

Seeing Eye to AI: Human Alignment via Gaze-Based Response Rewards for Large Language Models

Angela Lopez-Cardona ^{1,2} Carlos Segura ¹ Alexandros Karatzoglou ³ Sergi Abadal ² Ioannis Arapakis ¹

¹Telefónica Scientific Research, Barcelona, Spain ²Universitat Politècnica de Catalunya, Barcelona, Spain ³Amazon, Barcelona, Spain

4/2025



Table of Contents

- 1 Introduction

- 2 Background
 - Eye-tracking (ET) technology
 - Large Language Models (LLMs) and Human Alignment

- 3 Related work

- 4 Methodology

- 5 Results

- 6 Conclusions

1 Introduction

2 Background

- Eye-tracking (ET) technology
- Large Language Models (LLMs) and Human Alignment

3 Related work

4 Methodology

5 Results

6 Conclusions

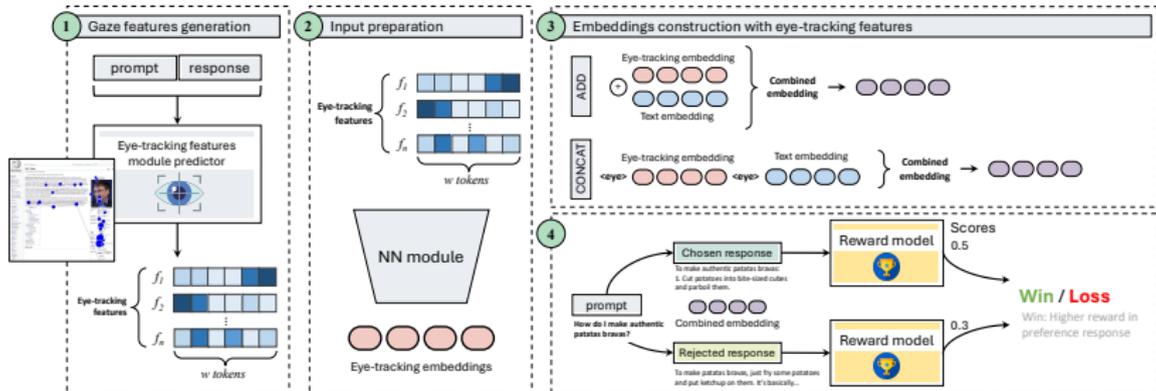


Figure: GazeReward Framework for using eye-tracking data for reward modelling. We use a generator model to compute ET features on a preference dataset D and we train the human preference by combining both text and ET embeddings

1 Introduction

2 Background

- Eye-tracking (ET) technology
- Large Language Models (LLMs) and Human Alignment

3 Related work

4 Methodology

5 Results

6 Conclusions

Eye-tracking (ET) technology

- Eye-tracking devices gather **raw data** about scanpaths, frequency and duration of fixations, as well as pupillary dilation variations.
- Algorithms to detect **fixations sequence**.
- From this **fixations sequence** → reading measures.



Figure: Eye-tracker

here's a limerick about cucumbers: there once
 was a green cucumber, that was really quite
 a late bloomer, it grew on the vine, and was
 oh so fine, until it became someone's consumer.

Figure: TRT per word. Deeper colour represents longer fixation.

Table: ET reading measures per word [10]

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

Training Large Language Models (LLMs)

- Unsupervised pretraining from raw text → **Language understanding**
- Large scale instruction tuning and reinforcement learning → **Human alignment**

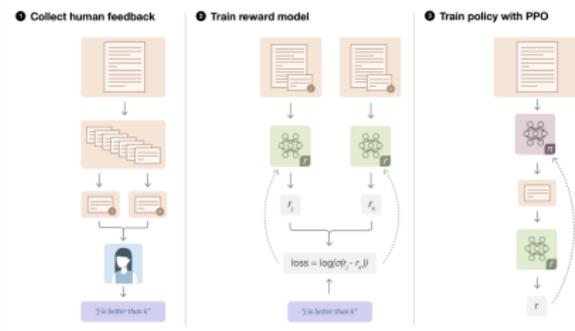


Figure: RLHF diagram [1]

RLHF phases:

- (1) collecting feedback
- (2) training a RM based on that feedback
- (3) optimising the LLMs using RL techniques

Other methods:

- Direct Preference Optimization (DPO) [25]
- Reinforcement Learning from AI Feedback (RLAIF) [2]

Reward Model (RM)

- The most **popular approach** to reward modeling follows the framework introduced by [23].
- Pretrained LLM + regression head that outputs a scalar **reward** [27].
- Reward indicates the quality of the model generation \Rightarrow proxy for human judgment.
- **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 [23].
- Loss function is defined in Equation 1, where y_w refers to the **preferred response** in a **pair** of completions y_w and y_l .

$$\text{loss}(\theta) = -E_{(x, y_w, y_l) \sim D} [\log(\sigma(r_\theta(x, y_w) - r_\theta(x, y_l)))] \quad (1)$$

Beyond its role in training, RMs are also used in **inference**, as seen in best-of-N sampling [4]. Additionally, RMs play a key role in generating **synthetic data** for preference alignment [2].

