relatively
 506 small datasets (Appendix A.1.2) that are not tailored to our tasks. Future work could benefit from
 507 directly collecting ET data specifically for LLM-generated responses, to offer insights into human
 508 reading comprehension and information processing of prompts, which could further improve model
 509 performance. Additionally, since the ET prediction models used in our experiments were trained
 510 on English corpora, the method’s generalizability to other languages requires further investigation.
 511 Moreover, we explored two methodologies for integrating ET features into the RM, but other ap-
 512 proaches could prove more effective. For instance, ET features could be used to modify the RM’s
 513 attention mask, as suggested by Zhang & Hollenstein (2024). Regarding dataset selection, both
 514 models used, Mistral 7B (Jiang et al., 2023a) and Llama 3 (Dubey et al., 2024), were fine-tuned
 515 on publicly available data, though specific details on the datasets are limited. Therefore, we cannot
 516 discount the possibility that the datasets we used may have been part of the models’ pretraining,
 517 particularly for Llama 3 7B Instruct, which has already undergone human alignment. However, as
 518 we compare against baselines using the same model checkpoints, any potential effects would be
 519 consistent across both conditions. **Finally, it is important to note that although the ET prediction**
 520 **models remain fixed during training, our solution has a slightly higher number of parameters than**
 521 **the baseline.**

522 **Training.** The scaling trends for the RM (Touvron et al., 2023) show that larger models or models
 523 trained on massive datasets perform better. A promising direction would be to validate our frame-
 524 work on larger models, without relying on PEFT methods, and on larger datasets. **We are confident**
 525 **that, despite the considerable computational costs this may entail, our framework is capable of scal-**
 526 **ing effectively.** Another direction is integrating the proposed RM into an alignment method like
 527 RLHF, or applying it in rejection sampling to generate synthetic preference datasets, ensuring that
 528 accuracy gains in the RM translate to improvements in the final LLMs.

529 REPRODUCIBILITY STATEMENT

531 All the code necessary to reproduce this work will be released in a GitHub repository once published.
 532 Both datasets used are publicly available (Köpf et al., 2023; Wang et al., 2024c). Additionally, both
 533 ET prediction models have been trained with public datasets (Cop et al., 2017; Hollenstein et al.,
 534 2018; 2020a).

536 IMPACT STATEMENT

538 Since our research uses only synthetic ET data, there are no privacy concerns or need for large-scale
 539 experiments involving human subjects. **We should also raise attention to the limitations of human**
feedback and ET prediction models bias, that inevitably reflect aspects of their training data.

540 REFERENCES
541

- 542 Seoyoung Ahn, Conor Kelton, Aruna Balasubramanian, and Greg Zelinsky. Towards Predicting
543 Reading Comprehension From Gaze Behavior. In *ACM Symposium on Eye Tracking Research*
544 *and Applications*, ETRA '20 Short Papers, pp. 1–5, New York, NY, USA, June 2020. Association
545 for Computing Machinery. ISBN 978-1-4503-7134-6. doi: 10.1145/3379156.3391335. URL
546 <https://dl.acm.org/doi/10.1145/3379156.3391335>.
- 547 Tarek Alakmeh, David Reich, Lena Jäger, and Thomas Fritz. Predicting code comprehension: a
548 novel approach to align human gaze with code using deep neural networks. In *Alakmeh, Tarek;*
549 *Reich, David; Jäger, Lena; Fritz, Thomas (2024). Predicting code comprehension: a novel ap-*
550 *proach to align human gaze with code using deep neural networks. In: ACM International Con-*
551 *ference on the Foundations of Software Engineering, Porto de Galinhas, Brazil, 17 July 2024 -*
552 *19 July 2024., Porto de Galinhas, Brazil, July 2024. University of Zurich. doi: 10.1145/3660795.*
553 *URL* <https://www.zora.uzh.ch/id/eprint/260042/>. Issue: 1 Number: 1.
- 554 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Con-
555 crete Problems in AI Safety, July 2016. URL <http://arxiv.org/abs/1606.06565>.
556 arXiv:1606.06565 [cs].
- 557 Anthropic. Claude 3: Introducing the next generation of Claude, April 2024. URL <https://www.anthropic.com/news/clause-3-family>.
- 558 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn
559 Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson
560 Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez,
561 Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario
562 Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan.
563 Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback,
564 April 2022a. URL <http://arxiv.org/abs/2204.05862>. arXiv:2204.05862 [cs].
- 565 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones,
566 Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Ols-
567 son, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-
568 Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse,
569 Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mer-
570 cado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna
571 Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Con-
572 erly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario
573 Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional AI:
574 Harmlessness from AI Feedback, December 2022b. URL <http://arxiv.org/abs/2212.08073>.
575 arXiv:2212.08073 [cs].
- 576 Maria Barrett, Joachim Bingel, Nora Hollenstein, Marek Rei, and Anders Søgaard. Sequence Clas-
577 sification with Human Attention. In Anna Korhonen and Ivan Titov (eds.), *Proceedings of the*
578 *22nd Conference on Computational Natural Language Learning*, pp. 302–312, Brussels, Bel-
579 gium, October 2018. Association for Computational Linguistics. doi: 10.18653/v1/K18-1030.
580 URL <https://aclanthology.org/K18-1030>.
- 581 Joshua Bensemann, Alex Peng, Diana Benavides-Prado, Yang Chen, Nesan Tan, Paul Michael
582 Corballis, Patricia Riddle, and Michael Witbrock. Eye Gaze and Self-attention: How Hu-
583 mans and Transformers Attend Words in Sentences. In Emmanuele Chersoni, Nora Hollenstein,
584 Cassandra Jacobs, Yohei Oseki, Laurent Prévot, and Enrico Santus (eds.), *Proceedings of the*
585 *Workshop on Cognitive Modeling and Computational Linguistics*, pp. 75–87, Dublin, Ireland,
586 May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.cmcl-1.9. URL
587 <https://aclanthology.org/2022.cmcl-1.9>.
- 588 Lena Bolliger, David Reich, Patrick Haller, Deborah Jakobi, Paul Prasse, and Lena Jäger. ScanDL:
589 A Diffusion Model for Generating Synthetic Scanpaths on Texts. In Houda Bouamor, Juan Pino,
590 and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural*
591 *Language Processing*, pp. 15513–15538, Singapore, December 2023. Association for Computa-
592 tional Linguistics. doi: 10.18653/v1/2023.emnlp-main.960. URL <https://aclanthology.org/2023.emnlp-main.960>.

- 594 Ralph Allan Bradley and Milton E. Terry. Rank Analysis of Incomplete Block Designs: I. The
 595 Method of Paired Comparisons. *Biometrika*, 39(3/4):324–345, 1952. ISSN 0006-3444. doi: 10.
 596 2307/2334029. URL <https://www.jstor.org/stable/2334029>. Publisher: [Oxford
 597 University Press, Biometrika Trust].
 598
- 599 Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémie Scheurer, Javier
 600 Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, Tony Wang, Samuel
 601 Marks, Charbel-Raphaël Segerie, Micah Carroll, Andi Peng, Phillip Christoffersen, Mehul
 602 Damani, Stewart Slocum, Usman Anwar, Anand Siththanjan, Max Nadeau, Eric J. Michaud,
 603 Jacob Pfau, Dmitrii Krasheninnikov, Xin Chen, Lauro Langosco, Peter Hase, Erdem Büyik,
 604 Anca Dragan, David Krueger, Dorsa Sadigh, and Dylan Hadfield-Menell. Open Problems and
 605 Fundamental Limitations of Reinforcement Learning from Human Feedback, July 2023. URL
 606 <http://arxiv.org/abs/2307.15217>. arXiv:2307.15217 [cs].
 607
- 607 Toon Colman, Margot Fonteyne, Joke Daems, Nicolas Dirix, and Lieve Macken. GECO-MT:
 608 The Ghent Eye-tracking Corpus of Machine Translation. In Nicoletta Calzolari, Frédéric
 609 Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hi-
 610 toshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Jan Odijk, and Stelios Piperidis
 611 (eds.), *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pp. 29–
 612 38, Marseille, France, June 2022. European Language Resources Association. URL <https://aclanthology.org/2022.lrec-1.4>.
 613
- 614 Uschi Cop, Nicolas Dirix, Denis Drieghe, and Wouter Duyck. Presenting GECO: An eyetracking
 615 corpus of monolingual and bilingual sentence reading. *Behavior Research Methods*, 49(2):602–
 616 615, April 2017. ISSN 1554-3528. doi: 10.3758/s13428-016-0734-0. URL <https://doi.org/10.3758/s13428-016-0734-0>.
 617
- 618 Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong
 619 Xie, Ruobing Xie, Yankai Lin, Zhiyuan Liu, and Maosong Sun. UltraFeedback: Boosting Lan-
 620 guage Models with Scaled AI Feedback, July 2024. URL <http://arxiv.org/abs/2310.01377>. arXiv:2310.01377 [cs].
 621
- 622 Shuwen Deng, Paul Prasse, David Reich, Tobias Scheffer, and Lena Jäger. Pre-Trained Language
 623 Models Augmented with Synthetic Scanpaths for Natural Language Understanding. In Houda
 624 Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empir-
 625 ical Methods in Natural Language Processing*, pp. 6500–6507, Singapore, December 2023a.
 626 Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.400. URL
 627 <https://aclanthology.org/2023.emnlp-main.400>.
 628
- 628 Shuwen Deng, David R. Reich, Paul Prasse, Patrick Haller, Tobias Scheffer, and Lena A. Jäger.
 629 Eyettention: An Attention-based Dual-Sequence Model for Predicting Human Scanpaths during
 630 Reading. *Proc. ACM Hum.-Comput. Interact.*, 7(ETRA), May 2023b. doi: 10.1145/3591131.
 631 URL <https://doi.org/10.1145/3591131>. Place: New York, NY, USA Publisher: As-
 632 sociation for Computing Machinery.
 633
- 634 Shuwen Deng, Paul Prasse, David Reich, Tobias Scheffer, and Lena Jäger. Fine-Tuning Pre-Trained
 635 Language Models with Gaze Supervision. In Lun-Wei Ku, Andre Martins, and Vivek Srikanth
 636 (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics
 637 (Volume 2: Short Papers)*, pp. 217–224, Bangkok, Thailand, August 2024. Association for Com-
 638 putational Linguistics. URL <https://aclanthology.org/2024.acl-short.21>.
 639
- 640 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. QLoRA: Efficient Fine-
 641 tuning of Quantized LLMs. November 2023. URL [https://openreview.net/forum?id=OUIFPHEgJU&referrer=%5Bthe%20profile%20of%20Ari%20Holtzman%5D\(%2Fprofile%3Fid%3D~Ari_Holtzman1\)](https://openreview.net/forum?id=OUIFPHEgJU&referrer=%5Bthe%20profile%20of%20Ari%20Holtzman%5D(%2Fprofile%3Fid%3D~Ari_Holtzman1)).
 642
- 643 Hanze Dong, Wei Xiong, Deepanshu Goyal, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum,
 644 and Tong Zhang. RAFT: Reward rAnked FineTuning for Generative Foundation Model Align-
 645 ment, May 2023a. URL <http://arxiv.org/abs/2304.06767>. arXiv:2304.06767 [cs,
 646 stat].

