

Speaker Identification and Dataset Construction Using LLMs: A Bilingual Case Study on Japanese and English Narratives

Anonymous ACL submission

Abstract

Speaker identification in narrative analysis is challenging due to complex dialogues, varying utterance patterns, and multiple characters with similar or ambiguous references. Accurately attributing utterances to the correct speakers is critical for understanding character interactions and the narrative structure.

To address these challenges, this study proposes a collaborative approach between humans and Large Language Models (LLMs) for dataset construction in speaker identification tasks. The process begins by manually extracting utterances and assigning speaker names to a small subset of the data. This labeled subset is then used to prompt-tune the LLM, enabling it to label speakers across the dataset. Subsequent manual corrections ensure accuracy while minimizing costs. Additionally, a paraphrased dataset is constructed to handle situations with multiple correct answers. Evaluation results indicate that models with larger parameter sizes, particularly those instruction-tuned in Japanese, achieve high accuracy in speaker identification.

1 Introduction

Narrative analysis is essential for understanding cultural values, psychological dynamics, and creative processes. By examining narrative structures and themes, we gain insights into societal norms and human behavior (Piper et al., 2021). Recent advancements in large language models (LLMs) (Zhao et al., 2023a) have opened new possibilities in fields like narrative analysis. LLMs can streamline tasks such as character emotion analysis and plot progression prediction.

Among narrative analysis tasks, speaker identification automatically attributing dialogue to the correct characters—is key. Accurate speaker identification is crucial for understanding character interactions and dynamics within a story, as it directly influences narrative interpretation.

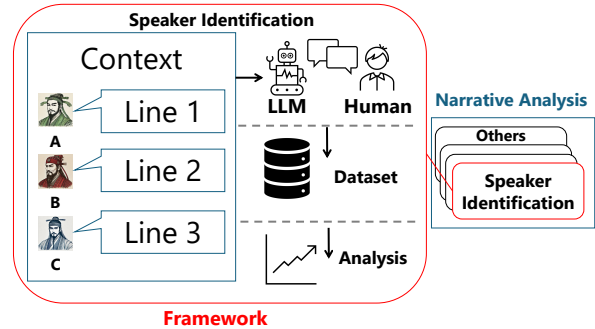


Figure 1: Method for constructing a dataset through collaboration between LLMs and human annotators for speaker identification in narrative analysis.

Traditionally, speaker identification has relied on machine learning models trained on datasets manually created by human annotators (Elson and McKeown, 2010; He et al., 2013; Muzny et al., 2017; Chen et al., 2019a; Vishnubhotla et al., 2022). However, creating high-quality datasets for speaker identification is labor-intensive and costly, as it requires careful consideration of consistency and paraphrase variations.

To address these challenges, this study applies a collaborative approach to dataset construction between LLMs and human annotators (Tan et al., 2024), specifically targeting the task of speaker identification in narrative analysis. By integrating LLMs for initial annotation followed by manual corrections, we aim to create a high-quality speaker identification dataset while significantly reducing the labor and costs. Our method annotates both main names and their paraphrased forms, inspired by the approach used in the PDNC dataset (Vishnubhotla et al., 2022), where both primary names and candidate paraphrases are manually annotated. This approach not only enhances the efficiency of speaker identification but also serves as a flexible framework applicable to other tasks in narrative analysis.

Existing speaker identification datasets have been limited to English and Chinese, restricting the generalizability of research findings to other languages. To overcome this limitation, our study introduces cross-lingual datasets developed from Wikisource¹ and Aozora Bunko², covering 14 diverse narratives across multiple languages. This approach not only advances the field of narrative analysis but is also valuable for evaluating the ability of LLMs to handle long contexts. Our results indicate that using only LLMs is possible to create datasets with approximately 80% accuracy in speaker identification, even across multiple languages (see Appendix O).

To further expand our dataset as a cross-lingual resource, we translated the Japanese version of “Romance of the Three Kingdoms” into English. This translation effort demonstrates an effective method for adapting datasets to multiple languages, thereby enhancing their applicability in cross-lingual studies.

2 Related Work

2.1 Dataset Construction

Elson and McKeown (2010) annotated speaker names and genders in 11 English narratives from the 19th century. He et al. (2013) treated separated lines in *Pride & Prejudice* as a single utterance for annotation. Muzny et al. (2017) expanded these datasets, creating the QuoteLi3 dataset, which includes annotations for all utterances in three narratives. Chen et al. (2019a) annotated utterances in the Chinese narrative World of Plainness (WP). Vishnubhotla et al. (2022) developed the Project Dialogism Novel Corpus (PDNC), annotating speakers, addressees, quote types, referring expressions, and mentions across 28 English novels, including main names and their variations.

Despite these advancements, existing datasets are primarily limited to English or Chinese, with no publicly available datasets for Japanese. Moreover, since these datasets depend on manual labor for annotation, they are inherently labor-intensive and costly to produce.

2.2 Speaker Identification

Feature-Based Approaches Several studies have employed linguistic features and manually crafted attributes for speaker identification (Elson

and McKeown, 2010; He et al., 2013; Bamman et al., 2014; Muzny et al., 2017).

Deep Learning Approaches With the advent of deep learning, more advanced methods for speaker identification have emerged. These include approaches that fine-tune models such as BERT (Bidirectional Encoder Representations from Transformers; (Devlin et al., 2019)), BART (Lewis et al., 2020) for speaker identification tasks (Cuesta-Lazaro et al., 2022; Vishnubhotla et al., 2023), and prompt tuning techniques with models such as GPT-3.5 (Ouyang et al., 2022) which have also demonstrated high accuracy on the Chinese WP dataset (Su et al., 2024).

These deep learning methods have improved the adaptability of speaker identification systems. However, they still face limitations that are related to the size of the context window. Michel et al. (2024) showed that while LLaMa3 (Dubey et al., 2024) expanded the context window and improved accuracy on the PDNC, their study was limited by the range of models and languages, leaving the evaluation incomplete.

3 Methods

Task Definition Speaker identification in narrative analysis involves determining which character or entity is responsible for a given utterance. This process requires analyzing both the utterance and its context to accurately attribute it to the correct speaker. In our approach, the set of possible speakers S is not predefined but derived from the context of the input text. Given a set of utterances $U = u_1, u_2, \dots, u_m$, we establish a mapping function $f : U \rightarrow S$ so that each utterance $u_i \in U$ is correctly attributed to a speaker $s_j \in S$. We annotated two types of speaker names: the “main name,” which is the most contextually appropriate (e.g., Elizabeth Bennet), and the “candidates,” which are variations or paraphrases of the same individual (e.g., Lizzy, Liz, Elizabeth). This dynamic speaker identification is crucial for capturing the fluid and complex nature of narrative interactions, enabling more accurate analysis of character relationships and narrative structure.

Prompt Tuning and Manual Correction To reduce the cost of creating a high-quality speaker identification dataset, we first manually created a small development dataset for speaker identification. We then applied prompt tuning using

¹https://wikisource.org/wiki/Main_Page

²<https://www.aozora.gr.jp/>

an LLM to generate speaker labels for the development data. Afterward, we manually corrected these LLM-generated speaker labels to ensure accuracy. This approach allowed us to maintain high data quality while significantly reducing the overall cost of dataset creation. We also used a chat template³ specifically designed for prompt tuning in a conversational format, employing a few-shot approach to enhance the LLM’s performance (see Appendix P). Additionally, this dataset included the identification of main names and candidate names for each speaker, ensuring comprehensive coverage of character references.

Cross-Lingual Dataset Creation We expanded our research to include cross-lingual datasets developed from Wikisource and Aozora Bunko, covering 14 diverse narratives across multiple languages. This approach offers a flexible and scalable framework for narrative analysis across various languages and cultural contexts, enhancing speaker identification by capturing the complexity of character references.

Robust Evaluation Metrics To ensure a robust evaluation of generation-based speaker identification systems like LLMs, we incorporated additional metrics such as substring match ratio and uncased evaluations. These metrics allow for a more flexible and accurate assessment of speaker identification performance by accounting for variations in text, thereby improving the reliability of the evaluation results.

4 Dataset Construction

The dataset construction was carried out according to the following steps, as shown in Figure 2.

STEP 1: Dialogue Extraction We initially extracted dialogues from *Aozora Bunko’s “Romance of the Three Kingdoms”* and Wikipedia sources by first tokenizing the data using the Llama-2 tokenizer and then extracting the surrounding 1,024-token contexts for each dialogue. This process yielded an initial dataset of 16,423 instances. The dataset is composed of 10 books, with book_id=52410 serving as the development data, and book_id=52411 to 52420 serving as the evaluation data (see Appendix O).

STEP 2: Speaker Labeling We utilized LLMs to identify and label the speakers in the extracted

dialogues. Speaker identification was performed on the dataset using a few-shot approach with Llama-3-70B-Instruct, which showed the highest performance on the development dataset (see Appendix B and M). During this phase, 1,011 instances were removed, resulting in a final dataset of 15,412 instances. The GPU was used for 200 hours for inference (see Appendix N).

STEP 3: Manual Correction Based on the identified labels, we manually corrected the speaker names following the annotation rules (see Appendix C). During this process, we corrected approximately 20% of the identified labels.

STEP 4: Translation We translated the dataset into English using GPT-4o-mini, specifically focusing on dialogues from “Romance of the Three Kingdoms” that were originally in Japanese (see Appendix E). Additionally, we used the GPT-4o-mini model for translation, including retry costs. The total translation cost was \$6.0 for processing 3,348 instances (book_id=52410, 52411), using a total of 30 million tokens.

This approach significantly reduced the time required for creating the evaluation data. While annotating 1,500 instances originally took approximately 10 hours, we reduced this time to 5 hours per 1,500 instances by focusing on correction tasks based on the identified speaker names. Our datasets are available at <https://huggingface.co/datasets/anonymized>.

4.1 Quality Assessment of Annotations

To verify the quality of the annotations, 100 samples from the evaluation dataset were reviewed by three independent annotators. They labeled the speaker names as “appropriate,” “inappropriate,” or “neutral,” and we calculated the agreement rates for the “appropriate” labels. The results showed high consistency, with two annotators achieving an agreement rate of 0.97 and one annotator achieving an agreement rate of 0.96 (see Appendix H).

4.2 Creation of Cross-Lingual Datasets Using Wikisource and Aozora Bunko

To facilitate cross-lingual analysis, we constructed cross-lingual datasets using texts from Wikisource and Aozora Bunko. We focused on 14 different stories and applied the same methodology as used in the dataset construction. Specifically, we annotated only the main names of characters in each story (see Appendix O).