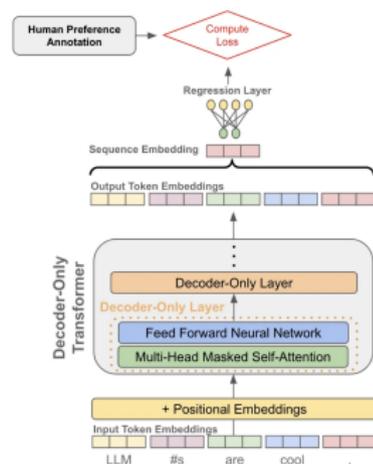


Figure: Reward model [1]

1 Introduction

2 Background

- Eye-tracking (ET) technology
- Large Language Models (LLMs) and Human Alignment

3 Related work

4 Methodology

5 Results

6 Conclusions

Main research questions → how to combine this two elements?

State of the Art

- Several studies have examined **alternative** versions for refining **RMs**. More **fine-grained** reward structures [2, 32, 6, 30]. Another line of research has focused on Process Based Reward Models (PRMs) [19, 28]. Some works have leveraged **synthetic** preference data for reward modelling [4, 14, 31] and self-training [24].
- Optimal methods for gathering **feedback** to align LLMs with human goals remain an **open question** [3].
- **Explicit feedback** is collected after users have reviewed model outputs. However, human decision making is inherently complex and involves a comprehensive evaluation of diverse information types before taking action.
- **ET** data has been shown to add value in various **NLP tasks**, as demonstrated by prior work [12, 15, 9, 33, 16, 5, 21, 22].

Contribution

- We augment the traditional **Reward Model**, by incorporating (artificial) **implicit feedback** through ET.

1 Introduction

2 Background

- Eye-tracking (ET) technology
- Large Language Models (LLMs) and Human Alignment

3 Related work

4 Methodology

5 Results

6 Conclusions

Method

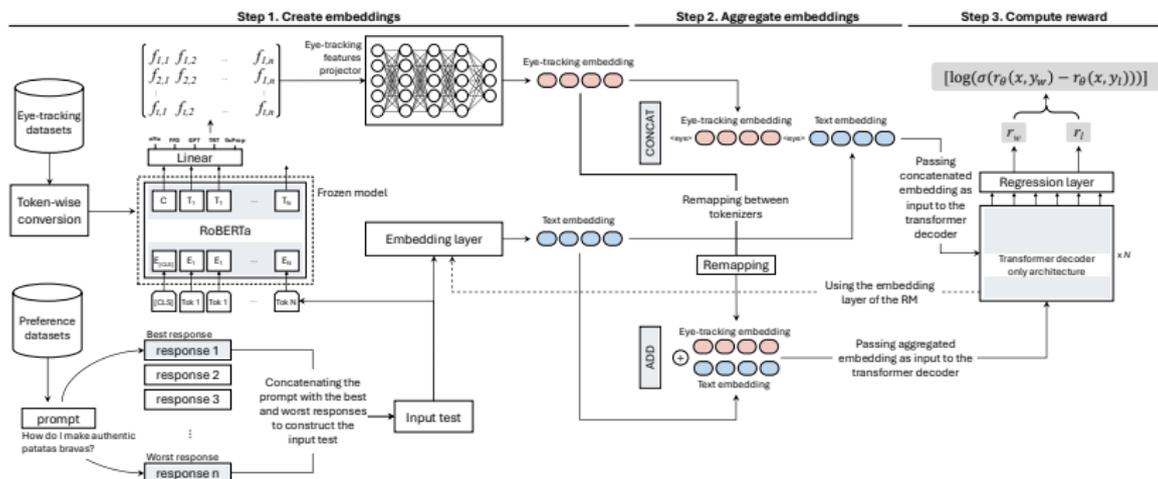


Figure: Overview of the **GazeReward** framework, incorporating eye-tracking features into the reward model. The architecture is illustrated in the figure using the second ET prediction model, but it would be identical if the first one were used instead

Experimental setup

2 Inclusion method:

- **GazeConcat:** The ET embedding, denoted as emb_{ETF} , is concatenated with the text embedding H to form the input for the RM: $(emb(\langle eye \rangle) \circ emb_{ETF} \circ emb(\langle / eye \rangle) \circ H)$.
- **GazeAdd:** The input to the RM consists of the ET embedding emb_{ETF} and the text embedding H , which are added in an elementwise fashion: $(emb_{ETF} + H)$.

2 Datasets:

[Table:](#) 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*	[17]
HelpSteer2	5938	1049	364	EN	[29]

3 open source backbone models:

Llama 3 8B, Llama 3 8B-instruct [7] and Mistral 7B [13].