- 648 Yi Dong, Zhilin Wang, Makesh Narsimhan Sreedhar, Xianchao Wu, and Oleksii Kuchaiev. SteerLM:
 649 Attribute Conditioned SFT as an (User-Steerable) Alternative to RLHF, October 2023b. URL
 650 <http://arxiv.org/abs/2310.05344>. arXiv:2310.05344 [cs].
 651
- 652 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
 653 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony
 654 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,
 655 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,
 656 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris
 657 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,
 658 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny
 659 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
 660 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael
 661 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-
 662 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah
 663 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan
 664 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-
 665 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy
 666 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,
 667 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-
 668 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,
 669 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der
 670 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,
 671 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-
 672 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova,
 673 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,
 674 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur
 675 Celebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-
 676 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
 677 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
 678 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-
 679 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,
 680 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,
 681 Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,
 682 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney
 683 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom,
 684 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,
 685 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-
 686 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang,
 687 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur,
 688 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre
 689 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha
 690 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,
 691 Ajay Menon, Ajay Sharma, Alex Boesenber, Alex Vaughan, Alexei Baevski, Allie Feinstein,
 692 Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples,
 693 Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco,
 694 Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman,
 695 Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto
 696 De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-
 697 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina
 698 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai,
 699 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,
 700 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana
 701 Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,
 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-
 702 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco
 703 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella
 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory
 704 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang,

- 702 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-
 703 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman,
 704 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer
 705 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe
 706 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie
 707 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun
 708 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal
 709 Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva,
 710 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian
 711 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson,
 712 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-
 713 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel
 714 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-
 715 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-
 716 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong,
 717 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,
 718 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,
 719 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,
 720 Rachel Rodriguez, Rafi Ayub, Raghatham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,
 721 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,
 722 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-
 723 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-
 724 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang
 725 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen
 726 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho,
 727 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser,
 728 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-
 729 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,
 730 Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu
 731 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-
 732 stable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu,
 733 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
 734 Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef
 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The Llama 3 Herd of Models, August
 2024. URL <http://arxiv.org/abs/2407.21783> [cs].
- 735 Google. Google Bard - Herramienta de IA Generativa y Bot Conversacional, 2023. URL <https://bard.google.com>.
 736
- 737 Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9
 738 (8):1735–1780, November 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL
 739 <https://doi.org/10.1162/neco.1997.9.8.1735>.
- 740 Nora Hollenstein and Lisa Beinborn. Relative Importance in Sentence Processing. In Chengqing
 741 Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting*
 742 *of the Association for Computational Linguistics and the 11th International Joint Conference*
 743 *on Natural Language Processing (Volume 2: Short Papers)*, pp. 141–150, Online, August 2021.
 744 Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-short.19. URL <https://aclanthology.org/2021.acl-short.19>.
- 745 Nora Hollenstein and Ce Zhang. Entity Recognition at First Sight: Improving NER with Eye Move-
 746 ment Information. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the*
 747 *2019 Conference of the North American Chapter of the Association for Computational Linguis-
 748 tics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 1–10, Minneapolis,
 749 Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1001.
 750 URL <https://aclanthology.org/N19-1001>.
- 751 Nora Hollenstein, Jonathan Rotsztejn, Marius Troendle, Andreas Pedroni, Ce Zhang, and Nicolas
 752 Langer. ZuCo, a simultaneous EEG and eye-tracking resource for natural sentence reading. *Sci-
 753 entific Data*, 5(1):180291, December 2018. ISSN 2052-4463. doi: 10.1038/sdata.2018.291. URL
 754 <https://doi.org/10.1038/sdata.2018.291>

- 756 <https://www.nature.com/articles/sdata2018291>. Publisher: Nature Publishing
 757 Group.
 758
- 759 Nora Hollenstein, Maria Barrett, Marius Troendle, Francesco Bigiolli, Nicolas Langer, and
 760 Ce Zhang. Advancing NLP with Cognitive Language Processing Signals, April 2019. URL
 761 <http://arxiv.org/abs/1904.02682>. arXiv:1904.02682 [cs].
- 762 Nora Hollenstein, Maria Barrett, and Lisa Beinborn. Towards Best Practices for Leveraging Hu-
 763 man Language Processing Signals for Natural Language Processing. In Emmanuele Chersoni,
 764 Barry Devereux, and Chu-Ren Huang (eds.), *Proceedings of the Second Workshop on Linguistic*
 765 *and Neurocognitive Resources*, pp. 15–27, Marseille, France, May 2020a. European Language
 766 Resources Association. ISBN 979-10-95546-52-8. URL <https://aclanthology.org/2020.lincr-1.3>.
- 767
- 768 Nora Hollenstein, Marius Troendle, Ce Zhang, and Nicolas Langer. ZuCo 2.0: A Dataset of Phys-
 769 iological Recordings During Natural Reading and Annotation. In Nicoletta Calzolari, Frédéric
 770 Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hi-
 771 toshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Asuncion Moreno, Jan Odijk, and
 772 Stelios Piperidis (eds.), *Proceedings of the Twelfth Language Resources and Evaluation Confer-
 773 ence*, pp. 138–146, Marseille, France, May 2020b. European Language Resources Association.
 774 ISBN 979-10-95546-34-4. URL <https://aclanthology.org/2020.lrec-1.18>.
- 775
- 776 Nora Hollenstein, Emmanuele Chersoni, Cassandra L. Jacobs, Yohei Oseki, Laurent Prévot, and
 777 Enrico Santus. CMCL 2021 Shared Task on Eye-Tracking Prediction. In Emmanuele Cher-
 778 soni, Nora Hollenstein, Cassandra Jacobs, Yohei Oseki, Laurent Prévot, and Enrico Santus (eds.),
 779 *Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics*, pp. 72–78,
 780 Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.cmcl-1.7.
 781 URL <https://aclanthology.org/2021.cmcl-1.7>.
- 782
- 783 Nora Hollenstein, Emmanuele Chersoni, Cassandra Jacobs, Yohei Oseki, Laurent Prévot, and En-
 784 rico Santus. CMCL 2022 Shared Task on Multilingual and Crosslingual Prediction of Human
 785 Reading Behavior. In Emmanuele Chersoni, Nora Hollenstein, Cassandra Jacobs, Yohei Oseki,
 786 Laurent Prévot, and Enrico Santus (eds.), *Proceedings of the Workshop on Cognitive Modeling*
 787 *and Computational Linguistics*, pp. 121–129, Dublin, Ireland, May 2022. Association for Com-
 788 putational Linguistics. doi: 10.18653/v1/2022.cmcl-1.14. URL <https://aclanthology.org/2022.cmcl-1.14>.
- 789
- 790 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
 791 and Weizhu Chen. LoRA: Low-Rank Adaptation of Large Language Models, October 2021. URL
 792 <http://arxiv.org/abs/2106.09685>. arXiv:2106.09685 [cs].
- 793
- 794 Xinting Huang and Nora Hollenstein. Long-Range Language Modeling with Selective Cache.
 795 In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Com-
 796 putational Linguistics: EMNLP 2023*, pp. 4838–4858, Singapore, December 2023. Associa-
 797 tion for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.321. URL <https://aclanthology.org/2023.findings-emnlp.321>.
- 798
- 799 Xinting Huang, Jiajing Wan, Ioannis Kritikos, and Nora Hollenstein. Longer Fixations, More Com-
 800 putation: Gaze-Guided Recurrent Neural Networks, October 2023. URL <http://arxiv.org/abs/2311.00159>. arXiv:2311.00159 [cs].
- 801
- 802 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chap-
 803 lot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,
 804 Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril,
 805 Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7B, October 2023a. URL
 806 <http://arxiv.org/abs/2310.06825>. arXiv:2310.06825 [cs].
- 807
- 808 Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. LLM-Blender: Ensembling Large Language Models
 809 with Pairwise Ranking and Generative Fusion, June 2023b. URL <http://arxiv.org/abs/2306.02561>. arXiv:2306.02561 [cs].

- 810 Tianhao Wu * Hanlin Zhu and Jiantao Banghua Zhu * Jiao, Evan Frick *. Starling-7B: Increasing
 811 LLM Helpfulness & Harmlessness with RLAIF, November 2023. URL <https://starling.cs.berkeley.edu>.
 812
 813 Alan Kennedy, Joël Pynte, Wayne Murray, and Shirley-Anne Paul. Frequency and predictability
 814 effects in the Dundee Corpus: An eye movement analysis. *Quarterly journal of experimental*
 815 *psychology* (2006), 66, March 2012. doi: 10.1080/17470218.2012.676054.
 816
 817 Varun Khurana, Yaman Kumar, Nora Hollenstein, Rajesh Kumar, and Balaji Krishnamurthy. Syn-
 818 thesizing Human Gaze Feedback for Improved NLP Performance. pp. 1895–1908, Dubrovnik,
 819 Croatia, May 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.
 820 eacl-main.139. URL <https://aclanthology.org/2023.eacl-main.139>.
 821
 822 Samuel Kiegeland, David Robert Reich, Ryan Cotterell, Lena Ann Jäger, and Ethan Wilcox. The
 823 Pupil Becomes the Master: Eye-Tracking Feedback for Tuning LLMs. July 2024. URL <https://openreview.net/forum?id=8oLUCBgKua>.
 824
 825 Chris L. Kleinke. Gaze and eye contact: A research review. *Psychological Bulletin*, 100(1):78–100,
 826 1986. ISSN 1939-1455. doi: 10.1037/0033-2909.100.1.78. Place: US Publisher: American
 827 Psychological Association.
 828
 829 Sigrid Klerke, Yoav Goldberg, and Anders Søgaard. Improving sentence compression by learn-
 830 ing to predict gaze. In Kevin Knight, Ani Nenkova, and Owen Rambow (eds.), *Proceed-
 831 ings of the 2016 Conference of the North American Chapter of the Association for Compu-
 832 tational Linguistics: Human Language Technologies*, pp. 1528–1533, San Diego, California,
 833 June 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-1179. URL
<https://aclanthology.org/N16-1179>.
 834
 835 Kristina Krasich, Robert McManus, Stephen Hutt, Myrthe Faber, Sidney K. D’Mello, and James R.
 836 Brockmole. Gaze-based signatures of mind wandering during real-world scene processing. *Jour-
 837 nal of Experimental Psychology. General*, 147(8):1111–1124, August 2018. ISSN 1939-2222.
 838 doi: 10.1037/xge0000411.
 839
 840 Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens,
 841 Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri,
 842 David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and
 843 Alexander Mattick. OpenAssistant Conversations – Democratizing Large Language Model Align-
 844 ment, April 2023. URL <http://arxiv.org/abs/2304.07327>. arXiv:2304.07327 [cs].
 845
 846 Nathan Lambert, Valentina Pyatkin, Jacob Morrison, L. J. Miranda, Bill Yuchen Lin, Khyathi
 847 Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Ha-
 848 jishirzi. RewardBench: Evaluating Reward Models for Language Modeling, March 2024. URL
<http://arxiv.org/abs/2403.13787>. arXiv:2403.13787 [cs].
 849
 850 M. F. Land and S. Furneaux. The knowledge base of the oculomotor system. *Philosophical Trans-
 851 actions of the Royal Society of London. Series B, Biological Sciences*, 352(1358):1231–1239,
 852 August 1997. ISSN 0962-8436. doi: 10.1098/rstb.1997.0105.
 853
 854 Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor
 855 Carbune, and Abhinav Rastogi. RLAIF: Scaling Reinforcement Learning from Human Feed-
 856 back with AI Feedback, September 2023. URL <http://arxiv.org/abs/2309.00267>.
 857 arXiv:2309.00267 [cs].
 858
 859 Ang Li, Qiugen Xiao, Peng Cao, Jian Tang, Yi Yuan, Zijie Zhao, Xiaoyuan Chen, Liang Zhang,
 860 Xiangyang Li, Kaitong Yang, Weidong Guo, Yukang Gan, Xu Yu, Daniell Wang, and Ying
 861 Shan. HRLAIF: Improvements in Helpfulness and Harmlessness in Open-domain Reinforcement
 862 Learning From AI Feedback, March 2024. URL <http://arxiv.org/abs/2403.08309>.
 863 arXiv:2403.08309 [cs].
 864
 865 Bai Li and Frank Rudzicz. TorontoCL at CMCL 2021 Shared Task: RoBERTa with Multi-Stage
 866 Fine-Tuning for Eye-Tracking Prediction. In Emmanuele Chersoni, Nora Hollenstein, Cassan-
 867 dra Jacobs, Yohei Oseki, Laurent Prévot, and Enrico Santus (eds.), *Proceedings of the Work-
 868 shop on Cognitive Modeling and Computational Linguistics*, pp. 85–89, Online, June 2021.

