

INDOPREF: A Multi-Domain Pairwise Preference Dataset for Indonesian

Vanessa Rebecca Wiyono^{1*}, David Anugraha^{2*}, Ayu Purwarianti¹, Genta Indra Winata³

¹Institut Teknologi Bandung ²Stanford University ³Capital One
vanessarebecca29@gmail.com, david.anugraha@stanford.edu

Abstract

Over 200 million people speak Indonesian, yet the language remains significantly underrepresented in preference-based research for large language models (LLMs). Most existing multilingual datasets are derived from English translations, often resulting in content that lacks cultural and linguistic authenticity. To address this gap, we introduce INDOPREF, the first fully human-authored and multi-domain Indonesian preference dataset designed to evaluate the naturalness and quality of LLM-generated text. The dataset contains 522 prompts and yields 4,099 human-annotated pairwise preferences from comparisons across five instruction-tuned LLMs. All annotations are natively written in Indonesian with strong inter-annotator agreement, measured by Krippendorff’s alpha. Our benchmark spans 10 diverse categories, enabling practitioners to identify LLMs’ fine-grained strengths and weaknesses.¹

1 Introduction

Despite being spoken by over 200 million people and ranking among the world’s ten most widely spoken languages, Indonesian remains significantly underrepresented in NLP research (Koto et al., 2020a; Aji et al., 2022; Winata et al., 2023). Predictive analyses of multilingual models show that under-represented languages often suffer systematically lower performance, which motivates efforts for such languages (Anugraha et al., 2025c). While benchmarks on Indonesian, such as IndoLEM (Koto et al., 2020b), IndoNLU (Wilie et al., 2020), and IndoNLG (Cahyawijaya et al., 2021) have advanced Indonesian NLP in classification and generation tasks, they do not address the critical area of preference modeling. This gap stems largely from the lack of annotated

*The authors contributed equally.

¹Our dataset is released at <https://huggingface.co/datasets/davidanugraha/IndoPref> under CC-BY-4.0.

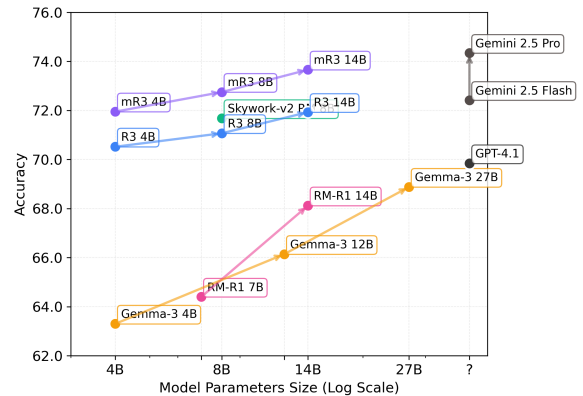


Figure 1: Model performance vs. model size on INDOPREF. The plot illustrates scaling trends across various model architectures, showing that larger models generally align better with human preferences.

preference datasets, limited language resources, and the absence of standardized evaluation benchmarks (Cahyawijaya et al., 2023; Lovenia et al., 2024), all of which hinder the development of models capable of capturing the linguistic and cultural nuances of Indonesian (Adilazuarda et al., 2024).

Preference datasets are essential for aligning model outputs with human expectations. Yet, no existing dataset offers native, human-authored preference annotations for Indonesian, leaving language models poorly equipped to reflect Indonesian-specific preferences. Prior datasets such as IndicXNLI (Aggarwal et al., 2022), M-RewardBench (Gureja et al., 2024), and Okapi (Lai et al., 2023) rely on translated content, which often introduces cultural mismatches and translation artifacts known to degrade model performance (Bizzone et al., 2020; Vanmassenhove et al., 2021). For instance, although M-RewardBench includes human-curated annotations, its prompts originate from translated English content, which may lack the naturalness and contextual relevance found in native Indonesian expressions.