³https://github.com/chujiezheng/chat_templates

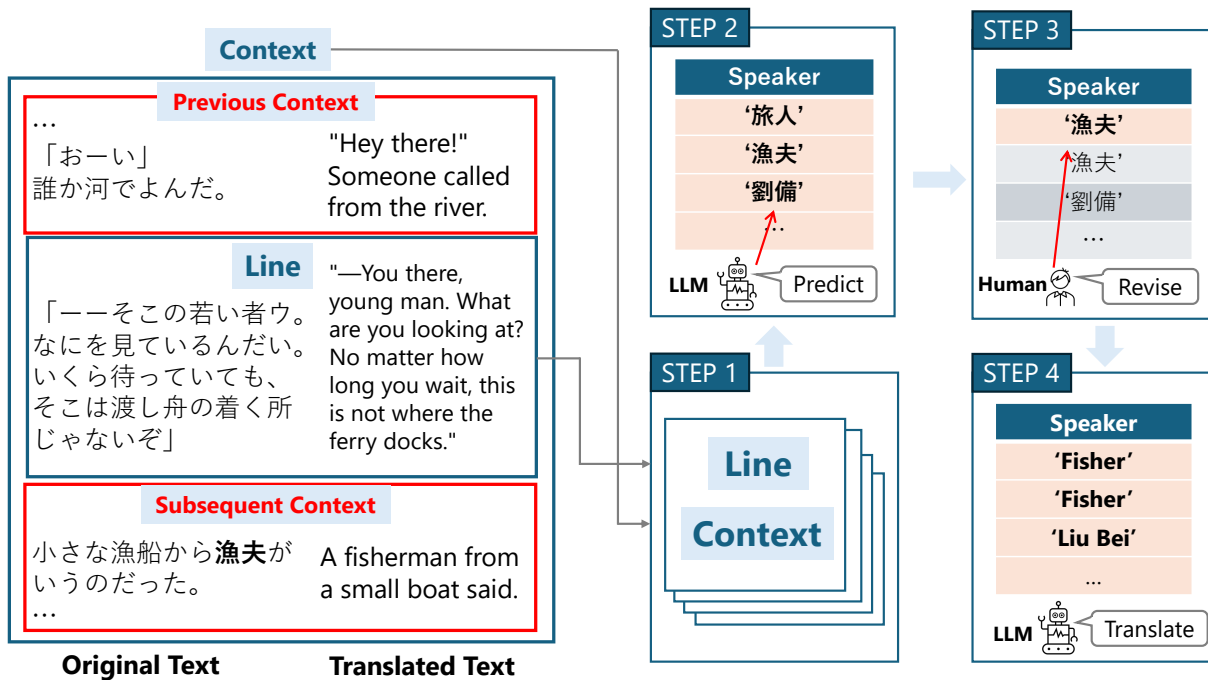


Figure 2: Workflow using Llama-3-70B-Instruct and GPT-4o-mini for labeling and translation.

5 Experiment

5.1 Prompt

We utilized LLMs to effectively perform speaker identification by providing few-shot examples through a chat template. In the chat template, we assigned the LLM a system role and guided it through the steps necessary to solve the task in a conversational format (see Appendix P).

5.2 Model

To compare model performance using LLMs, we selected LLaMa-3 (Dubey et al., 2024), a standard in LLM comparisons, along with Swallow-3 (Fujii, 2024), ELYZA-JP-8B (Hirakawa et al., 2024), and llama-3-youko-8B (Mitsuda et al.), all based on LLaMa-3 with additional Japanese training. For broader model evaluation, we included Mistral 7B (Jiang et al., 2023) and RakutenAI-7B (Group et al., 2024), which, like Mistral 7B, are trained on Japanese data. To assess the impact of training data composition on accuracy, we selected CALM-3-22B (Ishigami, 2024), primarily trained on Japanese data, and Karakuri-8x7B (Inc., 2024), which uses the Mixture of Experts technique (Jiang et al., 2024) (see Appendix M).

5.3 Evaluation Metrics

In this study, we used the following metrics to evaluate the accuracy of speaker attribution.

Exact Match Ratio This metric, commonly used in prior research (Vishnubhotla et al., 2023; Michel et al., 2024), measures the percentage of exact matches between the speakers identified in the generated text and those in the annotations.

Substring Match Ratio Given the variations in texts generated by LLMs, this metric recognizes partial matches in key elements of the speaker names (see Appendix A).

BERTScore (Zhang* et al., 2020) This metric assesses similarity based on embeddings, capturing cases where surface expressions differ but the underlying meaning remains the same.

Edit Distance (Levenshtein et al., 1966) Edit distance measures similarity by calculating the number of character insertions, deletions, and substitutions needed to convert one string into another.

Uncased Exact Match To account for case differences, the generated text is normalized for case-insensitive evaluation, treating “Old Woman” and “old woman” as equivalent. This metric is used only for English datasets.

5.4 Results

Overall Performance Figure 3 shows the overall books and the speaker identification accuracy for each model. In Fig. 3, the trend

Case	line	excerpt context	pred	true
A	Hahaha.	Yang Biao, harboring his secret plan, returned to his residence. As soon as he arrived, he went into his wife’s room and said, "So, how is it these days? Do you often meet with Lady Guo? I hear you ladies frequently have various gatherings." Placing his hands gently on his wife’s shoulders, he spoke with an unusual tenderness. Yang Biao’s wife, puzzled, teased him, "What’s gotten into you today? You’re never this sweet to me." "What’s the matter?" "Well, it’s just that you never act this way towards me normally." "Hahaha." "It actually makes me feel uneasy." "Is that so?"	Yang Biao	Yang Biao
B	Land of Jiangdong,	Wu is known as the "Land of Jiangdong," situated along the flow of the Great River.	Narration	Unknown
C	Diaochan, without showing any signs of agitation, immediately responded, "Yes. If it is the will of my lord, I am ready to give my life at any time." Wang Yun straightened his posture and said, "Then, I have something I wish to ask of you, trusting in your sincerity." "What is it?" "Dong Zhuo must be killed." "....." "If he is not removed, it will be as if the Han Emperor does not exist." "....."	Diaochan	Diaochan
D	The pleasures of life culminate here,	In the evening, a grand banquet was held with the slaughtering of cattle and horses for a feast. "The pleasures of life culminate here," said Guan Yu and Zhang Fei. "How could it end here? This is just the beginning," replied Xuande.	Guan Yu and Zhang Fei	Unknown
E	Lord Xuande, it is the fervent wish of both of us. Will you not consider it?	"It would be best." "Lord Xuande, it is the fervent wish of both of us. Will you not consider it?" From both sides,	Guan Yu	Guan Yu

Table 1: Case Study: 'Pred' indicates the predicted speaker, 'True' indicates the annotated speaker. Examples are translated into English; the original text is available in Appendix 2.

in results remained consistent across both the dev (book_id=52410) and eval (book_id=52411–52420) phases, maintaining an accuracy of approximately 90% (see Appendix B). The model that achieved the highest inference accuracy was the one that underwent continued pre-training on Japanese data using the base LLaMa-3 model. Applying instruction tuning after continued pre-training in Japanese appears to be effective for this task, suggesting that this combination was likely beneficial. The original LLaMa-3 model came next in performance.

Furthermore, comparing Swallow-3-8B-Instruct with Swallow-3-8B shows that instruction tuning improved performance by about 5%.

These findings suggest that while instruction tuning is effective for speaker identification, the performance of models trained with a sufficiently large parameter size approaches the upper limit for speaker identification accuracy.

Accuracy by Book To evaluate each model’s accuracy, we analyzed the substring match ratio for each book_id, focusing on the LLaMa-3-70B-Instruct model as an example. The LLaMa-3-70B-Instruct model identified

speaker names with an approximate accuracy of 0.9 across different book_ids, as shown in Figure 3, indicating consistent high accuracy in speaker identification.

For book_id=52419, the character “Sima Yi Zhongda” was referred to by different names, such as “Sima Yi” and “Zhongda”. During annotation, a rule prioritized the given name when present, leading to the frequent use of “Zhongda”. Consequently, the model sometimes identified the speaker as “Sima Yi,” who is the same individual. This suggests that the evaluation for this book_id may not fully reflect the model’s performance.

Relaxed Evaluation by Candidate Sets Using candidate sets for best matching allowed for relaxed evaluation, improving accuracy. In book_id=52419, “Sima Yi Zhongda” was referred to by various names, including “Sima Yi” and “Zhongda”.

According to annotation rules, “Zhongda” was used when it appeared in the context, and “Sima Yi” otherwise. Both names could serve as the main identifier. Following PDNC (Vishnubhotla et al., 2023), we prepared interchangeable candidate sets for “Zhongda,” including “Zhongda,” “Sima Yi,” “Sima Yi Zhongda,” and “Sima

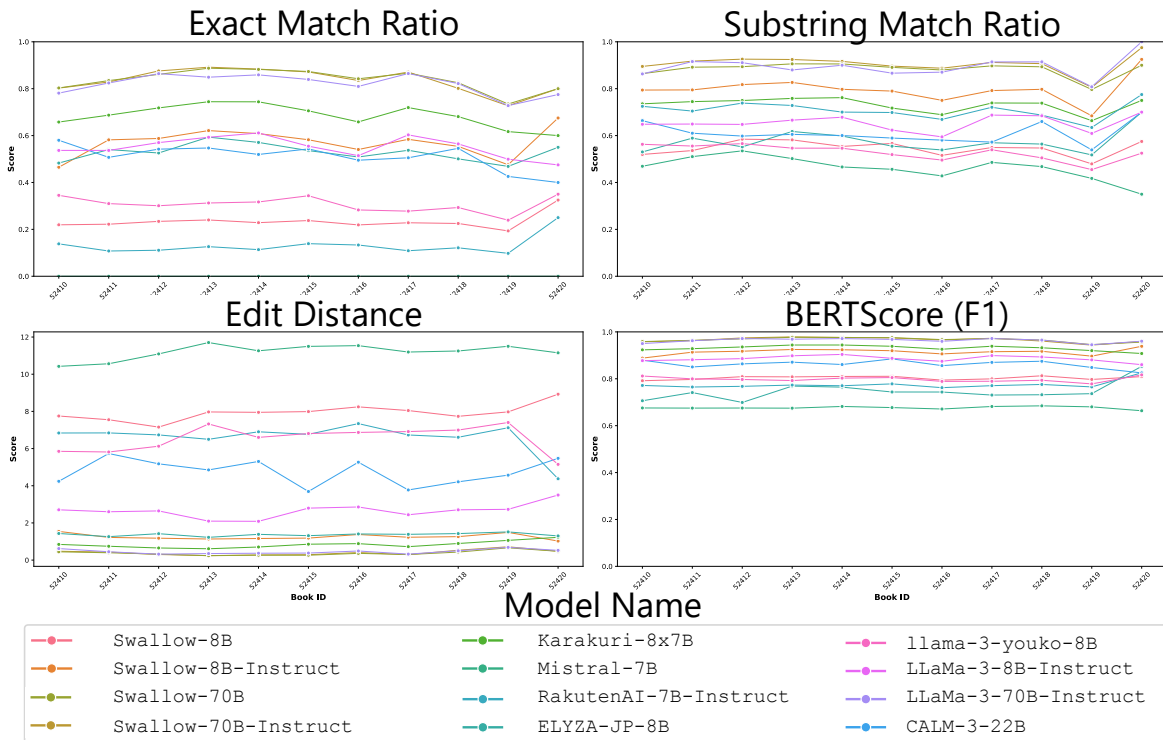


Figure 3: Results of Overall Performance for Various Models.