- 864 Association for Computational Linguistics. doi: 10.18653/v1/2021.cmcl-1.9. URL <https://aclanthology.org/2021.cmcl-1.9>.
- 865
- 866
- 867 Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Luke Zettlemoyer, Omer Levy, Jason Weston,
868 and Mike Lewis. Self-Alignment with Instruction Backtranslation, August 2023. URL <http://arxiv.org/abs/2308.06259>. arXiv:2308.06259 [cs].
- 869
- 870 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
871 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. LET'S VERIFY STEP BY STEP. 2024.
- 872
- 873 Tianqi Liu, Yao Zhao, Rishabh Joshi, Misha Khalman, Mohammad Saleh, Peter J. Liu, and Jialu Liu.
874 Statistical Rejection Sampling Improves Preference Optimization, January 2024. URL <http://arxiv.org/abs/2309.06657>. arXiv:2309.06657 [cs].
- 875
- 876 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy,
877 Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A Robustly Optimized
878 BERT Pretraining Approach, July 2019. URL <http://arxiv.org/abs/1907.11692>.
879 arXiv:1907.11692 [cs].
- 880
- 881 Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization, January 2019. URL
882 <http://arxiv.org/abs/1711.05101>. arXiv:1711.05101 [cs, math].
- 883
- 884 Steven G. Luke and Kiel Christianson. The Provo Corpus: A large eye-tracking corpus
885 with predictability norms. *Behavior Research Methods*, 50(2):826–833, April 2018. ISSN
886 1554-3528. doi: 10.3758/s13428-017-0908-4. URL <https://doi.org/10.3758/s13428-017-0908-4>.
- 887
- 888 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri
889 Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad
890 Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-
891 Refine: Iterative Refinement with Self-Feedback, May 2023. URL <http://arxiv.org/abs/2303.17651>. arXiv:2303.17651 [cs].
- 892
- 893 Kishan Maharaj, Ashita Saxena, Raja Kumar, Abhijit Mishra, and Pushpak Bhattacharyya. Eyes
894 Show the Way: Modelling Gaze Behaviour for Hallucination Detection. In Houda Bouamor, Juan
895 Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP*
896 2023, pp. 11424–11438, Singapore, December 2023. Association for Computational Linguistics.
897 doi: 10.18653/v1/2023.findings-emnlp.764. URL <https://aclanthology.org/2023.findings-emnlp.764>.
- 898
- 899 Sandeep Mathias, Diptesh Kanodia, Kevin Patel, Samarth Agrawal, Abhijit Mishra, and Pushpak
900 Bhattacharyya. Eyes are the Windows to the Soul: Predicting the Rating of Text Quality Using
901 Gaze Behaviour. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual
902 Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2352–
903 2362, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.
904 18653/v1/P18-1219. URL <https://aclanthology.org/P18-1219>.
- 905
- 906 Sandeep Mathias, Diptesh Kanodia, Abhijit Mishra, and Pushpak Bhattacharyya. A Survey on Using
907 Gaze Behaviour for Natural Language Processing, January 2022. URL <http://arxiv.org/abs/2112.15471>. arXiv:2112.15471 [cs].
- 908
- 909 Erik S. McGuire and Noriko Tomuro. Sentiment Analysis with Cognitive Attention Supervision.
910 *Proceedings of the Canadian Conference on Artificial Intelligence*, June 2021. doi: 10.21428/
911 594757db.90170c50. URL <https://caiac.pubpub.org/pub/kzxft4i8>.
- 912
- 913 Abhijit Mishra, Diptesh Kanodia, and Pushpak Bhattacharyya. Predicting Readers' Sarcasm Un-
914 derstandability by Modeling Gaze Behavior. *Proceedings of the AAAI Conference on Artifi-*
915 *cial Intelligence*, 30(1), March 2016. ISSN 2374-3468. doi: 10.1609/aaai.v30i1.9884. URL
916 <https://ojs.aaai.org/index.php/AAAI/article/view/9884>. Number: 1.
- 917
- 918 Abhijit Mishra, Diptesh Kanodia, Seema Nagar, Kuntal Dey, and Pushpak Bhattacharyya. Lever-
919 aging Cognitive Features for Sentiment Analysis, January 2017. URL <http://arxiv.org/abs/1701.05581>. arXiv:1701.05581 [cs].

- 918 Abhijit Mishra, Srikanth Tamilselvam, Riddhiman Dasgupta, Seema Nagar, and Kuntal Dey.
919 Cognition-Cognizant Sentiment Analysis With Multitask Subjectivity Summarization Based on
920 Annotators' Gaze Behavior. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
921 volume 32, April 2018. doi: 10.1609/aaai.v32i1.12068. URL <https://ojs.aaai.org/index.php/AAAI/article/view/12068>. ISSN: 2374-3468, 2159-5399 Issue: 1 Journal
922 Abbreviation: AAAI.
- 923
- 924 OpenAI. GPT-4 Technical Report, March 2023. URL <http://arxiv.org/abs/2303.08774>. arXiv:2303.08774 [cs].
- 925
- 926
- 927 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong
928 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-
929 ton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike,
930 and Ryan Lowe. Training language models to follow instructions with human feedback, March
931 2022. URL <http://arxiv.org/abs/2203.02155>. arXiv:2203.02155 [cs].
- 932
- 933 Alizée Pace, Jonathan Mallinson, Eric Malmi, Sebastian Krause, and Aliaksei Severyn. West-of-N:
934 Synthetic Preference Generation for Improved Reward Modeling, January 2024. URL <http://arxiv.org/abs/2401.12086>. arXiv:2401.12086 [cs].
- 935
- 936 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and
937 Chelsea Finn. Direct Preference Optimization: Your Language Model is Secretly a Reward
938 Model, December 2023. URL <http://arxiv.org/abs/2305.18290>. arXiv:2305.18290
939 [cs].
- 940
- 941 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
942 Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the Limits of Transfer Learning with a Uni-
943 fied Text-to-Text Transformer, July 2020. URL <http://arxiv.org/abs/1910.10683>.
944 arXiv:1910.10683 [cs, stat].
- 945
- 946 David Robert Reich, Paul Prasse, Chiara Tschirner, Patrick Haller, Frank Goldhammer, and Lena A.
947 Jäger. Inferring Native and Non-Native Human Reading Comprehension and Subjective Text
948 Difficulty from Scanpaths in Reading. In *2022 Symposium on Eye Tracking Research and*
949 *Applications*, ETRA '22, pp. 1–8, New York, NY, USA, June 2022. Association for Comput-
950 ing Machinery. ISBN 978-1-4503-9252-5. doi: 10.1145/3517031.3529639. URL <https://doi.org/10.1145/3517031.3529639>.
- 951
- 952 Yuqi Ren and Deyi Xiong. CogAlign: Learning to Align Textual Neural Representations to Cog-
953 nitive Language Processing Signals, November 2023. URL <http://arxiv.org/abs/2106.05544>. arXiv:2106.05544 [cs].
- 954
- 955 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Pol-
956 icy Optimization Algorithms, August 2017. URL <http://arxiv.org/abs/1707.06347>.
957 arXiv:1707.06347 [cs].
- 958
- 959 Lingfeng Shen, Sihao Chen, Linfeng Song, Lifeng Jin, Baolin Peng, Haitao Mi, Daniel Khashabi,
960 and Dong Yu. The Trickle-down Impact of Reward Inconsistency on RLHF. October 2023. URL
961 <https://openreview.net/forum?id=MeHmwCDifc>.
- 962
- 963 Ekta Sood, Simon Tannert, Diego Frassinelli, Andreas Bulling, and Ngoc Thang Vu. Interpreting
964 Attention Models with Human Visual Attention in Machine Reading Comprehension. In Raquel
965 Fernández and Tal Linzen (eds.), *Proceedings of the 24th Conference on Computational Natural*
966 *Language Learning*, pp. 12–25, Online, November 2020a. Association for Computational Lin-
967 guistics. doi: 10.18653/v1/2020.conll-1.2. URL <https://aclanthology.org/2020.conll-1.2>.
- 968
- 969 Ekta Sood, Simon Tannert, Philipp Mueller, and Andreas Bulling. Improving Natural
970 Language Processing Tasks with Human Gaze-Guided Neural Attention. In *Advances*
971 *in Neural Information Processing Systems*, volume 33, pp. 6327–6341. Curran Asso-
972 ciates, Inc., 2020b. URL <https://proceedings.neurips.cc/paper/2020/hash/460191c72f67e90150a093b4585e7eb4-Abstract.html>.