To address this gap, we introduce **INDOPREF**, a high-quality, fully human-annotated dataset for training and evaluating preference-aligned Indonesian language models. The dataset comprises 522 multi-domain instruction–response prompts, natively authored by fluent Indonesian speakers, yielding 4,099 pairwise human preference annotations derived from comparisons across five instruction-tuned LLMs. Our annotated data shows strong inter-annotator agreement on both relevance and fluency, based on Krippendorff’s α . We further benchmark 15 models spanning diverse architectures and parameter scales to demonstrate the utility of **INDOPREF**. This work provides a robust and culturally grounded resource to advance Indonesian LLM development and promote more equitable progress in multilingual NLP.

2 IndoPref

The **INDOPREF** dataset provides high-quality human preference data in Indonesian, specifically designed to support the preference tuning of LLMs. **INDOPREF** is composed entirely of prompts and annotations natively authored in Indonesian by fluent speakers, ensuring that the data better reflects the linguistic intuition, cultural context, and pragmatic norms of Indonesian users. The dataset encompasses various domains, including safety, logic, summarization, translation, and creative writing, aiming to reflect diverse real-world use cases and support robust model alignment across different task types. By focusing on native language elicitation and annotation, **INDOPREF** fills a critical gap in preference data for underrepresented languages.

2.1 Data Collection

The **INDOPREF** dataset provides high-quality pairwise human preference annotations for fine-tuning large language models in Indonesian. It comprises 522 prompts authored by fluent native speakers, designed to reflect natural, contextually appropriate, and culturally grounded language use. To generate candidate responses, each prompt is submitted to an instruction-tuned LLM. The resulting responses are anonymized, randomly shuffled, and indexed to minimize annotator bias. These pairwise preference judgments form the foundation for downstream evaluation and alignment tasks.

2.2 Topics

Prompts in **INDOPREF** are carefully curated categories designed to reflect real-world instruction-

following tasks and diverse linguistic phenomena. They include both objective tasks with single correct answers and subjective ones with multiple valid responses. Categories range from deterministic tasks (math, logic, programming) to generative ones (creative writing, brainstorming, open-ended questions). Others, like translation, summarization, and analysis, test comprehension and synthesis, while safety prompts evaluate ethical sensitivity and harm mitigation. To construct the prompt set, we utilize multiple sources. Structured tasks in domains like math and coding are adapted from educational resources and online platforms. Tasks involving summarization and analysis derive from real-world texts such as articles and essays. Additionally, many prompts are written from scratch by native-speaking prompt designers to ensure originality and cultural relevance.

2.3 Annotations

To construct reliable preference data for fine-tuning, we implement a structured annotation workflow involving two independent groups of native Indonesian annotators. Each group evaluates the same set of prompts, each accompanied by five model-generated responses. Annotators assess each response based on two criteria: relevance and fluency. Relevance measures how well the response addresses the intent of the prompt, while fluency reflects the grammatical correctness, coherence, and naturalness of the language. Both criteria are rated on a 5-point Likert scale, allowing detailed differentiation among outputs. In addition to these ratings, annotators select one response per prompt that they considered the most preferable overall. Summarized statistics about the relevance and fluency scores for each LLM’s responses rated by human annotators can be found in Table 12.

The annotation process is designed to ensure consistency and minimize cognitive load, with responses presented in randomized order and anonymized. Annotators follow clear guidelines with illustrative rating examples to maintain uniform judgments. Inter-annotator agreement, measured using Krippendorff’s α , shows high reliability with 0.965 for relevance and 0.862 for fluency, demonstrating strong consistency across annotator groups. The overall Krippendorff’s α for human pairwise rankings across annotators is 0.891, further confirming the robustness of the dataset. The final step involves converting the annotations into a pairwise preference format suitable for evaluating