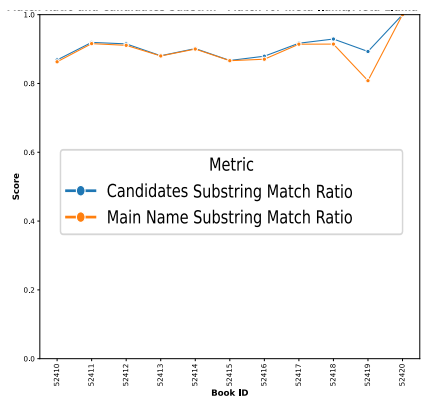


Figure 4: Comparison of the main name and its alternative candidates annotated through substring matching.

Zhongda”.

During evaluation, we matched predicted speaker names with the most corresponding name from the candidate sets. As shown in Figure 4, the substring match ratio using these sets was higher than in the initial evaluation. For book_id=52419, the evaluation became more consistent with the candidate sets. This suggests that flexibility in representing speaker names leads to more accurate evaluations.

Cross-Lingual Performance Figure 5 shows the substring match ratio for speaker identi-

cation using the LLaMa-3-70B-Instruct model on Japanese and English datasets. The model achieved higher accuracy on Japanese data, likely due to fewer label variations compared to English.

The Japanese dataset, composed mainly of simple folktales, exhibits fewer variations in referring terms. In contrast, the English dataset includes multiple synonyms for the same names, affecting the results. For example, the Japanese term “お母さん” in “matsuyama_kagami” is translated into various English terms, such as “Woman,” “Mother,” and “Wife”.

This suggests that, as noted in Section 5.4, preparing candidate sets for main names could reduce discrepancies. Additionally, to address case sensitivity issues in English, we used an Uncased Exact Match approach.

Evaluation of Models This section compares different models on the same story. Figure 6 shows the Llama-3-70B series demonstrated strong overall performance, with the Swallow-70B model showing particularly high accuracy, likely due to additional training on English data. However, when Swallow-70B is compared with Swallow-70B-Instruct, the former achieved better accuracy, suggesting that instruction tuning using Japanese data may have reduced

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404

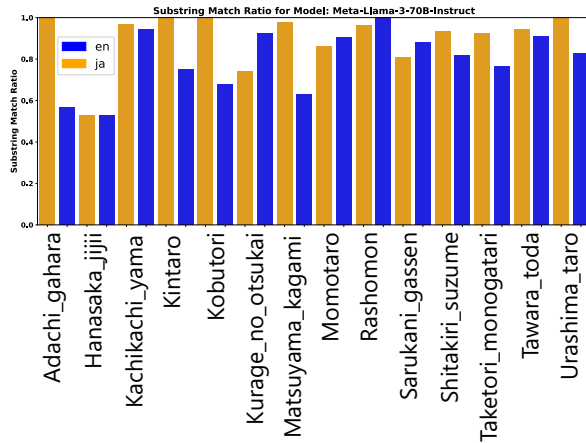


Figure 5: Comparison of substring match ratios for each story in both Japanese and English.

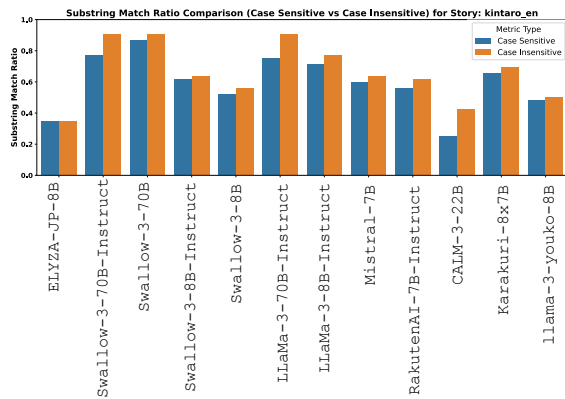


Figure 6: Comparison of Uncased Substring Match Ratio for story: kintaro_en.

the model’s generalization performance on English data.

Impact of Uncased Matching We address evaluation variations due to case sensitivity in English by incorporating an Uncased Exact Match metric. Figure 6 shows how case sensitivity affects evaluation by comparing uncased matching results for English and Japanese data, revealing that addressing it improves match ratio accuracy.

Applying the uncased match approach improved the substring match ratio for models like calm3-22b-chat and Llama-3-70B-Instruct. Additionally, Swallow-70B-Instruct closely matched Swallow-70B, suggesting that addressing case insensitivity reduces format variations, leading to more accurate model evaluation.

Performance on Translated Data In this section, we evaluate the performance of our model on the English version of “Romance of the Three

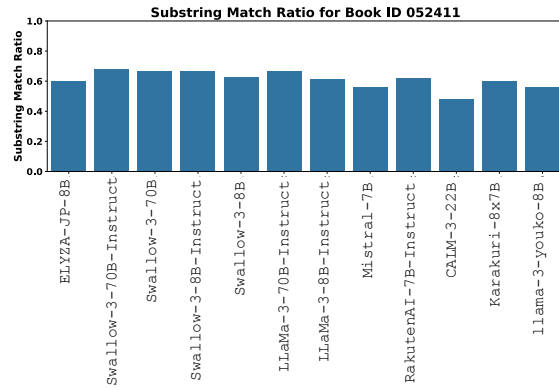


Figure 7: Substring Match Ratio Comparison Across Models for GPT-4o-mini Translated Data.

Kingdoms”. In Figure 7, the substring match ratio was approximately 70%, about 20% lower than the results obtained from the Japanese data. This reduction in accuracy is likely due to the inclusion of additional adjectives and other non-essential words in the English translation, making it more challenging to identify the core elements necessary for accurate speaker identification.

5.5 Analysis

Table 1 presents case study examples.

Case Study A: Long-Turn Dialogues The model generally identifies speakers accurately, even when relevant information is at the edges of the context. However, in case A, while the model correctly attributed “Hahaha.” to Yang Biao, it mistakenly attributed the following line, “Is that so?”, to his wife, indicating that errors were more likely in long-turn dialogues.

Case Study B: Narrator Identification We observed that the model correctly identifies the speaker as the narrator.

Case Study C: Silent Utterance Identification We confirmed the model can identify speaker names, even in implicit dialogues such as “.....”.

Case Study D: Multiple Speaker Identification We observed that the model correctly identifies the speaker even when multiple speakers are involved in the utterance.

Case Study E: Data Leak We analyzed potential data leakage by comparing ELYZA-JP-8B and LLaMa-3-70B-Instruct predictions with an 8-context length. While LLaMa-3-70B-Instruct inferred speaker names from the context,

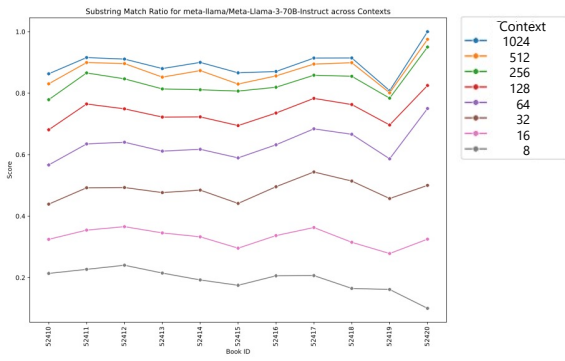


Figure 8: Variation in Substring Match Ratio by Context Length. This figure shows how the substring match ratio changes with different context lengths.

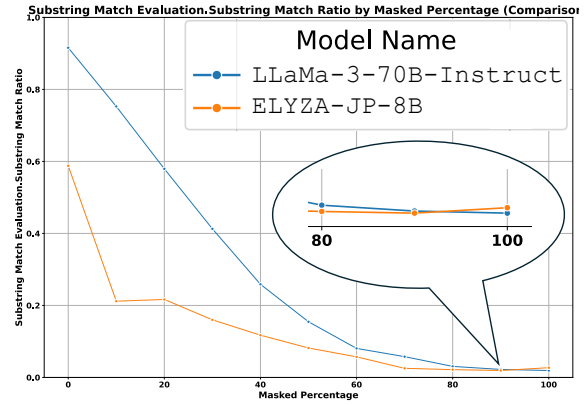


Figure 9: Substring Ratio by Mask Ratios

ELYZA-JP-8B correctly predicted speakers not explicitly mentioned. For example, ELYZA-JP-8B mistakenly identified “Guan Yu” as a speaker, likely due to reliance on prior knowledge triggered by the mention of “Xuande”.

Impact of Varying Context Lengths As illustrated in Figure 8, the LLaMa-3-70B-Instruct model’s identification accuracy improves with increasing context length but plateaus with minimal differences between 512 and 1024 lengths. Notably, models with smaller parameter sizes (8B or less) peaked at a context length of 512 (see Appendix Q).

This suggests that optimal context length is influenced by the model’s parameter size, indicating a dependency on computational capacity and design. Therefore, selecting an appropriate context length is crucial for maximizing performance, especially in resource-constrained environments (see Appendix O).

Impact of Context Masking We evaluated the effect of masking tokens within a 1,024-token context window on speaker identification accuracy. We tested the LLaMa-3-70B-Instruct model with mask ratios from 0% to 100% in 10% increments, replacing tokens with ‘<unk>’.