972 Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,
 973 Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh Mariooryad, Yifan Ding, Xinyang Geng,
 974 Fred Alcober, Roy Frostig, Mark Omernick, Lexi Walker, Cosmin Paduraru, Christina Sorokin,
 975 Andrea Tacchetti, Colin Gaffney, Samira Daruki, Olcan Sercinoglu, Zach Gleicher, Juliette Love,
 976 Paul Voigtlaender, Rohan Jain, Gabriela Surita, Kareem Mohamed, Rory Blevins, Junwhan Ahn,
 977 Tao Zhu, Kornraphop Kawintiranon, Orhan Firat, Yiming Gu, Yujing Zhang, Matthew Rahtz,
 978 Manaal Faruqui, Natalie Clay, Justin Gilmer, J. D. Co-Reyes, Ivo Penchev, Rui Zhu, Nobuyuki
 979 Morioka, Kevin Hui, Krishna Haridasan, Victor Campos, Mahdis Mahdieh, Mandy Guo, Sameer
 980 Hassan, Kevin Kilgour, Arpi Vezer, Heng-Tze Cheng, Raoul de Liedekerke, Siddharth Goyal,
 981 Paul Barham, D. J. Strouse, Seb Noury, Jonas Adler, Mukund Sundararajan, Sharad Vikram,
 982 Dmitry Lepikhin, Michela Paganini, Xavier Garcia, Fan Yang, Dasha Valter, Maja Trebacz, Kiran
 983 Vodrahalli, Chulayuth Asawaroengchai, Roman Ring, Norbert Kalb, Livio Baldini Soares, Sid-
 984 dhartha Brahma, David Steiner, Tianhe Yu, Fabian Mentzer, Antoine He, Lucas Gonzalez, Bibo
 985 Xu, Raphael Lopez Kaufman, Laurent El Shafey, Junhyuk Oh, Tom Hennigan, George van den
 986 Driessche, Seth Odoom, Mario Lucic, Becca Roelofs, Sid Lall, Amit Marathe, Betty Chan, San-
 987 tiago Ontanon, Luheng He, Denis Teplyashin, Jonathan Lai, Phil Crone, Bogdan Damoc, Lewis
 988 Ho, Sebastian Riedel, Karel Lenc, Chih-Kuan Yeh, Aakanksha Chowdhery, Yang Xu, Mehran
 989 Kazemi, Ehsan Amid, Anastasia Petrushkina, Kevin Swersky, Ali Khodaei, Gowoon Chen, Chris
 990 Larkin, Mario Pinto, Geng Yan, Adria Puigdomenech Badia, Piyush Patil, Steven Hansen, Dave
 991 Orr, Sebastien M. R. Arnold, Jordan Grimstad, Andrew Dai, Sholto Douglas, Rishika Sinha, Vikas
 992 Yadav, Xi Chen, Elena Gribovskaya, Jacob Austin, Jeffrey Zhao, Kaushal Patel, Paul Komarek,
 993 Sophia Austin, Sebastian Borgeaud, Linda Friso, Abhimanyu Goyal, Ben Caine, Kris Cao, Da-
 994 Woon Chung, Matthew Lamm, Gabe Barth-Maron, Thais Kagohara, Kate Olszewska, Mia Chen,
 995 Kaushik Shivakumar, Rishabh Agarwal, Harshal Godhia, Ravi Rajwar, Javier Snaider, Xerxes
 996 Dotiwalla, Yuan Liu, Aditya Barua, Victor Ungureanu, Yuan Zhang, Bat-Orgil Batsaikhan, Ma-
 997 teo Wirth, James Qin, Ivo Danihelka, Tulsee Doshi, Martin Chadwick, Jilin Chen, Sanil Jain,
 998 Quoc Le, Arjun Kar, Madhu Gurumurthy, Cheng Li, Ruoxin Sang, Fangyu Liu, Lampros Lam-
 999 prou, Rich Munoz, Nathan Lintz, Harsh Mehta, Heidi Howard, Malcolm Reynolds, Lora Aroyo,
 1000 Quan Wang, Lorenzo Blanco, Albin Cassirer, Jordan Griffith, Dipanjan Das, Stephan Lee, Jakub
 1001 Sygnowski, Zach Fisher, James Besley, Richard Powell, Zafarali Ahmed, Dominik Paulus, David
 1002 Reitter, Zalan Borsos, Rishabh Joshi, Aedan Pope, Steven Hand, Vittorio Selo, Vihan Jain, Nikhil
 1003 Sethi, Megha Goel, Takaki Makino, Rhys May, Zhen Yang, Johan Schalkwyk, Christina But-
 1004 terfield, Anja Hauth, Alex Goldin, Will Hawkins, Evan Senter, Sergey Brin, Oliver Woodman,
 1005 Marvin Ritter, Eric Noland, Minh Giang, Vijay Bolina, Lisa Lee, Tim Blyth, Ian Mackinnon,
 1006 Machel Reid, Obaid Sarvana, David Silver, Alexander Chen, Lily Wang, Loren Maggiore, Os-
 1007 car Chang, Nithya Attaluri, Gregory Thornton, Chung-Cheng Chiu, Oskar Bunyan, Nir Levine,
 1008 Timothy Chung, Evgenii Eltyshev, Xiance Si, Timothy Lillicrap, Demetra Brady, Vaibhav Ag-
 1009 garwal, Boxi Wu, Yuanzhong Xu, Ross McIlroy, Kartikeya Badola, Paramjit Sandhu, Erica Mor-
 1010 eira, Wojciech Stokowiec, Ross Hemsley, Dong Li, Alex Tudor, Pranav Shyam, Elahe Rahim-
 1011 toroghi, Salem Haykal, Pablo Sprechmann, Xiang Zhou, Diana Mincu, Yujia Li, Ravi Addanki,
 1012 Kalpesh Krishna, Xiao Wu, Alexandre Frechette, Matan Eyal, Allan Dafoe, Dave Lace, Jay
 1013 Whang, Thi Avrahami, Ye Zhang, Emanuel Taropa, Hanzhao Lin, Daniel Toyama, Eliza Ruther-
 1014 ford, Motoki Sano, HyunJeong Choe, Alex Tomala, Chalence Safranek-Shrader, Nora Kassner,
 1015 Mantas Pajarskas, Matt Harvey, Sean Sechrist, Meire Fortunato, Christina Lyu, Gamaleldin El-
 1016 sayed, Chenkai Kuang, James Lottes, Eric Chu, Chao Jia, Chih-Wei Chen, Peter Humphreys,
 1017 Kate Baumli, Connie Tao, Rajkumar Samuel, Cicero Nogueira dos Santos, Anders Andreassen,
 1018 Nemanja Rakicevic, Dominik Grewe, Aviral Kumar, Stephanie Winkler, Jonathan Caton, Andrew
 1019 Brock, Sid Dalmia, Hannah Sheahan, Iain Barr, Yingjie Miao, Paul Natsev, Jacob Devlin, Feryal
 1020 Behbahani, Flavien Prost, Yanhua Sun, Artiom Myaskovsky, Thanumalayan Sankaranarayana
 1021 Pillai, Dan Hurt, Angeliki Lazaridou, Xi Xiong, Ce Zheng, Fabio Pardo, Xiaowei Li, Dan Hor-
 1022 gan, Joe Stanton, Moran Ambar, Fei Xia, Alejandro Lince, Mingqiu Wang, Basil Mustafa, Albert
 1023 Webson, Hyo Lee, Rohan Anil, Martin Wicke, Timothy Dozat, Abhishek Sinha, Enrique Pi-
 1024 queras, Elahe Dabir, Shyam Upadhyay, Anudhyan Boral, Lisa Anne Hendricks, Corey Fry, Josip
 1025 Djolonga, Yi Su, Jake Walker, Jane Labanowski, Ronny Huang, Vedant Misra, Jeremy Chen, R. J.
 Skerry-Ryan, Avi Singh, Shruti Rijhwani, Dian Yu, Alex Castro-Ros, Beer Changpinyo, Romina
 Datta, Sumit Bagri, Arnar Mar Hrafnelsson, Marcello Maggioni, Daniel Zheng, Yury Sulsky,
 Shaobo Hou, Tom Le Paine, Antoine Yang, Jason Riesa, Dominika Rogozinska, Dror Marcus,
 Dalia El Badawy, Qiao Zhang, Luyu Wang, Helen Miller, Jeremy Greer, Lars Lowe Sjos, Azade
 Nova, Heiga Zen, Rahma Chaabouni, Mihaela Rosca, Jiepu Jiang, Charlie Chen, Ruibo Liu, Tara