Model	Analysis	Brainstorm.	Coding	Creative Writing	Logic	Math	Open Question	Safety	Summ.	Translation	Avg.
R3 4B	80.43	68.79	62.34	83.07	72.51	77.86	81.92	64.62	73.08	40.58	70.52
R3 8B	81.52	68.79	61.10	79.63	72.28	78.33	82.61	64.62	72.84	48.99	71.07
R3 14B	79.13	70.97	61.60	80.16	70.29	82.62	82.61	64.62	76.92	50.43	71.93
mR3 4B	81.30	79.53	76.56	79.37	68.29	65.24	86.27	67.69	75.00	40.29	71.95
mR3 8B	83.04	76.17	73.07	84.66	67.85	66.67	84.67	72.31	73.32	45.80	72.75
mR3 14B	83.26	79.70	76.31	80.95	66.30	72.38	82.84	73.33	75.48	46.09	73.66
RM-R1 7B	76.96	69.46	62.59	78.04	59.87	45.95	83.52	59.49	62.02	46.09	64.40
RM-R1 14B	81.96	71.31	63.34	80.42	67.85	54.29	82.61	64.62	70.91	42.90	68.12
Gemma-3 4B	78.48	65.94	53.87	76.46	52.33	55.95	75.74	58.46	66.83	48.99	63.30
Gemma-3 12B	71.96	64.60	61.60	75.40	60.98	71.67	79.86	61.03	64.18	50.14	66.14
Gemma-3 27B	83.91	67.62	63.84	80.69	65.41	71.67	79.18	60.00	67.79	48.70	68.88
Skywork-v2 RM 8B [†]	84.13	73.99	67.83	81.48	54.54	69.52	84.67	71.28	68.51	60.87	71.68
GPT-4.1	80.87	63.76	63.09	76.72	66.96	77.86	81.92	62.05	70.43	54.78	69.84
Gemini 2.5 Pro	83.04	70.30	69.33	81.75	74.50	80.95	82.84	71.28	71.39	57.97	74.34
Gemini 2.5 Flash	82.17	72.32	68.08	81.75	72.73	74.05	81.92	65.64	70.43	55.07	72.42

Table 1: Fine-grained accuracy (%) results on INDOPREF across various open-source and proprietary models.

Model	English	Indonesian
GPT-4.1	69.84	63.45
Gemini 2.5 Pro	74.34	73.90

Table 2: Average accuracy of model performance on English and Indonesian prompts.

LLMs’ generations as LLM-as-a-judge. Through careful annotation, rigorous validation, and structured formatting, the INDOPREF dataset provides a strong foundation for training and evaluating preference-aligned language models in Indonesian.

3 Experimental Setup

Models. We evaluate nine open-weight LLMs: R3 (4B, 8B, 14B) (Anugraha et al., 2025b), mR3 with English prompt and English reasoning (4B, 8B, 14B) (Anugraha et al., 2025a), RM-R1 (7B, 14B) (Chen et al., 2025), Gemma-3 (4B, 12B, 27B), and Skywork-v2 RM 8B (Liu et al., 2025),² alongside three proprietary models: GPT-4.1 (Achiam et al., 2023), Gemini 2.5 Pro, and Gemini 2.5 Flash (Comanici et al., 2025). We perform inference using their recommended generation settings across all models.

Dataset and Evaluation. We utilize the full set of 4,099 pairwise preference annotations spanning 10 diverse categories, featuring chosen and rejected responses generated by five instruction-tuned LLMs: Llama 3.1 8B Instruct (Dubey et al., 2024), GPT-4o, GPT-4o mini (Achiam et al., 2023), Gemini 1.5 Flash (Team et al., 2024), and Aya Expansive 8B (Dang et al., 2024b). To evaluate model alignment with human preferences, we cal-

²<https://huggingface.co/Skywork/Skywork-Reward-V2-Llama-3.1-8B-40M>.

culate the accuracy of each model in predicting the human-preferred response.

4 Results and Analysis

As shown in Table 1, Gemini 2.5 Pro achieves the highest average accuracy across all evaluated models, with a score of 74.34, followed closely by mR3 14B, demonstrating strong overall performance. Among open-sourced models, mR3 14B stands out as the best-performing reward models on INDOPREF. Furthermore, mR3’s smallest model, mR3 4B, performs better than other model families, including Gemini 2.5 Flash. This indicates that smaller models with strong reasoning capabilities can serve as effective LLM-as-a-judge systems, highlighting the importance of architectural specialization over sheer scale.

Figure 1 also shows model performance across varying sizes. Within the same architecture family, such as R3, mR3, RM-R1, and the Gemma series, we observe consistent improvements as model size increases. This trend suggests that larger models are better equipped to evaluate responses accurately and align more closely with human preferences.

Fine-Grained Performance. Among all categories, Creative Writing, Open Question, and Analysis consistently yield the highest scores, with several models scoring above 80. This suggests that models are particularly well-tuned for open-ended text generation. In contrast, Translation appears to be the most challenging task. Most models score considerably lower in this category, pointing to limitations in cross-lingual transfer or a lack of high-quality multilingual training data.