Figure 9 shows that accuracy declines as the Mask ratio increases. At 0% Mask, the model achieved 1.9% accuracy, which decreased as the Mask ratio increased.

The LLaMa-3-70B-Instruct model’s accuracy decreased with higher Mask ratios but still identified some speakers correctly. In contrast, the ELYZA-JP-8B model performed better at a 20% Mask ratio, indicating superior context retention. However, accuracy declined with excessive Masking due to reduced context.

At 100% Mask, the ELYZA-JP-8B model achieved a 2.7% match rate, surpassing the LLaMa-3-70B-Instruct model’s 1.9%. This suggests that the ELYZA-JP-8B model retains valuable contextual information even with full Masking (see Appendix I).

6 Conclusion

We collaborated with LLMs to create a speaker labeling dataset by annotating “Romance of the Three Kingdoms” from Aozora Bunko and 14 other stories in Japanese and English. The dataset included 15,412 entries and 1,017 annotations (517 in Japanese and 500 in English).

Using LLMs like LLaMa-3, we achieved a substring match ratio of approximately 90%. To handle multiple potential speakers, we developed a paraphrase dataset to improve evaluation accuracy.

We also used gpt-4o-mini for cross-lingual translation, enhancing annotation efficiency and reducing costs. This approach underscores the value of diverse datasets, adaptable evaluations, and LLM-assisted construction for effective, cost-efficient speaker identification, aiding narrative analysis and LLM development across languages.

7 Future Plans

We will advance narrative analysis by expanding multilingual datasets with advanced translation techniques and enhanced annotations, including Addressees and Quote Types, following the PDNC approach (Vishnubhotla et al., 2022). Additionally, we will develop improved speaker labeling methods and analyze complex stories with extensive character lists using enriched datasets. These efforts aim to enhance LLMs’ ability to handle intricate storytelling.

8 Limitations

Supported Languages This study can be extended to multiple languages by modifying the prompt methods and models used. The experiments, however, were limited to Japanese and English, and performance evaluations in other languages were not conducted.

While the results indicate high speaker identification performance in Japanese, comparing this with English presents certain challenges. English’s extensive use of synonyms and alternative expressions increases variations, making it difficult to draw direct comparisons between the two languages. These variations in English expressions might influence the results, highlighting the need for careful consideration when comparing performance across languages.

Translation In this study, we created a dataset translated using GPT-4o-mini for the purpose of cross-lingual evaluations. However, we only performed format checks on the translations (see Appendix E). To further enhance the quality of the dataset, human evaluation is deemed necessary.

9 Assurance of Research Ethics

We ensured adherence to research ethics by providing comprehensive explanations to the annotators about the study. Additionally, once the annotation was completed, we anonymized the collected data and paid careful attention to protecting personal information.

Furthermore, we verified the licenses for the artifacts, obtained the necessary approvals, and confirmed that our usage complies with the intended purposes.

References

David Bamman, Ted Underwood, and Noah A. Smith. 2014. [A Bayesian mixed effects model of literary character](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 370–379, Baltimore, Maryland. Association for Computational Linguistics.

Lasse Bergroth, Harri Hakonen, and Timo Raita. 2000. A survey of longest common subsequence algorithms. In *Proceedings Seventh International Symposium on String Processing and Information Retrieval. SPIRE 2000*, pages 39–48. IEEE.

Jia-Xiang Chen, Zhen-Hua Ling, and Li-Rong Dai. 2019a. [A Chinese Dataset for Identifying Speakers](#)

[in Novels](#). In *Proc. Interspeech 2019*, pages 1561–1565. 578 579

Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Wang, and Jia-Bin Huang. 2019b. A closer look at few-shot classification. In *International Conference on Learning Representations*. 580 581 582 583

Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46. 584 585 586

Carolina Cuesta-Lazaro, Animesh Prasad, and Trevor Wood. 2022. [What does the sea say to the shore? a BERT based DST style approach for speaker to dialogue attribution in novels](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5820–5829, Dublin, Ireland. Association for Computational Linguistics. 587 588 589 590 591 592 593 594

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics. 595 596 597 598 599 600 601 602 603

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esioibu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,

637	Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collet, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gouget, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwon Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papanikos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Barambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Laverder A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihalescu, Vladimir	701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764
-----	---	--

765	Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The llama 3 herd of models . <i>Preprint</i> , arXiv:2407.21783.	
775	David Elson and Kathleen McKeown. 2010. Automatic attribution of quoted speech in literary narrative . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 24(1):1013–1019.	
779	Kazuki Fujii. 2024. Llama-3-swallow: 日本語に強い継続事前学習モデル . https://zenn.dev/tokyotech_lm/articles/f65989d76baf2c .	
782	Kazuki Fujii, Taishi Nakamura, Mengsay Loem, Hiroki Iida, Masanari Ohi, Kakeru Hattori, Hirai Shota, Sakae Mizuki, Rio Yokota, and Naoaki Okazaki. 2024. Continual pre-training for cross-lingual llm adaptation: Enhancing japanese language capabilities . <i>Preprint</i> , arXiv:2404.17790.	
788	Ulrich Germann. 2003. Greedy decoding for statistical machine translation in almost linear time . In <i>Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology - Volume 1</i> , NAACL '03, page 1 – 8, USA. Association for Computational Linguistics.	
795	Rakuten Group, Aaron Levine, Connie Huang, Chenguang Wang, Eduardo Batista, Ewa Szymanska, Hongyi Ding, Hou Wei Chou, Jean-François Pesiot, Johannes Effendi, Justin Chiu, Kai Torben Ohlhus, Karan Chopra, Keiji Shinzato, Koji Murakami, Lee Xiong, Lei Chen, Maki Kubota, Maksim Tkachenko, Miroku Lee, Naoki Takahashi, Prathyusha Jwalapuram, Ryutaro Tatsushima, Saurabh Jain, Sunil Kumar Yadav, Ting Cai, Wei-Te Chen, Yandi Xia, Yuki Nakayama, and Yutaka Higashiyama. 2024. Rakutenai-7b: Extending large language models for japanese . <i>Preprint</i> , arXiv:2403.15484.	
808	Hua He, Denilson Barbosa, and Grzegorz Kondrak. 2013. Identification of speakers in novels . In <i>Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1312–1320, Sofia, Bulgaria. Association for Computational Linguistics.	
814	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding . In <i>International Conference on Learning Representations</i> .	
819	Masato Hirakawa, Shintaro Horie, Tomoaki Nakamura, Daisuke Oba, Sam Passaglia, and Akira Sasaki. 2024. elyza/llama-3-elyza-jp-8b .	
	KARAKURI Inc. 2024. KARAKURI LM 8x7B Instruct v0.1 .	822 823
	Ryosuke Ishigami. 2024. cyberagent/calm3-22b-chat .	824
	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b . <i>arXiv preprint arXiv:2310.06825</i> .	825 826 827 828 829
	Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L�elio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th�eophile Gervet, Thibaut Lavril, Thomas Wang, Timoth�ee Lacroix, and William El Sayed. 2024. Mixtral of experts . <i>Preprint</i> , arXiv:2401.04088.	830 831 832 833 834 835 836 837 838 839 840
	Vladimir I Levenshtein et al. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In <i>Soviet physics doklady</i> , volume 10, pages 707–710. Soviet Union.	841 842 843 844
	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880, Online. Association for Computational Linguistics.	845 846 847 848 849 850 851 852 853
	Gaspard Michel, Elena V. Epure, Romain Hennequin, and Christophe Cerisara. 2024. A realistic evaluation of llms for quotation attribution in literary texts: A case study of llama3 . <i>Preprint</i> , arXiv:2406.11380.	854 855 856 857 858
	Koh Mitsuda, Xinqi Chen, Toshiaki Wakatsuki, and Kei Sawada. rinna/llama-3-youko-8b .	859 860
	Grace Muzny, Michael Fang, Angel Chang, and Dan Jurafsky. 2017. A two-stage sieve approach for quote attribution . In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers</i> , pages 460–470, Valencia, Spain. Association for Computational Linguistics.	861 862 863 864 865 866 867
	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.	868 869 870 871 872 873
	Andrew Piper, Richard Jean So, and David Bamman. 2021. Narrative theory for computational narrative understanding . In <i>Proceedings of the 2021</i>	874 875 876

877 *Conference on Empirical Methods in Natural Lan-*
878 *guage Processing*, pages 298–311, Online and Punta
879 *Cana, Dominican Republic. Association for Compu-*
880 *tational Linguistics.*

881 Zhenlin Su, Liyan Xu, Jin Xu, Jiangnan Li, and
882 Mingdu Huangfu. 2024. [Sig: Speaker identifica-](#)
883 [tion in literature via prompt-based generation.](#) *Pro-*
884 [ceedings of the AAAI Conference on Artificial Intel-](#)
885 [ligence](#), 38(17):19035–19043.

886 Zhen Tan, Dawei Li, Song Wang, Alimohammad
887 Beigi, Bohan Jiang, Amrita Bhattacharjee, Man-
888 sooreh Karami, Jundong Li, Lu Cheng, and Huan
889 Liu. 2024. [Large language models for data annota-](#)
890 [tion: A survey.](#) *Preprint*, arXiv:2402.13446.

891 Krishnapriya Vishnubhotla, Adam Hammond, and
892 Graeme Hirst. 2022. [The project dialogism novel](#)
893 [corpus: A dataset for quotation attribution in lit-](#)
894 [erary texts.](#) In *Proceedings of the Thirteenth Lan-*
895 *guage Resources and Evaluation Conference*, pages
896 5838–5848, Marseille, France. European Language
897 Resources Association.

898 Krishnapriya Vishnubhotla, Frank Rudzicz, Graeme
899 Hirst, and Adam Hammond. 2023. [Improving au-](#)
900 [tomatic quotation attribution in literary novels.](#) In
901 *Proceedings of the 61st Annual Meeting of the As-*
902 *sociation for Computational Linguistics (Volume 2:*
903 *Short Papers)*, pages 737–746, Toronto, Canada. As-
904 sociation for Computational Linguistics.