1026 Sainath, Maxim Krikun, Alex Polozov, Jean-Baptiste Lespiau, Josh Newlan, Zeynep Cankara,
 1027 Soo Kwak, Yunhan Xu, Phil Chen, Andy Coenen, Clemens Meyer, Katerina Tsihlas, Ada Ma, Ju-
 1028 raj Gottweis, Jinwei Xing, Chenjie Gu, Jin Miao, Christian Frank, Zeynep Cankara, Sanjay Gana-
 1029 pathy, Ishita Dasgupta, Steph Hughes-Fitt, Heng Chen, David Reid, Keran Rong, Hongmin Fan,
 1030 Joost van Amersfoort, Vincent Zhuang, Aaron Cohen, Shixiang Shane Gu, Anhad Mohananey,
 1031 Anastasija Ilic, Taylor Tobin, John Wieting, Anna Bortsova, Phoebe Thacker, Emma Wang, Emily
 1032 Caveness, Justin Chiu, Eren Sezener, Alex Kaskasoli, Steven Baker, Katie Millican, Mohamed
 1033 Elhwaty, Kostas Aisopos, Carl Lebsack, Nathan Byrd, Hanjun Dai, Wenhao Jia, Matthew Wi-
 1034 ethoff, Elnaz Davoodi, Albert Weston, Lakshman Yagati, Arun Ahuja, Isabel Gao, Golan Pundak,
 1035 Susan Zhang, Michael Azzam, Khe Chai Sim, Sergi Caelles, James Keeling, Abhanshu Sharma,
 1036 Andy Swing, YaGuang Li, Chenxi Liu, Carrie Grimes Bostock, Yamini Bansal, Zachary Nado,
 1037 Ankesh Anand, Josh Lipschultz, Abhijit Karmarkar, Lev Proleev, Abe Ittycheriah, Soheil Has-
 1038 sas Yeganeh, George Polovets, Aleksandra Faust, Jiao Sun, Alban Rustemi, Pen Li, Rakesh
 1039 Shivanna, Jeremiah Liu, Chris Welty, Federico Lebron, Anirudh Baddepudi, Sebastian Krause,
 1040 Emilio Parisotto, Radu Soricut, Zheng Xu, Dawn Bloxwich, Melvin Johnson, Behnam Neyshabur,
 1041 Justin Mao-Jones, Renshen Wang, Vinay Ramasesh, Zaheer Abbas, Arthur Guez, Constant Segal,
 1042 Duc Dung Nguyen, James Svensson, Le Hou, Sarah York, Kieran Milan, Sophie Bridgers, Wiktor
 1043 Gworek, Marco Tagliasacchi, James Lee-Thorp, Michael Chang, Alexey Guseynov, Ale Jakse
 1044 Hartman, Michael Kwong, Ruizhe Zhao, Sheleem Kashem, Elizabeth Cole, Antoine Miech,
 1045 Richard Tanburn, Mary Phuong, Filip Pavetic, Sebastien Cevey, Ramona Comanescu, Richard
 1046 Ives, Sherry Yang, Cosmo Du, Bo Li, Zizhao Zhang, Mariko Iinuma, Clara Huiyi Hu, Aurko Roy,
 1047 Shaan Bijwadia, Zhenkai Zhu, Danilo Martins, Rachel Saputro, Anita Gergely, Steven Zheng,
 1048 Dawei Jia, Ioannis Antonoglou, Adam Sadovsky, Shane Gu, Yingying Bi, Alek Andreev, Sina
 1049 Samangooei, Mina Khan, Tomas Kociský, Angelos Filos, Chintu Kumar, Colton Bishop, Adams
 1050 Yu, Sarah Hodgkinson, Sid Mittal, Premal Shah, Alexandre Moufarek, Yong Cheng, Adam Blo-
 1051 niarz, Jaehoon Lee, Pedram Pejman, Paul Michel, Stephen Spencer, Vladimir Feinberg, Xuehan
 1052 Xiong, Nikolay Savinov, Charlotte Smith, Siamak Shakeri, Dustin Tran, Mary Chesus, Bernd
 1053 Bohnet, George Tucker, Tamara von Glehn, Carrie Muir, Yiran Mao, Hideto Kazawa, Ambrose
 1054 Sloane, Kedar Soparkar, Disha Srivastava, James Cobon-Kerr, Michael Sharman, Jay Pavagadhi,
 1055 Carlos Araya, Karolis Misiunas, Nimesh Ghelani, Michael Laskin, David Barker, Qiuja Li, An-
 1056 ton Briukhov, Neil Houldsby, Mia Glaese, Balaji Lakshminarayanan, Nathan Schucher, Yunhao
 1057 Tang, Eli Collins, Hyeontaek Lim, Fangxiaoyu Feng, Adria Recasens, Guangda Lai, Alberto
 1058 Magni, Nicola De Cao, Aditya Siddhant, Zoe Ashwood, Jordi Orbay, Mostafa Dehghani, Jenny
 1059 Brennan, Yifan He, Kelvin Xu, Yang Gao, Carl Saroufim, James Molloy, Xinyi Wu, Seb Arnold,
 1060 Solomon Chang, Julian Schrittweiser, Elena Buchatskaya, Soroush Radpour, Martin Polacek,
 1061 Skye Giordano, Ankur Bapna, Simon Tokumine, Vincent Hellendoorn, Thibault Sottiaux, Sarah
 1062 Cogan, Aliaksei Severyn, Mohammad Saleh, Shantanu Thakoor, Laurent Shefey, Siyuan Qiao,
 1063 Meenu Gaba, Shuo-yiin Chang, Craig Swanson, Biao Zhang, Benjamin Lee, Paul Kishan Ruben-
 1064 stein, Gan Song, Tom Kwiatkowski, Anna Koop, Ajay Kannan, David Kao, Parker Schuh, Axel
 1065 Stjerngren, Golnaz Ghiasi, Gena Gibson, Luke Vilnis, Ye Yuan, Felipe Tiengo Ferreira, Aish-
 1066 warya Kamath, Ted Klimentko, Ken Franko, Kefan Xiao, Indro Bhattacharya, Miteyan Patel, Rui
 1067 Wang, Alex Morris, Robin Strudel, Vivek Sharma, Peter Choy, Sayed Hadi Hashemi, Jessica
 1068 Landon, Mara Finkelstein, Priya Jhakra, Justin Frye, Megan Barnes, Matthew Mauger, Dennis
 1069 Daun, Khuslen Baatarsukh, Matthew Tung, Wael Farhan, Henryk Michalewski, Fabio Viola,
 1070 Felix de Chaumont Quirly, Charline Le Lan, Tom Hudson, Qingze Wang, Felix Fischer, Ivy
 1071 Zheng, Elspeth White, Anca Dragan, Jean-baptiste Alayrac, Eric Ni, Alexander Pritzel, Adam
 1072 Iwanicki, Michael Isard, Anna Bulanova, Lukas Zilka, Ethan Dyer, Devendra Sachan, Srivatsan
 1073 Srinivasan, Hannah Muckenhirn, Honglong Cai, Amol Mandhane, Mukarram Tariq, Jack W. Rae,
 1074 Gary Wang, Kareem Ayoub, Nicholas FitzGerald, Yao Zhao, Woohyun Han, Chris Alberti, Dan
 1075 Garrette, Kashyap Krishnakumar, Mai Gimenez, Anselm Levskaya, Daniel Sohn, Josip Matač,
 1076 Inaki Iturratxe, Michael B. Chang, Jackie Xiang, Yuan Cao, Nishant Ranka, Geoff Brown, Adrian
 1077 Hutter, Vahab Mirrokni, Nanxin Chen, Kaisheng Yao, Zoltan Egyed, Francois Galilee, Tyler
 1078 Liechty, Praveen Kallakuri, Evan Palmer, Sanjay Ghemawat, Jasmine Liu, David Tao, Chloe
 1079 Thornton, Tim Green, Mimi Jasarevic, Sharon Lin, Victor Cotruta, Yi-Xuan Tan, Noah Fiedel,
 Hongkun Yu, Ed Chi, Alexander Neitz, Jens Heitkaemper, Anu Sinha, Denny Zhou, Yi Sun,
 Charbel Kaed, Brice Hulse, Swaroop Mishra, Maria Georgaki, Sneha Kudugunta, Clement Fara-
 bet, Izhak Shafran, Daniel Vlasic, Anton Tsitsulin, Rajagopal Ananthanarayanan, Alen Carin,
 Guolong Su, Pei Sun, Shashank V, Gabriel Carvajal, Josef Broder, Iulia Comsa, Alena Repina,
 William Wong, Warren Weilun Chen, Peter Hawkins, Egor Filonov, Lucia Loher, Christoph

Hirnschall, Weiyi Wang, Jingchen Ye, Andrea Burns, Hardie Cate, Diana Gage Wright, Federico Piccinini, Lei Zhang, Chu-Cheng Lin, Ionel Gog, Yana Kulizhskaya, Ashwin Sreevatsa, Shuang Song, Luis C. Cobo, Anand Iyer, Chetan Tekur, Guillermo Garrido, Zhuyun Xiao, Rupert Kemp, Huaixiu Steven Zheng, Hui Li, Ananth Agarwal, Christel Ngani, Kati Goshvadi, Rebeca Santamaria-Fernandez, Wojciech Fica, Xinyun Chen, Chris Gorgolewski, Sean Sun, Roopal Garg, Xinyu Ye, S. M. Ali Eslami, Nan Hua, Jon Simon, Pratik Joshi, Yelin Kim, Ian Tenney, Sahitya Potluri, Lam Nguyen Thiet, Quan Yuan, Florian Luisier, Alexandra Chronopoulou, Salvatore Scellato, Praveen Srinivasan, Minmin Chen, Vinod Koverkathu, Valentin Dalibard, Yaming Xu, Brennan Saeta, Keith Anderson, Thibault Sellam, Nick Fernando, Fantine Huot, Junehyuk Jung, Mani Varadarajan, Michael Quinn, Amit Raul, Maigo Le, Ruslan Habalov, Jon Clark, Komal Jalan, Kalesha Bullard, Achintya Singhal, Thang Luong, Boyu Wang, Sujeevan Rajayogam, Julian Eisenschlos, Johnson Jia, Daniel Finchelstein, Alex Yakubovich, Daniel Balle, Michael Fink, Sameer Agarwal, Jing Li, Dj Dvijotham, Shalini Pal, Kai Kang, Jaclyn Konzelmann, Jennifer Beattie, Olivier Dousse, Diane Wu, Remi Crocker, Chen Elkind, Siddhartha Reddy Jonnalagadda, Jong Lee, Dan Holtmann-Rice, Krystal Kallarackal, Rosanne Liu, Denis Vnukov, Neera Vats, Luca Invernizzi, Mohsen Jafari, Huanjie Zhou, Lilly Taylor, Jennifer Prendki, Marcus Wu, Tom Eccles, Tianqi Liu, Kavya Kopparapu, Francoise Beaufays, Christof Angermueller, Andreea Marzoca, Shourya Sarcar, Hilal Dib, Jeff Stanway, Frank Perbet, Nejc Trdin, Rachel Sterneck, Andrey Khorlin, Dinghua Li, Xihui Wu, Sonam Goenka, David Madras, Sasha Goldstein, Willi Gierke, Tong Zhou, Yaxin Liu, Yannie Liang, Anais White, Yunjie Li, Shreya Singh, Sanaz Bahargam, Mark Epstein, Sujoy Basu, Li Lao, Adnan Ozturel, Carl Crous, Alex Zhai, Han Lu, Zora Tung, Neeraj Gaur, Alanna Walton, Lucas Dixon, Ming Zhang, Amir Globerzon, Grant Uy, Andrew Bolt, Olivia Wiles, Milad Nasr, Ilia Shumailov, Marco Selvi, Francesco Piccinno, Ricardo Aguilar, Sara McCarthy, Misha Khalman, Mrinal Shukla, Vlado Galic, John Carpenter, Kevin Villela, Haibin Zhang, Harry Richardson, James Martens, Matko Bosnjak, Shreyas Rammohan Belle, Jeff Seibert, Mahmoud Alnahlawi, Brian McWilliams, Sankalp Singh, Annie Louis, Wen Ding, Dan Popovici, Lenin Simicich, Laura Knight, Pulkit Mehta, Nishesh Gupta, Chongyang Shi, Saaber Fatehi, Jovana Mitrovic, Alex Grills, Joseph Pagadora, Dessie Petrova, Danielle Eisenbud, Zhishuai Zhang, Damion Yates, Bhavishya Mittal, Nilesh Tripuraneni, Yannis Assael, Thomas Brovelli, Prateek Jain, Mihajlo Velimirovic, Canfer Akbulut, Jiaqi Mu, Wolfgang Macherey, Ravin Kumar, Jun Xu, Haroon Qureshi, Gheorghe Comanici, Jeremy Wiesner, Zhitao Gong, Anton Ruddock, Matthias Bauer, Nick Felt, Anirudh GP, Anurag Arnab, Dustin Zelle, Jonas Rothfuss, Bill Rosgen, Ashish Shenoy, Bryan Seybold, Xinjian Li, Jayaram Mudigonda, Goker Erdogan, Jiawei Xia, Jiri Simsa, Andrea Michi, Yi Yao, Christopher Yew, Steven Kan, Isaac Caswell, Carey Radebaugh, Andre Elisseeff, Pedro Valenzuela, Kay McKinney, Kim Paterson, Albert Cui, Eri Latorre-Chimoto, Solomon Kim, William Zeng, Ken Durden, Priya Ponnappalli, Tiberiu Sosea, Christopher A. Choquette-Choo, James Manyika, Brona Robenek, Harsha Vashisht, Sebastien Pereira, Hoi Lam, Marko Velic, Denese Owusu-Afriyie, Katherine Lee, Tolga Bolukbasi, Alicia Parrish, Shawn Lu, Jane Park, Balaji Venkatraman, Alice Talbert, Lambert Rosique, Yuchung Cheng, Andrei Sozanschi, Adam Paszke, Praveen Kumar, Jessica Austin, Lu Li, Khalid Salama, Wooyeon Kim, Nandita Dukkipati, Anthony Baryshnikov, Christos Kaplanis, XiangHai Sheng, Yuri Chervonyi, Caglar Unlu, Diego de Las Casas, Harry Askham, Kathryn Tunyasuvunakool, Felix Gimeno, Siim Poder, Chester Kwak, Matt Miecnikowski, Vahab Mirrokni, Alek Dimitriev, Aaron Parisi, Dangyi Liu, Tomy Tsai, Toby Shevlane, Christina Kouridi, Drew Garmon, Adrian Goedeckemeyer, Adam R. Brown, Anitha Vijayakumar, Ali Elqursh, Sadegh Jazayeri, Jin Huang, Sara Mc Carthy, Jay Hoover, Lucy Kim, Sandeep Kumar, Wei Chen, Courtney Biles, Garrett Bingham, Evan Rosen, Lisa Wang, Qijun Tan, David Engel, Francesco Pongetti, Dario de Cesare, Dongseong Hwang, Lily Yu, Jennifer Pullman, Srini Narayanan, Kyle Levin, Siddharth Gopal, Megan Li, Asaf Aharoni, Trieu Trinh, Jessica Lo, Norman Casagrande, Roopali Vij, Loic Matthey, Bramandia Ramadhana, Austin Matthews, C. J. Carey, Matthew Johnson, Kremena Goranova, Rohin Shah, Shereen Ashraf, Kingshuk Dasgupta, Rasmus Larsen, Yicheng Wang, Manish Reddy Vuyyuru, Chong Jiang, Joana Ijazi, Kazuki Osawa, Celine Smith, Ramya Sree Boppana, Taylan Bilal, Yuma Koizumi, Ying Xu, Yasemin Altun, Nir Shabat, Ben Bariach, Alex Korchemniy, Kiam Choo, Olaf Ronneberger, Chimezie Iwuanyanwu, Shubin Zhao, David Soergel, Cho-Jui Hsieh, Irene Cai, Shariq Iqbal, Martin Sundermeyer, Zhe Chen, Elie Bursztein, Chaitanya Malaviya, Fadi Biadsy, Prakash Shroff, Inderjit Dhillon, Tejas Latkar, Chris Dyer, Hannah Forbes, Massimo Nicosia, Vitaly Nikolaev, Somer Greene, Marin Georgiev, Pidong Wang, Nina Martin, Hanie Sedghi, John Zhang, Praseem Banzal, Doug Fritz, Vikram Rao, Xuezhi Wang, Jiageng Zhang, Viorica Patraucean, Dayou Du, Igor