On the other hand, tasks such as Math, Logic, and Coding exhibit noticeable performance variation across models. While systems like Gemini 2.5

Pro, R3, and mR3 perform reasonably well, others, such as including RM-R1 and Gemma, consistently lag behind, indicating persistent challenges in complex reasoning. Safety is another area where many models struggle. Overall, the results in Table 1 indicate that although recent models achieve strong general performance, reasoning-oriented tasks and specialized domains such as Safety and Translation remain important targets for further improvement.

Prompts in Target Language. To evaluate the effect of instruction language, we re-run the evaluation using the same prompt structure as in Table 2, but with all instructions, task descriptions, and output formats translated into Indonesian. As shown in the results, Gemini 2.5 Pro maintains strong performance with only a slight drop in average score (from 74.34 to 73.90), suggesting robust generalization to Indonesian-language instructions. In contrast, GPT-4.1 exhibits a more pronounced decrease (from 69.84 to 63.45), suggesting that its ability to follow instructions is more sensitive to the language used in the prompt formatting. We hypothesize that this difference stems from the distribution of instruction-tuning data each model sees during training. The Gemini 2.5 Pro is likely exposed to a broader range of multilingual instructions, allowing it better to understand task setups in languages other than English.

5 Related Work

Multilingual Preference Datasets. Several multilingual datasets have been developed to support preference alignment in non-English languages. The Okapi dataset covers 26 languages using translated prompts and includes human preference annotations (Lai et al., 2023). M-RewardBench introduces a multilingual reward model benchmark across 23 languages, including Indonesian (Gureja et al., 2024). However, the preference annotations rely on semi-automatic methods, which may introduce bias, and the prompts rely primarily on automatic translation (via Google Translate), with only post hoc human filtering of poor-quality outputs. Furthermore, the prompts and/or responses may not be culturally relevant to Indonesians.

Indonesian NLP Benchmarks. While preference tuning for Indonesian remains underexplored, numerous datasets have improved downstream NLP performance. IndoLEM (Koto et al., 2020b) and IndoNLU (Wilie et al., 2020) introduced la-

beled datasets for tasks such as POS tagging, parsing, and entailment. IndoNLG (Cahyawijaya et al., 2021) provided benchmarks for summarization, QA, and translation. NusaCrowd (Cahyawijaya et al., 2023) further compiled a diverse Indonesian and local language datasets for multi-task and instruction tuning. However, none include human preference annotations.

Preference Tuning in Non-English Languages.

Recent work has extended preference tuning to languages beyond English. Dang et al. (2024a) demonstrates multilingual preference tuning across 23 languages can provide strong cross-lingual transfer. Preference tuning has also been applied to improve translation quality into different languages using implicit preferences from human-authored text (Xu et al., 2024). Other efforts include aligning a Chinese bilingual LLM through human feedback (Hou et al., 2024), and generating persona-consistent dialogue in Japanese using pseudo preference tuning (Takayama et al., 2025). On the reward modeling side, recent work such as mR3 explores multilingual reward reasoning, proposing a rubric-agnostic model trained across 72 languages (Anugraha et al., 2025a). These studies highlight the growing interest in scaling preference-based alignment techniques to multilingual and low-resource settings.

6 Conclusion

This paper introduces INDOPREF, a new benchmark for evaluating LLMs on Indonesian-language preference tasks using fully human-authored data. By focusing specifically on Indonesian, the dataset fills a critical gap left by existing multilingual resources, which often rely on translated content that lacks linguistic and cultural fidelity. Evaluation results show that models like Gemini 2.5, R3, and mR3 models perform well overall, particularly in open-ended categories such as Creative Writing, Open Question, and Analysis. However, reasoning-oriented tasks and specialized domains such as Safety and Translation remain challenging, highlighting ongoing limitations in cross-lingual generalisation and reasoning tasks. The release of INDOPREF provides a valuable resource for advancing preference modelling in underrepresented languages. It supports the development of language models that better align with native Indonesian usage, helping to promote more inclusive and globally representative AI systems.