905 Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q.
906 Weinberger, and Yoav Artzi. 2020. [Bertscore: Eval-](#)
907 [uating text generation with bert.](#) In *International*
908 *Conference on Learning Representations.*

909 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,
910 Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen
911 Zhang, Junjie Zhang, Zican Dong, et al. 2023a. [A](#)
912 [survey of large language models.](#) *arXiv preprint*
913 *arXiv:2303.18223.*

914 Zirui Zhao, Wee Sun Lee, and David Hsu. 2023b.
915 [Large language models as commonsense knowledge](#)
916 [for large-scale task planning.](#) In *Thirty-seventh Con-*
917 *ference on Neural Information Processing Systems.*

918 A Substring Match Ratio Evaluation

919 Method

920 The substring match ratio evaluates whether the
921 true speaker name, as annotated, exists as a sub-
922 string within the predicted speaker name. This
923 evaluation metric is mathematically formalized as
924 follows:

925 **Definitions** In a given dialogue dataset, we de-
926 fine the speaker names as follows:

- 927 • P_i : Predicted speaker name
- 928 • T_i : Annotated true speaker name

We define the match function M as:

$$M(P_i, T_i) = \begin{cases} 1 & \text{if there exists an integer } j \\ & \text{such that } 0 \leq j \leq |P_i| - |T_i| \\ & \text{and } P_i[j : j + |T_i|] = T_i \\ 0 & \text{otherwise} \end{cases}$$

Calculation of Substring Match Ratio The
931 substring match ratio for the entire dataset is cal-
932 culated as the proportion of dialogues where the
933 true speaker name is a substring of the predicted
934 speaker name. Formally, it is defined as:
935

$$r_s = \frac{1}{n} \sum_{i=1}^n M(P_i, T_i)$$

where $n \in N$ is the total number of dialogues.

938 Calculation Steps

- 939 1. For each dialogue i , check if the true speaker
940 name T_i is a substring of the predicted
941 speaker name P_i .
- 942 2. Assign $M(P_i, T_i) = 1$ if T_i is a substring of
943 P_i ; otherwise, assign $M(P_i, T_i) = 0$.
- 944 3. Calculate the sum of all $M(P_i, T_i)$ values and
945 divide by the total number of dialogues n .

Example Consider three dialogues with the fol-
946 lowing predicted and true speaker names:
947

- 948 • $P_1 = \text{“John Smith”}$, $T_1 = \text{“John”}$
- 949 • $P_2 = \text{“Alice”}$, $T_2 = \text{“Bob”}$
- 950 • $P_3 = \text{“Charlie Brown”}$, $T_3 = \text{“Charlie”}$

The substring matches are calculated as follows:

$$\begin{aligned} M(P_1, T_1) &= 1, \\ M(P_2, T_2) &= 0, \\ M(P_3, T_3) &= 1 \end{aligned}$$

Thus, the substring match ratio is calculated as:

$$r_s = \frac{1}{3}(1 + 0 + 1) = \frac{2}{3} \approx 0.67$$

955 Using the substring match ratio, we can eval-
956 uate how accurately the predicted speaker names
957 contain the true speaker names as substrings.

958 Particularly, LLMs often generate unnecessary
959 texts, such as special tokens like “[INST]” and un-
960 related tokens.

961	B Detailed Dataset Extraction and		
962	Segmentation Process		
963	Data Extraction	The data was meticulously extracted from <i>Aozora Bunko's "Romance of the Three Kingdoms"</i> using the Huggingface datasets ⁴ library. This curated dataset includes furigana and metadata, and was selected for its extensive character list and the potential to extract complex relationships.	
964			
965			
966			
967			
968			
969			
970	Development and Evaluation Sets	The dataset was split into development and evaluation sets as follows:	
971			
972			
973		• Volume 02: Peach Garden Oath (Shinjitai, Book ID: 52410) served as the development set.	
974			
975			
976		• Volume 03: Among the Stars (Shinjitai, Book ID: 52411) to Volume 11: Wuzhang Plains (Shinjitai, Book ID: 52419) constituted the evaluation set.	
977			
978			
979			
980	C Annotation Rules		
981	The following annotation rules were applied for label assignment:		
982			
983		1. As a general principle, the smallest constituent part of a character's name used in the narrative text is considered the correct label. (Example: For "劉備玄德", "玄德" is the correct label.)	
984			
985			
986			
987			
988		2. When multiple candidates exist, the given name is preferred if it is present in the context.	
989			
990			
991		3. If the text is not a dialogue, label it as 'Unknown'. (Examples: characters, narrator, book titles)	
992			
993			
994		4. If multiple speakers are indicated for a single utterance, label it as 'Unknown'. (Examples: Guan Yu, Zhao Yun, Liu Bei)	
995			
996			
997		5. Due to the high preparation cost, dynamic generation based on reading the context is preferred, as annotators had prior access to speaker information.	
998			
999			
1000			
			6. Each utterance, along with the preceding and following 1,024 tokens, is set as the context. Only the names found within this context are subject to annotation. The number of tokens is calculated based on the Llama-2 Tokenizer ⁵ .
			1001
			1002
			1003
			1004
			1005
			1006
		7. If multiple names representing a single person appear in the context, the most appropriate one is labeled as the "main name," while other possible names are labeled as "candidates."	1007
			1008
			1009
			1010
			1011
		8. List candidates for each main name in a dictionary format. Include various expressions, such as courtesy names or official titles, in the candidates list.	1012
			1013
			1014
			1015
		For each main name, the presence of candidates in the context is checked, and a set of potential names is automatically generated.	1016
			1017
			1018
	D Comparison of Paraphrase Set Acquisition with Wikipedia Redirects		1019
			1020
		The Romance of the Three Kingdoms is well-known, resulting in many of its characters having dedicated Wikipedia pages. Therefore, using Wikipedia Redirects ⁶ to acquire paraphrase expressions is conceivable. However, our attempts revealed that paraphrase expressions could only be acquired for some characters.	1021
			1022
			1023
			1024
			1025
			1026
			1027
		Specifically, excluding the names we extracted as Main Name from our created paraphrase sets, only 1.83% of paraphrase candidates could be obtained using Wikipedia Redirects. Notably, expressions corresponding to "劉備" as "青年" or "應德" as "旅人" could not be obtained.	1028
			1029
			1030
			1031
			1032
			1033
		These results indicate the limitations of using Wikipedia Redirects for acquiring paraphrase expressions. Hence, combining other methods and data sources is essential for comprehensive paraphrase collection.	1034
			1035
			1036
			1037
			1038
	E Translation of Annotated Datasets		1039
		To align the annotated datasets with their English counterparts, we utilized OpenAI's GPT-4o-mini ⁷ model following a detailed translation process. The procedure involved several steps:	1040
			1041
			1042
			1043
		⁵ https://huggingface.co/meta-llama/Llama-2-7b-hf	
		⁶ https://en.wikipedia.org/wiki/Wikipedia:Redirect	
		⁷ https://platform.openai.com/docs/models	

1044	Translation Procedure	We employed the GPT-4o-mini model to translate the annotated datasets.	F Expenses for Translation	1092
1045				
1046		For book_id=052410, the translation resulted	Format checks are crucial in the translation process, ensuring that each translation conforms to	1093
1047		in 1,574 entries, while book_id=052411 pro-	the specified format. If a translation fails the format check, up to five retries are performed, po-	1094
1048		duced 1,528 entries. In the evaluation phase	tentially increasing the number of tokens used for	1095
1049		for book_id=052410, we skipped portions of the	input and output. The number of tokens used di-	1096
1050		prompts that included contextual information.	rectly impacts the expenses. In this study, we used	1097
1051	Prompt Design	For the datasets created in Section 4, we used few-shot prompts to translate dia-	the gpt-4o-mini model for translation, consuming	1098
1052		logues, contexts, and speaker information (see Ap-	30 million tokens, resulting in a total translation	1099
1053		pendix P). These prompts were crafted to guide the	cost of \$6.0.	1100
1054		model in generating accurate translations by offer-		1101
1055		ing relevant examples.		1102
1056			G Original Japanese Text of Case Study	1103
1057	Quality Assurance	The translations were evaluated based on several criteria:	Table 2 presents the original Japanese text of the	1104
1058			case study discussed (see Section 5.5).	1105
1059			H Detailed Quality Assessment of	1106
1060			Annotations	1107
1061				
1062			In this study, all annotations were independently	1108
1063			performed by the first author, making it impossi-	1109
1064			ble to directly evaluate inter-annotator agreement.	1110
1065			To verify the quality of the created annotations, we	1111
1066			randomly selected 100 samples from the evalua-	1112
1067			tion dataset and asked three independent annota-	1113
1068			tors to review them.	1114
1069			The annotators were tasked with evaluating the	1115
1070			labeled speaker names as “appropriate,” “inappro-	1116
1071			priate,” or “cannot judge”. We assigned weights	1117
1072			to these evaluations: 3 points for “appropriate,” 2	1118
1073			points for “cannot judge,” and 1 point for “inap-	1119
1074			propriate”. The agreement was calculated based	1120
1075			on these weighted scores using a three-point Lik-	1121
1076			ert scale.	1122
1077			The results showed that two annotators had an	1123
1078			agreement rate of 0.97, and one annotator had an	1124
1079			agreement rate of 0.96, indicating a very high level	1125
1080			of consistency. This suggests that the dataset con-	1126
1081			structed in this study is of high quality.	1127
1082			Typically, Cohen’s kappa coefficient (Cohen,	1128
1083			1960) is used to evaluate inter-annotator agree-	1129
1084			ment. However, in this case, the agreement rates	1130
1085			were so high that setting the original data labels	1131
1086			to 3 when calculating the kappa coefficient could	1132
1087			lead to undefined values. Therefore, we report	1133
1088			only the agreement rate and its variance (see Ap-	1134
1089			pendix J for details).	1135
1090			Additionally, the annotation task required an av-	1136
1091			erage of 2 hours per annotator, with a compensa-	1137
			tion rate set at 1,000 yen per hour. The annotations	1138
			were performed by three native Japanese graduate	1139

Case	line	excerpt context	pred	true
A	あははは	楊彪は秘策を胸にねりながら、わが邸へ帰って行った。帰るとすぐ、彼は妻の室へは行って、「どうだな。この頃は、郭汜の令夫人とも、時々お目にかかるかね。……おまえたち奥さん連ばかりで、よく色々な会があるとのことだが」と、両手を妻の肩にのせながら、いつになく優しい良人になって云った。二 楊彪の妻は怪しんで、良人を揶揄した。「あなた。どうしたんですか、いったい今日は」「なにが?」「だって、常には、私に対して、こんなに機嫌をとるあなたではありませんもの」「あははは」「かえって、気味が悪い」「そうかい」	楊彪	楊彪
B	江東の地	呉は、大江の流れに沿うて、「江東の地」と称われている。	不明 (ナレーション)	Unknown
C	……………	貂蟬は、さわぐ色もなく、すぐ答えた。「はい。大人のおたのみなら、いつでもこの生命は捧げます」王允は、座を正して、「では、おまえの真心を見込んで頼みたいことがあるが」「なんですか」「董卓を殺さねばならん」「……………」「彼を除かなければ、漢室の天子はあってもないのと同じだ」「……………」	貂蟬	貂蟬
D	人生の快、ここに尽くる	夜は、牛馬を宰して、聚議の大歓宴が設けられた。「人生の快、ここに尽くる」関羽、張飛がいうと、「何でこれに尽きよう。これからである」と、玄德はいった。	関羽、張飛	Unknown
E	玄德様、ふたりの熱望です。ご承知くださるまいか	たほうがよい 「玄德様、ふたりの熱望です。ご承知くださるまいか」 左右から	関羽	関羽

Table 2: Original Case Study in Japanese. ‘pred’ indicates the predicted speaker label, and ‘true’ indicates the annotated speaker label.