- 1134 Mordatch, Ivan Jurin, Lewis Liu, Ayush Dubey, Abhi Mohan, Janek Nowakowski, Vlad-Doru Ion,
 1135 Nan Wei, Reiko Tojo, Maria Abi Raad, Drew A. Hudson, Vaishakh Keshava, Shubham Agrawal,
 1136 Kevin Ramirez, Zhichun Wu, Hoang Nguyen, Ji Liu, Madhavi Sewak, Bryce Petrini, DongHyun
 1137 Choi, Ivan Philips, Ziyue Wang, Ioana Bica, Ankush Garg, Jarek Wilkiewicz, Priyanka Agrawal,
 1138 Xiaowei Li, Danhao Guo, Emily Xue, Naseer Shaik, Andrew Leach, Sadh MNM Khan, Julia
 1139 Wiesinger, Sammy Jerome, Abhishek Chakladar, Alek Wenjiao Wang, Tina Ornduff, Folake Abu,
 1140 Alireza Ghaffarkhah, Marcus Wainwright, Mario Cortes, Frederick Liu, Joshua Maynez, Andreas
 1141 Terzis, Pouya Samangouei, Riham Mansour, Tomasz Kepa, François-Xavier Aubet, Anton Al-
 1142 gymr, Dan Banica, Agoston Weisz, Andras Orban, Alexandre Senges, Ewa Andrejczuk, Mark
 1143 Geller, Niccolo Dal Santo, Valentin Anklin, Majd Al Merey, Martin Baeuml, Trevor Strohman,
 1144 Junwen Bai, Slav Petrov, Yonghui Wu, Demis Hassabis, Koray Kavukcuoglu, Jeffrey Dean, and
 1145 Oriol Vinyals. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of
 1146 context, August 2024. URL <http://arxiv.org/abs/2403.05530>. arXiv:2403.05530
 [cs].
- 1147
- 1148 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
 1149 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher,
 1150 Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy
 1151 Fu, Wenjin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,
 1152 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madijan Khabsa, Isabel
 1153 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,
 1154 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,
 1155 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,
 1156 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh
 1157 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen
 1158 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,
 1159 Sergey Edunov, and Thomas Scialom. Llama 2: Open Foundation and Fine-Tuned Chat Models,
 1160 July 2023. URL <http://arxiv.org/abs/2307.09288>. arXiv:2307.09288 [cs].
- 1161
- 1162 Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia
 1163 Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process- and
 1164 outcome-based feedback, November 2022. URL <http://arxiv.org/abs/2211.14275>.
 arXiv:2211.14275 [cs].
- 1165
- 1166 Bingbing Wang, Bin Liang, Jiachen Du, Min Yang, and Ruifeng Xu. SEMGraph: Incorporating
 1167 Sentiment Knowledge and Eye Movement into Graph Model for Sentiment Analysis. In Yoav
 1168 Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on
 1169 Empirical Methods in Natural Language Processing*, pp. 7521–7531, Abu Dhabi, United Arab
 1170 Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.
 emnlp-main.510. URL <https://aclanthology.org/2022.emnlp-main.510>.
- 1171
- 1172 Bingbing Wang, Bin Liang, Lanjun Zhou, and Ruifeng Xu. Gaze-infused BERT: Do human
 1173 gaze signals help pre-trained language models? *Neural Computing and Applications*, April
 1174 2024a. ISSN 0941-0643, 1433-3058. doi: 10.1007/s00521-024-09725-8. URL <https://link.springer.com/10.1007/s00521-024-09725-8>.
- 1175
- 1176 Tianlu Wang, Ilia Kulikov, Olga Golovneva, Ping Yu, Weizhe Yuan, Jane Dwivedi-Yu,
 1177 Richard Yuanzhe Pang, Maryam Fazel-Zarandi, Jason Weston, and Xian Li. Self-Taught Evalu-
 1178 ators, August 2024b. URL <http://arxiv.org/abs/2408.02666>. arXiv:2408.02666
 [cs].
- 1179
- 1180 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and
 1181 Hannaneh Hajishirzi. Self-Instruct: Aligning Language Models with Self-Generated Instructions,
 1182 May 2023a. URL <http://arxiv.org/abs/2212.10560>. arXiv:2212.10560 [cs].
- 1183
- 1184 Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang,
 1185 Xin Jiang, and Qun Liu. Aligning Large Language Models with Human: A Survey, July 2023b.
 1186 URL <http://arxiv.org/abs/2307.12966>. arXiv:2307.12966 [cs].
- 1187
- 1188 Zhilin Wang, Yi Dong, Jiaqi Zeng, Virginia Adams, Makesh Narsimhan Sreedhar, Daniel Egert,
 1189 Olivier Delalleau, Jane Polak Scowcroft, Neel Kant, Aidan Swope, and Oleksii Kuchaiev.

- 1188 HelpSteer: Multi-attribute Helpfulness Dataset for SteerLM, November 2023c. URL <http://arxiv.org/abs/2311.09528>. arXiv:2311.09528 [cs].
1189
1190
- 1191 Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy J. Zhang,
1192 Makesh Narsimhan Sreedhar, and Oleksii Kuchaiev. HelpSteer2: Open-source dataset for training
1193 top-performing reward models, June 2024c. URL <http://arxiv.org/abs/2406.08673>.
1194 arXiv:2406.08673 [cs].
1195
1196
- 1197 Guojun Wu, Lena Bolliger, David Reich, and Lena Jäger. An Eye Opener Regarding Task-Based
1198 Text Gradient Saliency. In Tatsuki Kurabayashi, Giulia Rambelli, Ece Takmaz, Philipp Wicke,
1199 and Yohei Oseki (eds.), *Proceedings of the Workshop on Cognitive Modeling and Computational
1200 Linguistics*, pp. 255–263, Bangkok, Thailand, August 2024a. Association for Computational Lin-
1201 guistics. URL <https://aclanthology.org/2024.cmcl-1.22>.
1202
1203
- 1204 Tianhao Wu, Banghua Zhu, Ruoyu Zhang, Zhaojin Wen, Kannan Ramchandran, and Jiantao Jiao.
1205 Pairwise Proximal Policy Optimization: Harnessing Relative Feedback for LLM Alignment, Oc-
1206 tober 2023a. URL <http://arxiv.org/abs/2310.00212>. arXiv:2310.00212 [cs].
1207
1208
- 1209 Tianhao Wu, Weizhe Yuan, Olga Golovneva, Jing Xu, Yuandong Tian, Jiantao Jiao, Jason We-
1210 ston, and Sainbayar Sukhbaatar. Meta-Rewarding Language Models: Self-Improving Alignment
1211 with LLM-as-a-Meta-Judge, July 2024b. URL <http://arxiv.org/abs/2407.19594>.
1212 arXiv:2407.19594 [cs].
1213
1214
- 1215 Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A.
1216 Smith, Mari Ostendorf, and Hannaneh Hajishirzi. Fine-Grained Human Feedback Gives Bet-
1217 ter Rewards for Language Model Training, October 2023b. URL <http://arxiv.org/abs/2306.01693>. arXiv:2306.01693 [cs].
1218
1219
- 1220 Duo Yang and Nora Hollenstein. PLM-AS: Pre-trained Language Models Augmented with Scan-
1221 paths for Sentiment Classification. *Proceedings of the Northern Lights Deep Learning Workshop*,
1222 4, January 2023. ISSN 2703-6928. doi: 10.7557/18.6797. URL <https://septentrio.uitt.no/index.php/nldl/article/view/6797>.
1223
1224
- 1225 Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing
1226 Yin, and Xia Hu. Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond,
1227 April 2023. URL <http://arxiv.org/abs/2304.13712>. arXiv:2304.13712 [cs].
1228
1229
- 1230 Kevin Yang, Dan Klein, Asli Celikyilmaz, Nanyun Peng, and Yuandong Tian. RLCD: Reinforce-
1231 ment Learning from Contrastive Distillation for Language Model Alignment, March 2024. URL
1232 <http://arxiv.org/abs/2307.12950>. arXiv:2307.12950 [cs].
1233
1234
- 1235 Jie Yu, Wenya Kong, and Fangfang Liu. CeER: A Nested Name Entity Recognition Model Incor-
1236 porating Gaze Feature. In Wenjie Zhang, Anthony Tung, Zhonglong Zheng, Zhengyi Yang, Xi-
1237 aoyang Wang, and Hongjie Guo (eds.), *Web and Big Data*, pp. 32–45, Singapore, 2024. Springer
1238 Nature Singapore. ISBN 978-981-9772-32-2.
1239
1240
- 1241 Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. RRHF:
1242 Rank Responses to Align Language Models with Human Feedback without tears, May 2023.
1243 URL <http://arxiv.org/abs/2304.05302>. arXiv:2304.05302 [cs].
1244
1245
- 1246 Leran Zhang and Nora Hollenstein. Eye-Tracking Features Masking Transformer Attention in
1247 Question-Answering Tasks. In Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste, Alessan-
1248 dro Lenci, Sakriani Sakti, and Nianwen Xue (eds.), *Proceedings of the 2024 Joint Interna-
1249 tional Conference on Computational Linguistics, Language Resources and Evaluation (LREC-
1250 COLING 2024)*, pp. 7057–7070, Torino, Italy, May 2024. ELRA and ICCL. URL <https://aclanthology.org/2024.lrec-main.619>.
1251
1252
- 1253 Yifan Zhang, Jiliang Li, Zachary Karas, Aakash Bansal, Toby Jia-Jun Li, Collin McMillan,
1254 Kevin Leach, and Yu Huang. EyeTrans: Merging Human and Machine Attention for Neu-
1255 ral Code Summarization, February 2024. URL <http://arxiv.org/abs/2402.14096>.
1256 arXiv:2402.14096 [cs].
1257
1258

1242 Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J. Liu.
 1243 SLiC-HF: Sequence Likelihood Calibration with Human Feedback, May 2023. URL <http://arxiv.org/abs/2305.10425> [cs].
 1244
 1245

1246 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
 1247 Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
 1248 Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena, December 2023. URL <http://arxiv.org/abs/2306.05685> [cs].
 1249
 1250

1251 A APPENDIX

1252 A.1 IMPLEMENTATION DETAILS

1253 This section provides further details on the implementation of our method. Subsection A.1.1 pro-
 1254 vides more details on the datasets used and their preprocessing steps. In subsection A.1.2, further
 1255 information is given about the models used for generating ET features, along with the specific pre-
 1256 processing required for each. Subsection A.1.3 explains the process of mapping the fixations from
 1257 the tokenizer used by the generation model to the tokenizer used by the Reward Model. Subsection
 1258 A.1.4 give more details on the checkpoints used for the RM backbone models. Finally, additional
 1259 implementation details are discussed in subsection A.1.5.