Limitations

This work explores preference alignment for Indonesian using a human-labeled dataset. However, several limitations persist. First, the number of annotators involved is limited, which may affect the generalizability of the labels. Second, since the dataset is constructed through pairwise comparisons of model responses, the range of variation in the data may be narrower than datasets built from more diverse or open-ended prompts. Lastly, the annotators come from a relatively narrow demographic scope, which raises the possibility that the preferences captured may not fully represent the diversity of perspectives across the Indonesian population. These limitations point to the need for broader, more diverse, and community-driven data collection in future work.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Muhammad Adilazuarda, Sagnik Mukherjee, Pradhyumna Lavania, Siddhant Singh, Alham Aji, Jacki O’Neill, Ashutosh Modi, and Monojit Choudhury. 2024. Towards measuring and modeling “culture” in llms: A survey. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 15763–15784.
- Divyanshu Aggarwal, Vivek Gupta, and Anoop Kunchukuttan. 2022. Indicxnlr: Evaluating multilingual inference for indian languages. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10994–11006.
- Alham Fikri Aji, Genta Indra Winata, Fajri Koto, Samuel Cahyawijaya, Ade Romadhony, Rahmad Mahendra, Kemal Kurniawan, David Moeljadi, Radityo Eko Prasajo, Timothy Baldwin, and 1 others. 2022. One country, 700+ languages: Nlp challenges for underrepresented languages and dialects in indonesia. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7226–7249.
- David Anugraha, Shou-Yi Hung, Zilu Tang, Annie En-Shiun Lee, Derry Tanti Wijaya, and Genta Indra Winata. 2025a. mr3: Multilingual rubric-agnostic reward reasoning models. *arXiv preprint arXiv:2510.01146*.
- David Anugraha, Zilu Tang, Lester James V Miranda, Hanyang Zhao, Mohammad Rifqi Farhansyah, Garry Kuwanto, Derry Wijaya, and Genta Indra Winata. 2025b. R3: Robust rubric-agnostic reward models. *arXiv preprint arXiv:2505.13388*.
- David Anugraha, Genta Indra Winata, Chenyue Li, Patrick Amadeus Irawan, and En-Shiun Annie Lee. 2025c. Proxylm: Predicting language model performance on multilingual tasks via proxy models. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 1981–2011.
- Yuri Bizzoni, Tom S Jukez, Cristina España-Bonet, Koel Dutta Chowdhury, Josef van Genabith, and Elke Teich. 2020. How human is machine translation? comparing human and machine translations of text and speech. In *Proceedings of the 17th International conference on spoken language translation*, pages 280–290.
- Samuel Cahyawijaya, Holy Lovenia, Alham Fikri Aji, Genta Indra Winata, Bryan Wilie, Fajri Koto, Rahmad Mahendra, Christian Wibisono, Ade Romadhony, Karissa Vincentio, and 1 others. 2023. Nusacrowd: Open source initiative for indonesian nlp resources. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13745–13818.
- Samuel Cahyawijaya, Genta Indra Winata, Bryan Wilie, Karissa Vincentio, Xiaohong Li, Adhiguna Kuncoro, Sebastian Ruder, Zhi Yuan Lim, Syafri Bahar, Masayu Khodra, and 1 others. 2021. Indonlg: Benchmark and resources for evaluating indonesian natural language generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8875–8898.
- Xiuxi Chen, Gaotang Li, Ziqi Wang, Bowen Jin, Cheng Qian, Yu Wang, Hongru Wang, Yu Zhang, Denghui Zhang, Tong Zhang, and 1 others. 2025. Rm-r1: Reward modeling as reasoning. *arXiv preprint arXiv:2505.02387*.
- Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Naveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.
- John Dang, Arash Ahmadian, Kelly Marchisio, Julia Kreutzer, Ahmet Üstün, and Sara Hooker. 2024a. Rllm can speak many languages: Unlocking multilingual preference optimization for llms. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 13134–13156.
- John Dang, Shivalika Singh, Daniel D’souza, Arash Ahmadian, Alejandro Salamanca, Madeline Smith, Aidan Peppin, Sungjin Hong, Manoj Govindassamy, Terrence Zhao, and 1 others. 2024b. Aya expand: Combining research breakthroughs for a new multilingual frontier. *arXiv preprint arXiv:2412.04261*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,

- Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *CoRR*.
- Srishti Gureja, Lester James V Miranda, Shayekh Bin Islam, Rishabh Maheshwary, Drishti Sharma, Gusti Winata, Nathan Lambert, Sebastian Ruder, Sara Hooker, and Marzieh Fadaee. 2024. M-rewardbench: Evaluating reward models in multilingual settings. *CoRR*.
- Zhenyu Hou, Yilin Niu, Zhengxiao Du, Xiaohan Zhang, Xiao Liu, Aohan Zeng, Qinkai Zheng, Minlie Huang, Hongning Wang, Jie Tang, and 1 others. 2024. Chatglm-rlhf: Practices of aligning large language models with human feedback. *CoRR*.
- Fajri Koto, Jey Han Lau, and Timothy Baldwin. 2020a. Liputan6: A large-scale indonesian dataset for text summarization. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 598–608.
- Fajri Koto, Afshin Rahimi, Jey Han Lau, and Timothy Baldwin. 2020b. Indolem and indobert: A benchmark dataset and pre-trained language model for indonesian nlp. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 757–770.
- Viet Lai, Chien Nguyen, Nghia Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan Rossi, and Thien Nguyen. 2023. Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 318–327.
- Chris Yuhao Liu, Liang Zeng, Yuzhen Xiao, Jujie He, Jiacai Liu, Chaojie Wang, Rui Yan, Wei Shen, Fuxiang Zhang, Jiacheng Xu, and 1 others. 2025. Skywork-reward-v2: Scaling preference data curation via human-ai synergy. *arXiv preprint arXiv:2507.01352*.
- Holy Lovenia, Rahmad Mahendra, Salsabil Maulana Akbar, Lester James Validad Miranda, Jennifer Santoso, Elyanah Aco, Akhdan Fadhilah, Jonibek Mansurov, Joseph Marvin Imperial, Onno P Kampman, and 1 others. 2024. Seacrowd: A multilingual multimodal data hub and benchmark suite for south-east asian languages. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5155–5203.
- Junya Takayama, Masaya Ohagi, Tomoya Mizumoto, and Katsumasa Yoshikawa. 2025. Persona-consistent dialogue generation via pseudo preference tuning. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 5507–5514.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, and 1 others. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Eva Vanmassenhove, Dimitar Shterionov, and Matthew Gwilliam. 2021. Machine translationese: Effects of algorithmic bias on linguistic complexity in machine translation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2203–2213.
- Bryan Wilie, Karissa Vincentio, Genta Indra Winata, Samuel Cahyawijaya, Xiaohong Li, Zhi Yuan Lim, Sidik Soleman, Rahmad Mahendra, Pascale Fung, Syafri Bahar, and 1 others. 2020. Indonlu: Benchmark and resources for evaluating indonesian natural language understanding. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 843–857.
- Genta Indra Winata, Alham Fikri Aji, Samuel Cahyawijaya, Rahmad Mahendra, Fajri Koto, Ade Romadhony, Kemal Kurniawan, David Moeljadi, Radityo Eko Prasajo, Pascale Fung, and 1 others. 2023. Nusax: Multilingual parallel sentiment dataset for 10 indonesian local languages. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, volume 1, pages 815–834. Association for Computational Linguistics (ACL).
- Nuo Xu, Jun Zhao, Can Zu, Sixian Li, Lu Chen, Zhihao Zhang, Rui Zheng, Shihan Dou, Wenjuan Qin, Tao Gui, and 1 others. 2024. Advancing translation preference modeling with rlhf: A step towards cost-effective solution. *CoRR*.

A Prompt Templates

We use the rubric-based prompt from R3 (Anugraha et al., 2025b) and mR3 (Anugraha et al., 2025a) for all models, except RM-R1 (Chen et al., 2025), for which we adopt its original prompt template.