Metric	Annotator ID		
	A	B	C
Agreement Rate	0.97	0.97	0.96
Count (3)	97	97	96
Count (2)	3	2	3
Count (1)	0	1	1
Total	100	100	100
Weighted Average Score	2.97	2.96	2.95

Table 3: Annotation agreement and evaluation distribution by annotator. The "Agreement Rate" represents the proportion of cases where independent evaluators marked the data as "appropriate" (3) when the author had labeled it as 3 in the dataset. The "Count (x)" rows indicate the number of times each annotator selected "appropriate" (3), "neutral" (2), or "inappropriate" (1). The "Total" row indicates that each annotator evaluated 100 cases. The "Weighted Average Score" reflects the average score calculated by assigning weights of 3, 2, and 1 to the respective categories.

students, selected for their advanced language proficiency, further contributing to the reliability and accuracy of the data.

I Further Case Study

Table 4 shows that ELYZA-JP-8B had already read these datasets during the training steps.

This finding indicates that the ELYZA-JP-8B

model may have leveraged learned patterns or relationships to make accurate predictions even when the context is heavily Masked.

J Challenging Cases in Annotation Judgment

Table 5 presents examples where annotation decisions were particularly challenging.

Examining the final portion of the context in Table A, it is evident that the character “張飛” strongly asserts that “呂布” must be defeated. This suggests that the preceding conversation was primarily conducted by “玄德” and “張飛”. Therefore, considering the immediate context, it is highly likely that the line in question was spoken by “張飛”.

However, reading the previous tokens reveals that the line “何事を曹操からいってよこしたのですか” could be attributed to both “張飛” and “関羽”. Consequently, there is a slight possibility that “関羽” could have responded to “玄德”’s statement, “まあ、これを見るがいい”.

Two of the independent annotators employed to assess annotation quality provided feedback suggesting that the possibility of “関羽” being the speaker could not be entirely ruled out. Such cases, where reaching a consensus on the speaker annotation was extremely difficult, were reported by the annotators three or four times per 100 cases.

id	line	excerpt context	pred	true
1869	ですから、父上のお顔で、富豪を紹介して下さい。曹家は、財産こそないが、遠くは夏侯氏の流れを汲み、漢の丞相曹参の末流です。この名門の名を利用して、富豪から金を出させて下さい		曹操	曹操

Table 4: Correct Identification of an Absent Name : ELYZA-JP-8B accurately predicts the name “曹操,” despite it not being present in the context.

id	line	excerpt context	true	corr	incor	neu
3818	呂布を殺せという密命ですな	<p>何度も、繰り返し繰り返し読み直していると、後ろに立っていた張飛、関羽のふたりが、「何事を曹操からいつてよこしたのですか」と、訊ねた。</p> <p>「まあ、これを見るがいい」</p> <p>「呂布を殺せという密命ですな」</p> <p>「そうじゃ」</p> <p>「呂布は、兇勇のみで、もともと義も欠けている人間ですから、曹操のさしずをよい機として、この際、殺してしまうがよいでしょう」</p> <p>「いや、彼はたのむ所がなく、わが懐に投じてきた窮鳥だ。それを殺すは、飼禽を縊るようなもの。玄德こそ、義のない人間といわれよう」</p> <p>「――が、不義の漢を生かしておけば、ろくなことはしませんぞ。国に及ぼす害は、誰が責めを負いますか」</p> <p>「次第に、義に富む人間となるように、温情をもって導いてゆく」</p> <p>「そうやすやす、善人になれるものですか」</p> <p>張飛は、あくまでも、呂布討つべしと主張したが、玄德は、従う色もなかった。</p>	張飛	1	0	2

Table 5: Challenging Annotation Example. ‘true’ indicates the predicted speaker label. ‘corr’ indicates the number of annotators who judged the annotated label to be correct, ‘incor’ indicates those who judged it to be incorrect, and ‘neu’ indicates those who judged it to be neutral. This example illustrates a difficult case where the three independent annotators had differing opinions, highlighting the complexity and subjectivity involved in the annotation process.

K Token Count Variations

Figure 10 shows the maximum input token count per book_id, confirming that the actual number of input tokens in this study falls within 8,192 tokens when converted using the Llama 3 Tokenizer. As illustrated in Figure 10, this study employed the Llama 2 Tokenizer to extract the preceding and following 1,024 tokens, thereby creating context tokens. Among the tokenizers used in the comparative models, the most commonly utilized base tokenizer was the Llama 3 Tokenizer.

Furthermore, Figure 11 demonstrates the variation in token count per index for book_id=052415, which had the highest number of input tokens. Excluding a few exceptionally long dialogue examples, almost all token counts were distributed around 2,250 tokens using the Llama 2 Tokenizer and around 1,500 tokens using the Llama 3 Tokenizer.

Reducing the length of the input context or randomly masking it was confirmed to significantly decrease identification accuracy (see Section 5.5

and Section 5.5). Therefore, to solve this task with high accuracy, it is necessary to process a sufficiently long context of at least 1,500 tokens using the Llama 3 Tokenizer.

This indicates that the number of tokens handled is extremely large compared to the methods used for evaluating the performance of existing LLMs, such as MMLU (Hendrycks et al., 2021) and Commonsense (Zhao et al., 2023b). By addressing this task, it is believed that we can measure the inference performance of LLMs with respect to long contexts.

Additionally, in this study, the dataset length was set to fit within the maximum input token count of 8,192 tokens, which is the limit for the models used in comparison. For identification tasks using similar methods, simply increasing the length of the input context or simultaneously targeting multiple dialogues for speaker identification could easily extend the evaluation to tasks requiring longer contexts, such as those involving 100,000 tokens.

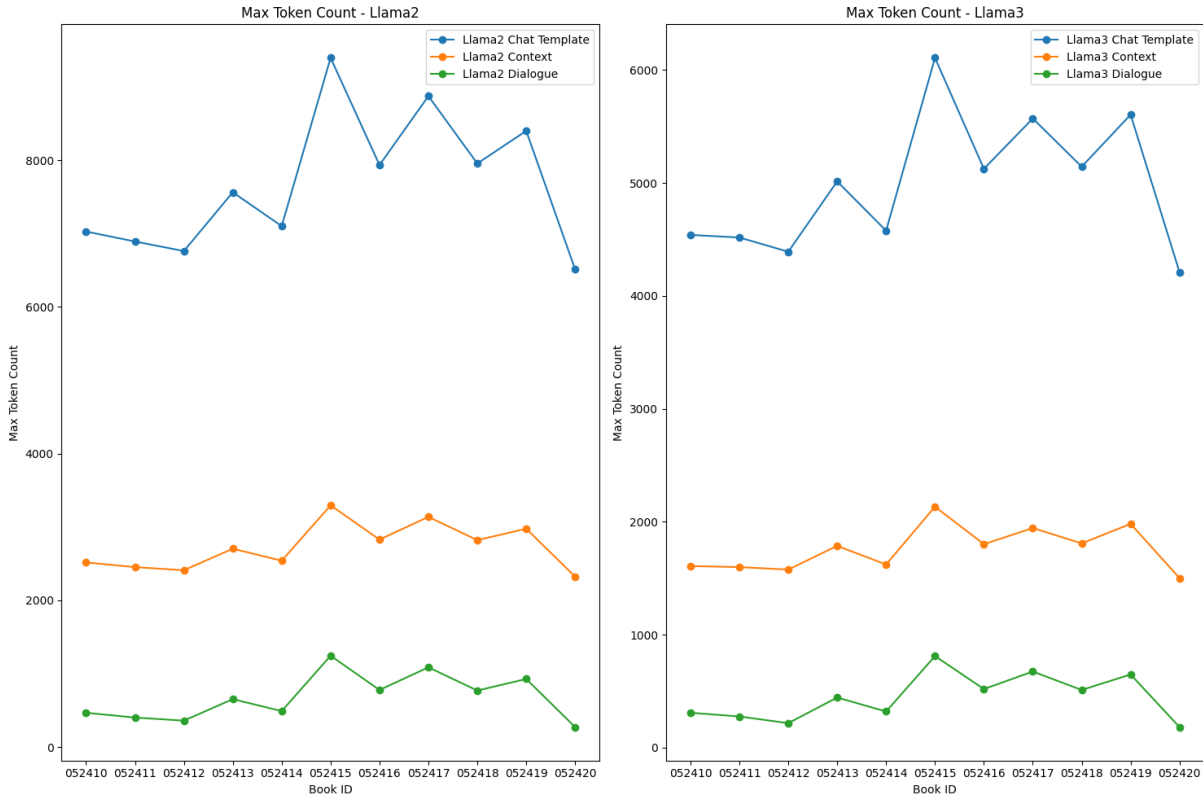


Figure 10: The Chat Template indicates the maximum token count when including tokens that control few-shots and prompt format. Context shows the maximum token count when inferring speaker names and combining the target dialogue with the preceding and following 1,024 tokens. Dialogue shows the maximum token count for the dialogue itself.

L Content Warning for Violent Expressions

This dataset contains stories written several decades ago, during a period when violent expressions and provocative language, including depictions of murder and aggressive behavior, were more commonplace. Users are advised to exercise caution and be mindful of the potentially disturbing content when utilizing this dataset.

M Model Description

The selection criteria for each model aim to comprehensively evaluate performance across various languages and tasks, adaptation to Japanese data, and differences between architectures. This allows for a multifaceted assessment of LLM performance.

In this study, we selected 12 models for comparison, organized into six categories. Below is a description of each model and the rationale for its selection.