1260 A.1.1 DATASET PROCESSING

1261 In this subsection, we provide more details about the datasets used and the preprocessing to train
 1262 the RM. We use two different datasets: OpenAssistant Conversations dataset's (OASST1)³ (Köpf
 1263 et al., 2023) and HelpSteer2⁴ (Wang et al., 2024c).

1264 **OASST1.** A human-generated, human-annotated assistant-style conversation corpus consisting of
 1265 161,443 messages in 35 different languages, resulting in over 10,000 complete and fully annotated
 1266 conversation trees. The basic data structure is a Conversation Tree (CT), with nodes representing
 1267 written messages in a conversation. A CT's root node represents an initial prompt, given by the
 1268 prompter. The data was collected using a web-app interface as a product of a worldwide crowd-
 1269 sourcing effort involving over 13,500 volunteers, dividing the collection of each tree into five sepa-
 1270 rate steps: prompting, labelling prompts, adding reply messages as prompter or assistant, labelling
 1271 replies, and ranking assistant replies.

1272 **HelpSteer2.** A CC-BY-4.0-licensed open-source helpfulness dataset, designed to train state-of-the-
 1273 art RM consisting on 10,000 response pairs. It collects prompts mainly from ShareGPT⁵, focusing
 1274 on user inputs and filtering out non-English and programming-related prompts for quality. The
 1275 prompts are clustered into topics and sampled based on complexity to ensure diversity. Multi-turn
 1276 prompts are generated using an in-house model, with responses sourced from various internal mod-
 1277 els and human annotators. For each response, they annotate five attributes (helpfulness, correctness,
 1278 coherence, complexity, and verbosity) on a Likert-5 scale involving multiple annotators for each
 1279 response, ensuring high-quality ratings across five attributes.

1280 Conversation format and dataset preparation.

1281 To fine-tune LLMs for human-AI interaction, we need to define a chat protocol. We use a multi-
 1282 message chat setup with a special header and termination tokens, similar to the one in Llama 3
 1283 Dubey et al. (2024). The header tokens differentiate the turns between the user and the system. For
 1284 this, we use the *apply_chat_template*⁶ feature from *FastTokenizers* in the *transformers* library.

1285 The tokenizer used by the Meta-Llama-3-8B-Instruct model already incorporates this chat format
 1286 since this model has already undergone human alignment. Therefore, we use this format in our ex-
 1287 periments. For the other two models, we employ the default chat format provided by their respective
 1288 tokenizers. We add new tokens in the embeddings layer for these chat formats and we train them as
 1289 part of our process. Below, we provide an example of the template for each model.

1290³<https://huggingface.co/datasets/OpenAssistant/oasst1>

1291⁴<https://huggingface.co/datasets/nvidia/HelpSteer2>

1292⁵<https://huggingface.co/datasets/RyokoAI/ShareGPT52K>

1293⁶https://huggingface.co/docs/transformers/main/en/chat_templating

- **Meta-Llama-3-8B-Instruct:** <begin_of_text><start_header_id>user<end_header_id>
Example Prompt <eot_id><start_header_id>assistant<end_header_id> Example Response <eot_id>
- **Meta-Llama-3-8B:** <im_start>user Example Prompt <im_end> <im_start>assistant
Example Response <im_end>
- **Mistral-7B:** <s>[INST] Example Prompt [/INST] Example Response </s>

A.1.2 EYE-TRACKER FEATURES GENERATION MODELS

Special tokens are removed from the text before it is tokenized with the corresponding tokenizer used for the ET generator model. This is done to ensure that special tokens related to the chat format are not included in the input and are not assigned ET features to them, since these tokens are just for the RM to understand the chat format.

First model: Model presented in Huang & Hollenstein (2023). The code for the model along with the weights is publicly available, so we used the pre-trained checkpoint and we adapted their code for our implementation. This model was trained on several eye-tracking datasets, including Dundee (Kennedy et al., 2012), GECO (Cop et al., 2017), ZuCo1 (Hollenstein et al., 2018), ZuCo2 (Hollenstein et al., 2020a). More detailed information about this datasets is presented in Table 7. **The best model achieves an mean squared error (MSE) of 4.02 on a randomly held-out test set (25% of all data). This model has 17.5M pararemetes and remains frozen during training.** In Figure 3 we show an example of the synthetic total reading time (TRT) generated for the chosen and rejected response to a prompt.

For training this model, since the fixation duration is distributed differently across corpora, the authors normalize the fixation duration for each corpus, by dividing it by the mean duration of the corpus. Moreover, they map the duration values to discrete space $[1, 2, \dots, K]$. Using K-quantiles, the fixation values were partitioned into K subsets of nearly equal sizes, and each value was assigned to the index of the corresponding subset. The model is then trained in a multi-task setting, computing the mean and variance of the fixation duration. This quartile-based processing is used in other works (Huang et al., 2023) that use ET data to improve performance in NLP tasks, and we believe it is part of the reason why we obtained better results with this combination when training the RM. The authors proposed a method specifically for converting word-level TRT to token-level fixation data during model training. Initially, the TRT of a word is assigned to its characters, then a small number is assigned to the last character of the word (mainly to give small values to punctuation). After tokenizing the word the span of each subword is obtained, and the maximum value in each span is taken as the final token-level fixation data.

To use this model in our setup, we need to reverse this conversion process and recompute the features from token-level back to word-level, allowing us to remap the features to a different tokenizer. This is done by summing the original features for all tokens corresponding to a word and then distributing them across the tokens mapped to the same word in the other tokenizer. More details of this conversion are explained in Appendix A.1.3 and an example of the process in Table 8.

Second model: Model presented in Li & Rudzicz (2021). The code and training data are also publicly available, so we trained it following their original methodology and we adapted their code to incorporate it into our implementation. The model was trained using the ZuCo 1 (Hollenstein et al., 2018) and ZuCo 2 (Hollenstein et al., 2020a) and PROVO (Luke & Christianson, 2018) datasets. For the ZuCo datasets, 800 sentences (15.7 tokens) were provided as training data and 191 sentences (3.5k tokens) were held out for evaluation. Information about the datasets used to train these models is in Table 7. The model is based on RoBERTa (Liu et al., 2019) with a regression head on each token. This head is a linear layer that outputs five features: FFD, fixProp, GPT, TRT, and nFix (Table 1). **mean absolute error (MAE) for each feature is presented in Table 6. This model has 125M parameters and remains frozen during training.**

In this generative model, the conversion of word-level features to token-level features during training is done by assigning the features of a word to the first token and it is assumed that the rest of the tokens of this word do not have features. We reversed this process similarly during inference by forcing the predictions for tokens that are not the first in a word to be zero. Since the maximum sequence length for RoBERTa is 512 and we are dealing with longer sequences, we implemented a

1350 here's a limerick about cucumbers: there once there once was a cucumber so green, the biggest
 1351 was a green cucumber, that was really quite you've ever seen, it was fresh and so cool,
 1352 a late bloomer, it grew on the vine, and was but ended up in a pool, now it's a pickle,
 1353 oh so fine, until it became someone's consumer. not so serene.

1354 TRT per word in chosen response.

1355 TRT per word in rejected response.

1356 Figure 3: TRT generated by first model (Huang & Hollenstein, 2023) of the chosen and rejected
 1357 response to prompt 'Create a limerick about cucumbers'. Deeper colour represents longer fixation.
 1358

1360 sliding window approach. We split the input into sequences of 512 tokens with a 50-token overlap.
 1361 After processing, we combine the results using a linear combination.
 1362

1363 Table 6: MAE performance of the model reported in Li & Rudzicz (2021).

nFix	FFD	GPT	TRT	fixProp	All (Dev)
3.984	0.713	2.424	1.556	10.781	3.892

1369 Table 7: Overview of different corpora used to train the ET features generator models.

Corpus	Lang.	Sents.	Tokens	Subjects	Reference
Dundee	EN	2367	58598	20	Kennedy et al. (2012)
Provo	EN	189	2659	84	Luke & Christianson (2018)
ZuCo 1	EN	300	6588	12	Hollenstein et al. (2018)
ZuCo 2	EN	349	6828	18	Hollenstein et al. (2020b)
Geco	EN*	2449	-	23	Cop et al. (2017)

1379 A.1.3 MAPPING ET FEATURES BETWEEN DIFFERENT TOKENIZERS

1380 Both ET features generator models used are based on different tokenizers, which are also different
 1381 from the tokenizers employed by the based models used as RM. As a result, the number of tokens
 1382 n in the input for the reward model and the number of tokens w for the ET features may not be
 1383 the same. For **GazeAdd**, to be able to combine elementwise the ET feature embedding and the
 1384 text embedding, they must have the same temporal dimensions. Therefore, we need to map the ET
 1385 features per token from the ET tokenizer to the tokens of the RM tokenizer. Specifically, we convert
 1386 our $f_{et} \in \mathbb{R}^{w \times f}$ (w is the number of tokens, and f is the number of features) to $f_{et}^{mapped} \in \mathbb{R}^{n \times f}$
 1387 where n is the number of tokens in the RM input. For that, we perform a mapping between the two
 1388 tokenizers to obtain the mapped features f_{et}^{mapped} .
 1389

1390 To map tokens generated by two different tokenizers, we use our *EyeTrackPy* python library that
 1391 will be publicly released. First, we perform an initial mapping of tokens to the words they belong to
 1392 in each tokenizer with some properties of *FastTokenizers* from the *transformers* library⁷. Then, we
 1393 map words from one tokenizer to the words in the other and finally, we assume that the combination
 1394 of the tokens that are mapped to a word in one tokenizer correspond to the tokens that are mapped
 1395 to the word that is mapped to the initial word in the other tokenizer. Each row in Table 8, refers to a
 1396 step in this process.

1397 For each predictor, we reverse the method used to convert word-level features into token-level fea-
 1398 tures (more details in Appendix A.1.2) but passing from tokens in the first one, to tokens in the
 1399 second tokenizer. For example, if for the first ET features predictor models tokens t_1, t_2 are mapped
 1400 to tokens t_1, t_2, t_3 in another second tokenizer, the values sum for all the tokens in the first list and
 1401 distribute them equally across all the tokens in the second list: being t_1 (1s TRT) and t_2 (2s TRT)
 1402 each of t_1, t_2, t_3 are assigned a TRT of $(1 + 2)/3 = 1s$. In Table 8 is represented an example of this
 1403 process where row TRT(1) are the final TRT mapped for the first ET predictor and TRT(2) for the

1404 ⁷https://huggingface.co/docs/transformers/main_classes/tokenizer

second one. Finally, because special chat tokens were removed when generating the ET features, we assign value 0 for all features in this tokens. At the time of publishing this work, some of the tokenizer functionalities needed for alignment between tokenizers were not available in Mistral 7B.