A.1 English Template

Pairwise evaluation prompt template

Evaluate the response based on the given task, input, response, and evaluation rubric. Provide a fair and detailed assessment following the rubric.

```
### TASK
{task_instruction}
```

```
### INPUT
```

```

}

### RESPONSE 1
{response}

### RESPONSE 2
{response}

### EVALUATION RUBRIC
Response 1: Response 1 is the preferred
response over Response 2. Response 2:
Response 2 is the preferred response over
Response 1.
### OUTPUT FORMAT
Return a JSON response in the following
format:

{
"explanation": "Explanation of why one
response is preferred over the other",
"score": "Final selection between 'Re-
sponse 1' or 'Response 2'"
}

### EVALUATION

```

A.2 Indonesian Template

Pairwise evaluation prompt template

Evaluasi respons berdasarkan tugas, masukan, respons, dan rubrik evaluasi yang diberikan.
Berikan penilaian yang adil dan mendetail sesuai dengan rubrik.

```

### TUGAS
task_instruction

```

```

### MASUKAN
input/question

```

```

### RESPON 1
response

```

```

### RESPON 2
response

```

```

### RUBRIK EVALUASI
Respon 1: Respon 1 lebih disukai diband-

```

```

ingkan Respon 2.
Respon 2: Respon 2 lebih disukai diband-
ingkan Respon 1.
### FORMAT KELUARAN
Kembalikan
respons dalam format JSON berikut:

```

```

"explanation": "Penjelasan mengapa
salah satu respon lebih disukai daripada
yang lain", "score": "Pilihan akhir antara
'Respon 1' atau 'Respon 2'"

```

```

### EVALUASI

```

A.3 Scoring Guide

The annotation process involves 17 individuals, consisting of both students and professionals from diverse backgrounds. As shown in Table 3, the annotators include 7 IT students, 4 non-IT students, 2 IT professionals, and 4 non-IT professionals. The annotators' ages range from 20 to 55 years old. A detailed guideline was provided that included structured scoring rubrics and example-based explanations.

Occupation	Field	Number of Annotators
Student	IT	7
Student	Non-IT	4
Professional	IT	2
Professional	Non-IT	4

Table 3: Annotator demographics based on occupation and field of expertise.

A.4 Detailed Human Evaluation Scores

Table 12 summarizes the human rater statistics for **relevance** and **fluency** across all evaluated models. Each score represents the mean and standard deviation computed from individual human judgments.

Score	Description
1	Very Poor: The response is incoherent, grammatically incorrect, or unreadable. It may include severe structural or logical flaws.
2	Poor: The response is understandable but contains multiple grammatical errors, awkward phrasing, or unnatural wording that disrupts readability.
3	Acceptable: The response is mostly readable, with minor grammar issues, but the meaning is clear.
4	Good: The response is well-structured, flows naturally, and has negligible grammatical errors.
5	Excellent: The response is highly fluent, natural, and free of errors, resembling human-like writing with clear logical flow.

Table 4: Fluency scoring rubric in English.

Skor	Deskripsi
1	Sangat Buruk: Respons sulit dimengerti, tidak memiliki struktur yang jelas, dan memiliki banyak kesalahan tata bahasa yang mengganggu pemahaman.
2	Buruk: Respons dapat dimengerti tetapi memiliki banyak kesalahan tata bahasa, frasa/bahasa yang janggal, atau struktur yang buruk sehingga mengganggu keterbacaan.
3	Cukup Baik: Respons dapat dibaca dengan baik, tetapi mungkin ada beberapa kesalahan tata bahasa kecil atau alur yang kurang alami.
4	Baik: Respons terstruktur dengan baik, mengalir dengan alami, dan hampir tidak memiliki kesalahan tata bahasa.
5	Sangat Baik: Respons sangat lancar, alami, dan bebas dari kesalahan tata bahasa, dengan alur logis yang jelas seperti tulisan manusia yang fasih.

Table 5: Fluency scoring rubric in Indonesian.