LLaMa-3 (Dubey et al., 2024) LLaMa-3 is an LLM that considers human preferences, demonstrating high performance in various tasks such as multilingual support, coding, and mathematics. It is also used as a base model for many other models, making it suitable for comparative validation.

Swallow-3 (Fujii et al., 2024) Swallow-3 is a model based on LLaMa-3 that has undergone continual pretraining and instruction tuning with Japanese data. It was selected to analyze changes in Japanese performance and potential performance degradation in English data relative to LLaMa-3.

ELYZA-JP-8B (Hirakawa et al., 2024) ELYZA-JP-8B is a model based on LLaMa-3 that has undergone continual pretraining and instruction tuning with Japanese data. We selected this model to evaluate whether instruction tuning leads to differences when compared to Swallow-3.

llama-3-youko-8B (Mitsuda et al.) llama-3-youko-8B is a model based on LLaMa-3 that has

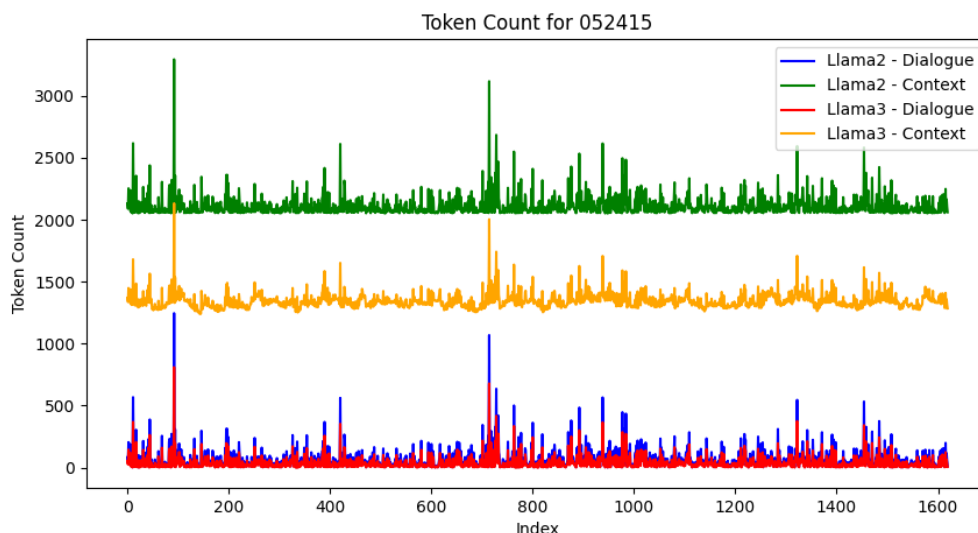


Figure 11: Variation in token count per index for book_id=052415. Excluding exceptionally long dialogues, most token counts are distributed around 2,250 tokens based on the Llama 2 Tokenizer and around 1,500 tokens based on the Llama 3 Tokenizer.

1260 undergone continual pretraining using a mixture
1261 of Japanese and English datasets.

1262 **Mistral-7B (Jiang et al., 2023)** Mistral-7B, like
1263 LLaMa-3, is frequently used for comparisons with
1264 other models and is known for its high perform-
1265 ance despite its smaller size. It was selected
1266 to compare a model from a different lineage to
1267 LLaMa-3.

1268 **RakutenAI-7B (Group et al., 2024)**
1269 RakutenAI-7B is a model fine-tuned with
1270 Japanese data based on Mistral 7B. It was selected
1271 to compare the performance of models fine-tuned
1272 with Japanese data, similar to Swallow-3.

1273 **CALM-3-22B (Ishigami, 2024)** CALM-3-22B
1274 is an LLM primarily trained on proprietary
1275 Japanese data. It was selected to compare the per-
1276 formance of models that mainly handle Japanese
1277 data with those that support multiple languages,
1278 primarily focusing on English.

1279 **Karakuri-8x7B (Inc., 2024)** Karakuri-8x7B is
1280 a model that uses a Mixture of Experts (MoE) ap-
1281 proach by combining multiple models for more ef-
1282 fective inference, specifically Mixtral-8x7B (Jiang
1283 et al., 2024), and has undergone continual pre-
1284 training and fine-tuning with Japanese data. It
1285 was selected to compare MoE models with other
1286 LLMs.

N Inference and Evaluation Setup 1287

1288 In this study, we set the random seed at 42 and per-
1289 formed 4-bit quantization for model inference. We
1290 used the Greedy Decoding Algorithm (Germann,
1291 2003) for decoding. Inference was conducted us-
1292 ing an A6000 GPU, with a total inference time of
1293 approximately 200 hours.

1294 During evaluation, unnecessary strings, such
1295 as special tokens [INST] generated by the LLM,
1296 were removed using regular expressions wherever
1297 possible.

1298 Additionally, various libraries were utilized for
1299 inference, evaluation, and visualization. For ex-
1300 ample, we employed scikit-learn⁸, transformers⁹,
1301 beautifulsoup4¹⁰, tiktoken¹¹, openai¹², evaluate¹³,
1302 accelerate¹⁴, torch¹⁵, datasets¹⁶, and matplotlib¹⁷.

O Number of Tokens and Speakers 1303

1304 Table 6 shows the number of tokens (based on the
1305 Llama-2 and Llama-3 base models), lines, unique
1306 speakers, skip, and line_ids for each book_id. Ad-
1307 ditionally, the number of unique speakers, exclud-

⁸<https://scikit-learn.org/>

⁹<https://github.com/huggingface/transformers>

¹⁰<https://beautiful-soup-4.readthedocs.io/>

¹¹<https://github.com/openai/tiktoken>

¹²<https://github.com/openai/openai-python>

¹³<https://github.com/huggingface/evaluate>

¹⁴<https://github.com/huggingface/accelerate>

¹⁵<https://github.com/pytorch/pytorch>

¹⁶<https://github.com/huggingface/datasets>

¹⁷<https://matplotlib.org/>

ing duplicates in the annotated speaker names, is confirmed to be 856.

Table 7 summarizes the number of tokens, utterances, and characters for each story.

In this table, “Tokens (Llama-3, JA)” and “Tokens (Llama-3, EN)” indicate the number of tokens in the Japanese and English versions of each story, respectively. Similarly, “Lines (JA)” and “Lines (EN)” represent the number of utterances in Japanese and English, respectively.

P Prompt Configuration

Predict Quoted Speech Listings 1 and 2 show the prompts used for speaker identification. As shown in these listings, we provide several few-shot examples in a chat format. The prompt consists of text extracted from the beginning of book_id=052410 included in Aozora Bunko. In Listings 2, few-shot examples (Chen et al., 2019b) related to the story in Listing 1, along with the target story (Context) and are provided the utterance line (Line) for speaker identification.

Using these prompts, we constructed a dataset to evaluate the accuracy of speaker identification and conducted speaker identification based on this dataset.

Translation Similar to speaker identification, we configured these prompts, including few-shot examples, for translation. Additionally, we incorporated prompts that included failure cases (see Table 9).

Q Impact of Varying Context Lengths with Other Models

Figures 12–13 illustrate the accuracy of substring matches when varying the input context length across different models.

As shown in these figures, models with approximately 70B parameters exhibited improved speaker identification accuracy as the context length increased. Conversely, for models with 8B parameters or fewer, accuracy plateaued when the context length was extended from 256 to 512 tokens. Beyond this point, providing additional context resulted in a performance decline due to the introduction of noise, with the extent of the decline varying across models.

These observations suggest that the effective context length for input varies depending on the model’s parameter size and training methodology.

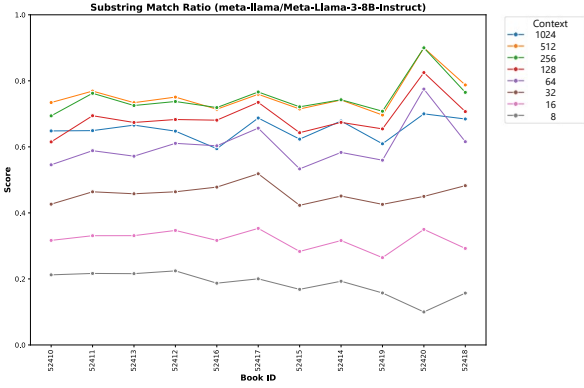


Figure 12: Variation in Substring Match Ratio by Context Length for LLaMa-8B-Instruct. This figure shows how the substring match ratio changes with different context lengths.

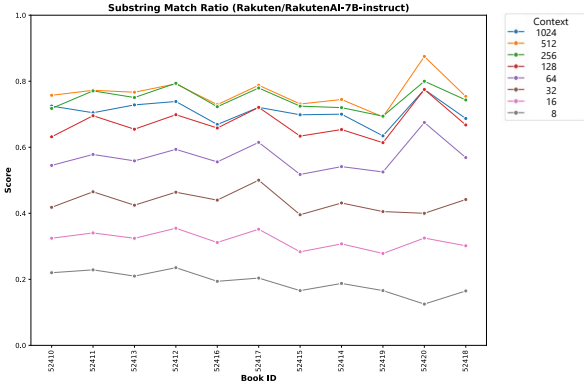


Figure 13: Variation in Substring Match Ratio by Context Length for RakutenAI-7B-Instruct. This figure shows how the substring match ratio changes with different context lengths.

book_id	tokens (Llama-2)	tokens (Llama-3)	lines	unique speakers	skip	line_id
052409	1,866	1,129	0	0	2	1-2
052410	195,226	124,143	1,686	113	70	3-1,758
052411	195,589	124,772	1,662	157	108	1,759-3,528
052412	193,973	124,364	1,649	136	129	3,529-5,306
052413	201,042	129,000	1,616	123	82	5,307-7,004
052414	205,799	131,796	1,461	159	89	7,005-8,554
052415	209,759	133,797	1,532	117	88	8,555-1,0174
052416	204,514	130,989	1,598	153	83	10,175-11,855
052417	222,992	143,735	1,433	171	95	11,856-13,383
052418	249,258	159,547	1,426	186	96	13,384-14,905
052419	223,710	143,901	1,308	122	130	14,906-16,343
052420	27,050	16,968	40	26	40	16,344-16,423
Total	2,130,778	1,364,141	15,411	1,463	1,012	1-16,423