Table 8: Example of mapping TRT between two different tokenizers. TRT (1) represents the process used for the first ET predictor, and TRT(2) for the second ET predictor.

	Tokenizer 1	Tokenizer 2
Words	astrophotography	astrophotography
Tokens str	['Astro', 'photo', 'graphy']	[‘Ast’, ‘raph’, ‘ot’, ‘ography’]
Tokens idx	[22, 23, 24]	[23, 24, 25, 26, 271]
Tokens IDs	[15001, 17720, 16369]	[198, 62152, 22761, 354, 5814]
TRT (1)	[11.23, 11.49, 10.16]	[6.58, 6.58, 6.58, 6.58, 6.58]
TRT (2)	[24.53, 0, 0]	[24.53, 0, 0, 0, 0]

A.1.4 MODELS

As RM base models we use the pretrained checkpoint of Hugging Face for Llama 3 8B⁸ (Dubey et al., 2024), Llama 3 8B-instruct⁹(Dubey et al., 2024) and Mistral 7B¹⁰ (Jiang et al., 2023a).

A.1.5 TRAINING DATAILS

During training, the ET features predictor model remains frozen, while the ET features projector is normally trained without adapters (Figure 2). The RM is fine-tuned on top of the open-source models using QLoRA (Dettmers et al., 2023) based on Low-Rank Adaptation (LoRA) (Hu et al., 2021), which fine-tunes select dense layers by optimizing low-rank decomposition matrices representing weight changes, instead of directly updating pre-trained weights. QLoRA introduces memory optimization techniques such as the 4-bit NormalFloat (NF4), a novel data type, to improve performance without increasing memory usage. Following Dettmers et al. (2023) we use hyperparameters: r=8, alpha=32, and dropout=0.1.

We also fine-tune the RM embedding layer, since we are adding new tokens for the chat format and special separators tokens in our **RewardConcat** method (section 4). Also, the last layer added to the RM for the scalar reward is trained from scratch without adapters. Our implementation is based in *pytorch* and we use *transformers*¹¹ from Hugging Face.

Hardware. We trained the models on servers equipped with 2x Intel Xeon Platinum 8470 CPUs, 1TB of RAM, and either 2x NVIDIA H100 (80GB) or 4x NVIDIA A100 (80GB) GPUs. We always train using only GPU at a time per each model and training times were between 20 and 50 hours depending mainly on the number of steps between model evaluations.

B REWARD BENCHMARK

As we described in section 6, a future direction would be to train a Reward Model on a larger model with more data. It has been proven the scaling trends for the reward model; More data and a larger-size model generally improve accuracy (Touvron et al., 2023). Nevertheless, as a complement to our results, we evaluated our trained models on the dataset with the best results, OAAST1, in this RewardBench, a benchmark for Reward Models. RewardBench, proposed in Lambert et al. (2024), is a benchmark designed to evaluate the performance and safety of reward models. It consists of a set of datasets intended for measuring how reward models perform on challenging prompts across chat, reasoning, and safety domains, using a trio structure of prompt-chosen-rejected pairs. It comprises 2985 diverse tasks, each sample is formatted as a prompt with a manual or human-verified chosen and rejected completion. Due to its diversity of tasks (4 categories and 23 sub-categories) this

⁸<https://huggingface.co/meta-llama/Meta-Llama-3-8B>

⁹[meta-llama/Meta-Llama-3-8B-Instruct](https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct)

¹⁰<https://huggingface.co/mistralai/Mistral-7B-v0.3>

¹¹<https://huggingface.co/docs/transformers/index>

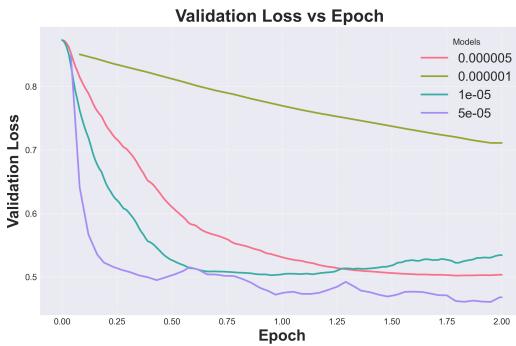


Figure 4: Validation loss with different LR on ConcatReward, batch size 8, features: f_1 and Meta-Llama-3-8B-Instruct base model

benchmark minimizes the likelihood of overfitting. Task accuracy is calculated based on whether the chosen response receives a higher reward than the rejected response. We directly evaluate their open dataset reward-bench¹²

C HYPER PARAMETER TUNING

We performed hyperparameter tuning for the **GazeConcat** method and the baseline, and we replicated them in the **GazeAdd** method, as testing all combinations is computationally very expensive. For each dataset, 15% is reserved for validation to perform hyperparameter tuning. The best-performing checkpoints are selected based on the lowest validation loss and are subsequently used for performance evaluation in the test split. We trained for two epochs, as described in Wang et al. (2024c) and in line with trends observed in Touvron et al. (2023) where they found that training longer can lead to over-fitting. We perform a grid search for the optimal batch size, testing {8, 16, 32} values. AdamW optimizer (Loshchilov & Hutter, 2019) is used and the learning rate is tuned within the range of {1, 5, 10, 50}e-6, inspired by the values reported in Touvron et al. (2023); Cui et al. (2024); Wang et al. (2024c). Additionally, we explored different scheduler configuration, comparing constant, linear, and cosine with a minimum learning rate.

In general, the parameter that most affected the validation results was the learning rate. For the others, we ended up choosing values that worked well across all combinations. We achieved better results in both the baseline and the models concatenating the ET features with a learning rate of 0.0005 (Figure 4). In Figure 6a and Figure 6b, the validation loss and learning rate with different schedulers are shown. For lower learning rates, such as 0.00001, the scheduler had little effect. However, with higher learning rates, using a scheduler helped to mitigate overfitting. In Figure 7a and Figure 7b, it represents validation loss and learning rates with a higher learning rate of 0.0005. We opted to use this 0.00005 learning rate for all experiments, employing a cosine learning rate scheduler with a minimum learning rate of 0.7, in line with other studies such as Wang et al. (2024c); Touvron et al. (2023). We did not find a significant effect of training batch size on validation accuracy, but we opted for a value of 8, which often (especially with high learning rates and without a scheduler) was the one that tended to overfit the least (Figure 5).

ET features projector The PyTorch architecture of our ET features projector is shown below. $num_features$ varies between 1, 5, and 2 (depending on the configuration used f_{comb1} , $f_{comb2.5}$, $f_{comb2.2}$ subsection 5.2). p_1, p_2 are dropout values. Finally, after testing different combinations, we used 0.1 and 0.3. **This model has 0.53M parameters.**

¹²<https://huggingface.co/datasets/allenai/reward-bench>

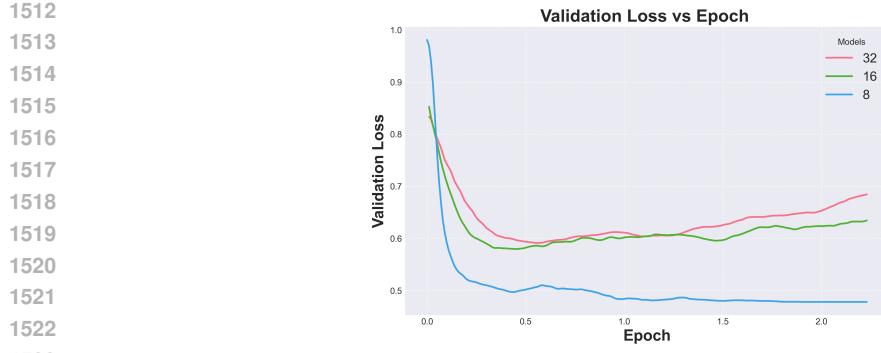


Figure 5: Validation loss with different batch size, learning rate: 5e-5, features: f_1 and Meta-Llama-3-8B-Instruct base model

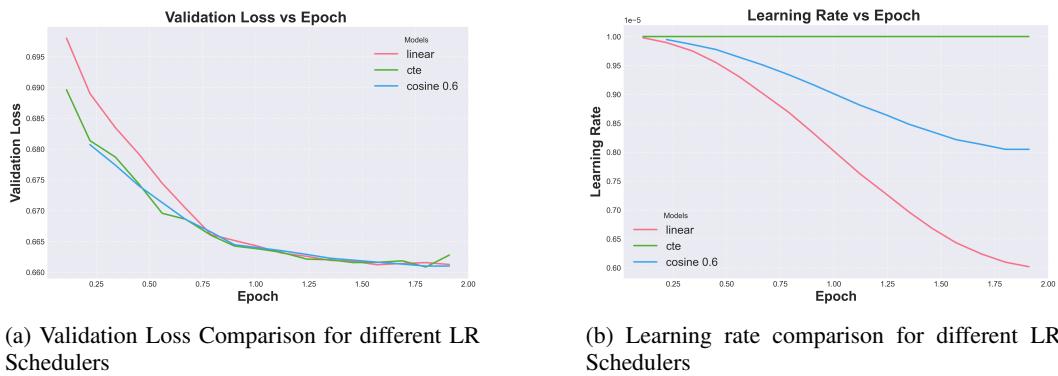


Figure 6: ConcatReward, LR: 0.00001, features: f_1 and Meta-Llama-3-8B-Instruct base model

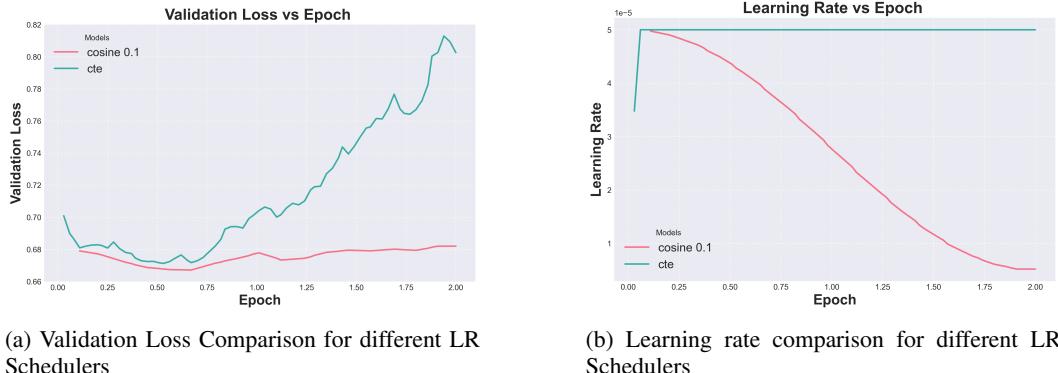


Figure 7: ConcatReward, LR: 0.00005, features: f_1 and Meta-Llama-3-8B-Instruct base model

Listing 1: PyTorch architecture of our gaze features projector

```
self.fixations_embedding_projector = nn.Sequential(
    nn.Linear(num_features, 128),
    nn.LayerNorm(128),
    nn.ReLU(),
    nn.Dropout(p=p_1),
    nn.Linear(128, hidden_size),
    nn.Dropout(p=p_2),
)
```