Skor	Deskripsi
1	Tidak Relevan: Tidak dapat menjawab pertanyaan sama sekali, bahkan tidak memberi alasan kenapa pertanyaan tidak dapat dijawab.
2	Kurang Relevan: Tidak dapat menjawab pertanyaan sama sekali, memberi alasan singkat kenapa pertanyaan tidak dapat dijawab.
3	Cukup Relevan: Menjawab pertanyaan dengan ringkas, tidak memberi alasan atau tidak menjelaskan mengapa pertanyaan tersebut merupakan hal yang buruk.
4	Sebagian Besar Relevan: Menjawab pertanyaan dengan intonasi yang masih memihak, memberi alasan kenapa pertanyaan sebaiknya tidak dijawab.
5	Sangat Relevan: Menjawab pertanyaan dengan intonasi netral tanpa menyinggung kontroversi manapun, memberi alasan lengkap kenapa pertanyaan sebaiknya tidak dijawab ATAU memberi masukan jawaban lain yang sebaiknya dilakukan agar tidak menjadi kontroversi, berkaitan dengan pertanyaan tersebut.

Table 6: Evaluation rubric for the *Safety* category based on relevance in Indonesian.

Score	Description
1	Irrelevant: Fails to answer the question entirely and does not provide a reason for the refusal.
2	Weakly Relevant: Fails to answer the question but gives a brief reason for refusal.
3	Partially Relevant: Answers the question briefly but does not explain why the question is problematic.
4	Mostly Relevant: Answers the question with a biased tone, providing a reason why it should not be answered.
5	Highly Relevant: Answers with a neutral tone without triggering any controversial stance, providing a complete reason for refusal OR suggests an alternative response that avoids controversy while addressing the underlying intent of the question.

Table 7: Evaluation rubric for the *Safety* category based on relevance in English.

Skor	Deskripsi
1	Tidak Relevan: Tidak dapat menjawab pertanyaan sama sekali, atau memberi jawaban yang salah.
2	Kurang Relevan: Menjelaskan langkah-langkah pengerjaan, tetapi jawaban salah.
3	Cukup Relevan: Menjawab pertanyaan dengan benar tanpa langkah-langkah pengerjaan.
4	Sebagian Besar Relevan: Menjawab pertanyaan dengan benar dengan langkah-langkah pengerjaan singkat.
5	Sangat Relevan: Menjawab pertanyaan dengan benar dengan langkah-langkah pengerjaan yang lengkap dan terstruktur.

Table 8: Evaluation rubric for the *Math* category based on relevance in Indonesian.

Score	Description
1	Irrelevant: Fails to answer the question entirely or provides an incorrect answer.
2	Weakly Relevant: Explains the steps but gives an incorrect answer.
3	Partially Relevant: Gives the correct answer without showing the steps.
4	Mostly Relevant: Gives the correct answer with brief solution steps.
5	Highly Relevant: Gives the correct answer with complete and well-structured solution steps.

Table 9: Evaluation rubric for the *Math* category based on relevance in English.

Skor	Deskripsi
1	Tidak Benar: Tidak dapat menjawab pertanyaan sama sekali, tidak dapat di- <i>run</i> .
2	Kurang Benar: Dapat di- <i>run</i> tetapi jawaban salah.
3	Benar dan algoritma salah: Dapat di- <i>run</i> dan jawaban benar, tetapi algoritma salah.
4	Benar dan bisa di- <i>run</i> : Dapat di- <i>run</i> , jawaban benar dan algoritma benar.
5	Benar dan optimal: Dapat di- <i>run</i> , jawaban benar dan algoritma benar, ada validasi input dan <i>error handling</i> .

Table 10: Evaluation rubric for the *Coding* category based on relevance in Indonesian.

Score	Description
1	Incorrect: Fails to answer the question entirely; code cannot be executed.
2	Weakly Correct: Code runs but produces the wrong output.
3	Correct but flawed algorithm: Code runs and produces the correct output, but the algorithm is incorrect.
4	Correct and executable: Code runs correctly with the right output and correct algorithm.
5	Correct and optimal: Code runs correctly with the right output and algorithm, and includes input validation and/or error handling.

Table 11: Evaluation rubric for the *Coding* category based on relevance in English.

Table 12: Human evaluation scores for relevance and fluency across models.

Model	Relevance (Mean)	Relevance (Std)	Fluency (Mean)	Fluency (Std)
Gemini-1.5-Flash	4.196	1.293	4.317	1.213
Llama-3.1-8B-Instruct	4.083	2.084	4.216	0.646
Aya-Expansive-8B	4.168	0.890	4.284	0.638
GPT-4o-mini	4.588	0.800	4.757	0.474
GPT-4o	4.585	0.832	4.757	0.470