3 ET feature importance:

- $fcomb_1$ – TRT generated by [11].
- $fcomb_{2.5}$ – nFix, FFD, GPT, TRT, fixProp generated by [18].
- $fcomb_{2.2}$ – TRT and FFD generated by [18].

Experimental setup

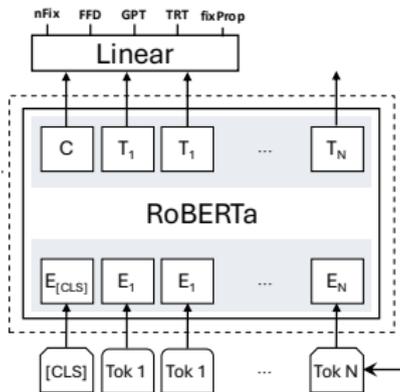


Figure: Generative model architecture [18].

2 Generative models:

- [18]: Based on RoBERTa [20] with a regression head on each token. This head is a linear layer that outputs five features: FFD, fixProp, GPT, TRT, and nFix.
- [11]: Based on T5 embedding layer [26], a two-layer BiLSTM [8], and a one-hidden-layer MLP. Predicts total reading time (TRT) per token.

1 Introduction

2 Background

- Eye-tracking (ET) technology
- Large Language Models (LLMs) and Human Alignment

3 Related work

4 Methodology

5 Results

6 Conclusions

Results

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

		Llama-3-8B-Instruct		Llama-3-8B		Mistral-7B	
baseline		65.9 ± 0.5	diff (%)	65.5 ± 2.1	diff (%)	66.3 ± 0.1	diff (%)
GazeConcat	<i>fcomb₁</i>	69.0 ± 0.4*	4.7	69.3 ± 0.6	5.9	67.6 ± 1.7	2.1
	<i>fcomb_{2.5}</i>	<u>70.2 ± 0.3**</u>	6.6	71.5 ± 0.5	9.2	<u>70.2 ± 0.4*</u>	5.9
	<i>fcomb_{2.2}</i>	70.0 ± 0.4**	6.3	<u>71.2 ± 0.8</u>	8.8	71.0 ± 1.0	7.1
GazeAdd	<i>fcomb₁</i>	68.9 ± 0.9	4.6	68.9 ± 1.0	5.3	-	-
	<i>fcomb_{2.5}</i>	70.2 ± 0.1*	6.6	69.5 ± 0.3	6.1	-	-
	<i>fcomb_{2.2}</i>	69.0 ± 0.4*	4.7	68.3 ± 0.7	4.4	-	-

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

		Llama-3-8B-Instruct		Llama-3-8B		Mistral-7B	
baseline		54.7 ± 0.7	diff (%)	53.3 ± 0.8	diff (%)	54.1 ± 0.3	diff (%)
GazeConcat	<i>fcomb₁</i>	<u>61.1 ± 1.2*</u>	11.8	59.1 ± 0.2*	10.8	57.6 ± 2.3	6.5
	<i>fcomb_{2.5}</i>	58.5 ± 1.6	7.0	<u>60.3 ± 0.5**</u>	13.2	58.7 ± 2.4	8.5
	<i>fcomb_{2.2}</i>	60.6 ± 3.3	10.9	<u>57.9 ± 2.0</u>	8.6	56.0 ± 2.4	3.4
GazeAdd	<i>fcomb₁</i>	62.3 ± 0.6**	13.9	62.4 ± 1.0**	17.0	-	-
	<i>fcomb_{2.5}</i>	59.6 ± 1.1*	9.0	58.6 ± 1.2*	10.0	-	-
	<i>fcomb_{2.2}</i>	60.3 ± 0.5**	10.2	59.3 ± 0.1*	11.3	-	-

1 Introduction

2 Background

- Eye-tracking (ET) technology
- Large Language Models (LLMs) and Human Alignment

3 Related work

4 Methodology

5 Results

6 Conclusions

Conclusions and future direction

- **Implicit Feedback for Reward Models:** The study presents a new framework integrating implicit feedback into the RM to enhance LLMs alignment and synthetic data generation.
- **Validation with Open-Source Models:** The approach was tested using Llama 3 and Mistral models, demonstrating consistent improvements in predicting user preferences in state-of-the-art models.
- **ET Features for Scalability:** The method leverages ET features generated by models, making it fully scalable and applicable to various human alignment techniques, including those using synthetic datasets.
- **Limitations in Data and Generalization:** The ET prediction models were trained on small, English-only datasets, raising questions about generalizability. Future work should focus on collecting task-specific ET data and exploring alternative feature integration methods.
- **Future Directions in Scaling and Application:** Further validation on larger models and datasets and integration into RLHF or rejection sampling for preference dataset generation could enhance alignment and overall LLM performance.

Contact information

Ángela López Cardona

angela.lopezcardona@telefonica.com



Figure: GitHub Respository

Acknowledgments



SYMBIOTIK Horizon Europe's European Innovation Council, under
Pathfinder Grant 101071147. Grant **AGAUR** Generalitat de
Catalunya, under Grant AGAUR 2023 DI060