Table 6: Number of Tokens and Speakers by Dataset

Story	Tokens	Tokens	Lines	Lines	Skip	Skip
	(Llama-3, JA)	(Llama-3, EN)	(JA)	(EN)	(JA)	(EN)
Shita-kiri Suzume	2,838	3,256	46	22	1	2
Tawara Toda	2,035	2,823	18	11	0	1
Urashima Taro	4,036	5,272	36	69	0	3
Kachikachi Yama	3,175	2,842	58	17	1	0
Kintaro	2,816	3,920	30	52	1	6
Taketori Monogatari	5,452	6,680	27	17	0	0
Matsuyama Kagami	2,839	6,219	40	46	0	0
Adachigahara	2,479	2,083	17	23	0	0
Hanasaka Jijii	2,237	3,339	19	19	2	2
Kurage no Otsukai	2,837	3,728	58	67	0	0
Saru Kani Kassen	2,498	3,256	42	17	0	0
Momotaro	4,031	5,361	58	83	9	1
Rashomon	2,176	2,730	26	32	4	0
Kubu-tori	3,539	2,579	42	25	0	0
Total	42,988	54,088	517	500	18	15

Table 7: Summary of token and utterance counts for both Japanese and English versions of each story. Annotation was performed on the main names of characters, following the methodology used in constructing the dataset for the Japanese version of “Romance of the Three Kingdoms” (see Section 4).

type	prompt
Japanese Example Story	<p>後漢の建寧元年のころ。今から約千七百八十年ほど前のことである。一人の旅人があった。腰に、一剣を佩いているほか、身なりはいたって見すばらしいが、眉は秀で、唇は紅く、とりわけ聡明そうな眸や、豊かな頬をしていて、つねにどこかに微笑をふくみ、総じて賤しげな容子がなかった。年の頃は二十四、五。草むらの中に、ぼつねんと坐って、膝をかかえこんでいた。悠久と水は行く――微風は爽やかに鬢をなでる。涼秋の八月だ。そしてそこは、黄河の畔の――黄土層の低い断り岸であった。「おーい」誰か河でよんだ。「――そこの若い者ウ。なにをえているんだい。いくら待っていても、そこは渡し舟の着く所じゃないぞ」小さな漁船から漁夫がいうのだった。青年は笑くぼを送って、「ありがとう」と、少し頭を下げた。漁船は、下流へ流れ去った。けれど青年は、同じ所に、同じ姿をしていた。膝をかかえて坐ったまま遠心的な眼をうごかさなかった。「おい、おい、旅の者」こんどは、後ろを通った人間が呼びかけた。近村の百姓であろう。ひとりは鶏の足をつかんでさげ、ひとりは農具をかついでいた。「――そんな所で、今朝からなにを待っているんだね。このごろは、黄巾賊とかいう悪徒が立ち廻るからな。役人衆に怪しまれるぞよ」青年は、振りかえって、「はい、どうも」おとなしい会釈をかえした。</p>
English Example Story	<p>In the first year of the Jianning era of the Later Han Dynasty. This was about one thousand seven hundred and eighty years ago. There was a traveler. Apart from wearing a sword at his waist, his appearance was quite shabby. However, he had prominent eyebrows, red lips, especially intelligent-looking eyes, and full cheeks that always seemed to hold a smile, overall giving him an air that was not at all lowly. He appeared to be around twenty-four or twenty-five years old. He was sitting alone in a patch of grass, hugging his knees. Time flows like the eternal river—A gentle breeze brushed his sideburns. It was August, a cool autumn month. And this was the bank of the Yellow River—on a low clay cliff. "Hey there!" Someone called from the river. "—You there, young man. What are you looking at? No matter how long you wait, this is not where the ferry docks." A fisherman from a small boat said. The young man smiled and, "Thank you," he said with a slight nod. The fishing boat drifted downstream. But the young man stayed in the same spot, in the same posture, his eyes still looking into the distance. "Hey, hey, traveler." This time, someone passing by from behind called out. It seemed to be a farmer from a nearby village. One was holding a chicken by its feet, and the other was carrying farming tools. "—What have you been waiting for since this morning in a place like this? Lately, there have been bandits called the Yellow Turbans around. The officials will get suspicious of you." The young man turned and, "Yes, thank you," he replied with a gentle nod.</p>

Table 8: Example Stories

Listing 1: Example Chat Template (JA)

```

chat = [
  {"role": "user", "content": "次の物語 (# Example) 中のセリフに対して、誰が発言したのかを推測して、話者の名前のみを教えてください。Story"},
  {"role": "assistant", "content": "了解しました。以下の物語とセリフに基づいて答えを提供します。"},
  {"role": "user", "content": f"#{Example_Story}"},
  {"role": "assistant", "content": "物語を確認しました。では、セリフごとに誰が発言したのかを答えます。"},
  {"role": "user", "content": "次の発話は誰が発言しましたか?"},
  {"role": "assistant", "content": "セリフを教えてください。"},
  {"role": "user", "content": "おーい"},
  {"role": "assistant", "content": "漁夫"},
  {"role": "user", "content": "—そこの若い者ウ。なにを見ているんだい。いくら待っていても、そこは渡し舟の着く所じゃないぞ"},
  {"role": "assistant", "content": "漁夫"},
  {"role": "user", "content": "ありがとう"},
  {"role": "assistant", "content": "青年"},
  {"role": "user", "content": "おい、おい、旅の者"},
  {"role": "assistant", "content": "百姓"},
  {"role": "user", "content": "—そんな所で、今朝からなにを待っているんだね。このごろは、黄巾賊とかいう悪徒が立ち廻るからな。役人衆に怪しまれるぞよ"},
  {"role": "assistant", "content": "百姓"},
  {"role": "user", "content": "同様にして、次の物語 (# Target) 中のセリフに対して、誰が発言したのかを推測して、話者の名前のみを教えてください。Story"},
  {"role": "assistant", "content": "了解しました。以下の物語とセリフに基づいて答えを提供します。"},
  {"role": "user", "content": f"#{Target_Story}"},
  {"role": "assistant", "content": "物語を確認しました。では、セリフごとに誰が発言したのかを答えます。"},
  {"role": "user", "content": "次の発話は誰が発言しましたか?"},
  {"role": "assistant", "content": "セリフを教えてください。"},
  {"role": "user", "content": f"#{Line}"},
]

```

Listing 2: Example Chat Template (EN)

```

chat = [
  {"role": "user", "content": "Please guess who is speaking each line of dialogue in the following story (# Example Story) and provide only the speaker's name."},
  {"role": "assistant", "content": "Understood. I will provide answers based on the story and dialogues below."},
  {"role": "user", "content": f"#{Example_Story}"},
  {"role": "assistant", "content": "I have reviewed the story. Now, I will identify the speaker for each line of dialogue."},
  {"role": "user", "content": "Who said the following line?"},
  {"role": "assistant", "content": "Please provide the line of dialogue."},
  {"role": "user", "content": "Hey there!"},
  {"role": "assistant", "content": "Fisherman"},
  {"role": "user", "content": "—You there, young man. What are you looking at? No matter how long you wait, this is not where the ferry docks."},
  {"role": "assistant", "content": "Fisherman"},
  {"role": "user", "content": "Thank you,"},
  {"role": "assistant", "content": "Young Man"},
  {"role": "user", "content": "Hey, hey, traveler."},
  {"role": "assistant", "content": "Farmer"},
  {"role": "user", "content": "—What have you been waiting for since this morning in a place like this? Lately, there have been bandits called the Yellow Turbans around. The officials will get suspicious of you."},
  {"role": "assistant", "content": "Farmer"},
  {"role": "user", "content": "Similarly, guess who is speaking each line of dialogue in the following story (# Target Story) and provide only the speaker's name."},
  {"role": "assistant", "content": "Understood. I will provide answers based on the story and dialogues below."},
  {"role": "user", "content": f"#{Target_Story}"},
  {"role": "assistant", "content": "I have reviewed the story. Now, I will identify the speaker for each line of dialogue."},
  {"role": "user", "content": "Who said the following line?"},
  {"role": "assistant", "content": "Please provide the line of dialogue."},
  {"role": "user", "content": f"#{Line}"},
]

```

R Use of AI Tools in Writing and Coding

We used AI tools to assist in the writing and coding processes for this project. Specifically, we employed ChatGPT¹⁸ to help draft and refine the text, and we utilized GitHub Copilot¹⁹ for code completion and suggestions during the coding tasks. These tools were incorporated into our workflow to support the efficient completion of the project.

¹⁸<https://openai.com/chatgpt/>

¹⁹<https://docs.github.com/en/copilot>

type	prompt
Speaker	<p>Translate the following speaker's name into English, using terms that appear in the translated context. Provide the translation only:</p> <p>Example 1: Translated context: "The farmer walked through his fields, greeting the old man sitting by the road." Output: old man</p> <p>Example 2: Translated context: "In the small village, the young woman was known for her kindness." Output: young woman</p> <p>Example 3: Translated context: "The wise elder spoke to the gathered crowd with great wisdom." Output: wise elder</p>
Dialogue	<p>Extract the entire line that is most similar to this dialogue: 'original_dialogue', excluding the quotation marks. Ensure to extract the full sentence from the start to the end.</p> <p>Example 1: Original dialogue: "これからどうする?" Translated context: "They looked at each other, wondering about the next steps. One of them asked, 'What are we going to do now?' Another responded, 'We need to think carefully.'" Extracted line: What are we going to do now?</p> <p>Example 2: Original dialogue: "何を言えればいいかわからない。" Translated context: "He scratched his head, lost for words. He finally said, 'I have no idea what to say.' Another person nodded in agreement, 'It's a tough situation.'" Extracted line: I have no idea what to say.</p> <p>Failure Example 1: Original dialogue: "こっちへ行こう。" Translated context: "They were considering their options. One said, 'Let's go this way.' Another said, 'I think we should stay here.'" Extracted line: I think we should stay here. # The extracted line is incorrect as it does not match the original dialogue's intent to move.</p>
Context	<p>Translate the following context into English, ensuring consistency and that the provided dialogue is included. The translation should maintain a coherent narrative flow. Provide the translation only:</p> <p>Example 1: Original context: "彼は暗闇の中で独り、静かな夜の音を聞いていた。その時、彼は『おい、誰かいるのか?』と呼びかけた。" Translated dialogue: "Hey, is anyone there?" Translated context: "He sat alone in the darkness, listening to the quiet sounds of the night. At that moment, he called out, 'Hey, is anyone there?'"</p> <p>Example 2: Original context: "彼女は辺りを見回し、そして『ここに何かあるの?』と尋ねた。周りには何もないようだった。" Translated dialogue: "What's here?" Translated context: "She looked around and then asked, 'What's here?' There seemed to be nothing around."</p>

Table 9: Prompts for translation