

MaLLaM - Malaysia Large Language Model

Husein Zolkepli* Aisyah Razak† Kamarul Adha‡ Ariff Nazhan§

January 26, 2024

Abstract

Addressing the gap in Large Language Model pretrained from scratch with Malaysian context, We trained models with 1.1 billion, 3 billion, and 5 billion parameters on a substantial 349GB dataset, equivalent to 90 billion tokens based on our pretrained Byte Pair Encoding (BPE) tokenizer for a single epoch. MaLLaM contributes to enhanced natural language understanding and generation tasks in the Malay language.

Although trained on a smaller dataset of 90 billion tokens, our instruction-tuned MaLLaM models perform competitively. When compared to ChatGPT3.5 and Malaysian Mistral, MaLLaM's instruction-tuned models demonstrate notable proficiency, underscoring the effectiveness of our approach in capturing and understanding the nuances of the Malaysian language.

MaLLaM models mark a significant contribution to the field, providing comprehensive language representations grounded in Malaysian context. This endeavor aims to pave the way for enhanced natural language understanding and generation tasks specific to the linguistic nuances present in Malaysia. We discuss the training methodology, dataset composition, and the potential impact of MaLLaM in advancing the capabilities of large language models within the context of the Malay language.

All models released at [HuggingFace Mesolitica MaLLaM Collection](#).

1 Introduction

The landscape of large language models (LLMs) has predominantly been shaped by models trained on English, with subsequent adaptations and fine-tunings for languages like Tamil and Malaysian. Existing models such as Tamil-LLama [1] and Malaysian-Mistral [2] have emerged as valuable assets, leveraging the foundation of English-based LLMs for further optimization in non-English languages. However, despite their commendable contributions, these models still carry traces of English nuances, presenting a unique challenge in achieving a fully language-specific representation.

This introduction underscores the significance of mitigating residual English influences in language models tailored for specific linguistic contexts. While the adaptation of English-based LLMs has facilitated advancements in various languages, it falls short of capturing the intricacies and nuances unique to languages like Tamil and Malaysian. To address this gap, we embark on a novel approach by pre-training a large language model entirely from scratch, with a specific focus on the Malaysian language. This deliberate endeavor aims to establish a model that is inherently attuned to the linguistic subtleties and idiosyncrasies of Malaysian, overcoming the limitations posed by models with English-centric origins.

The presence of such inherent English biases becomes particularly pertinent when envisioning the creation of a dedicated large language model for Malaysian linguistic contexts. The prevailing influence of English-centric sources, prevalent in public English news and articles, poses a potential challenge to the development of a truly indigenous language model. The need for a linguistic model devoid of these biases

*husein@mesolitica.com

†aisyahrazak171@gmail.com

‡kamarul.adha360@gmail.com

§ariffnzhn@gmail.com

is underscored by the desire to accurately capture the unique nuances and intricacies of the Malaysian language.

In light of these considerations, our initiative introduces MaLLaM, a large language model specifically pre-trained from scratch using a robust dataset equivalent to 90 billion tokens sourced from Malaysian contexts. The distinctiveness of MaLLaM lies in its genesis, free from the shadow of English-centric biases pervasive in existing language models. By cultivating a language model that is inherently attuned to Malaysian linguistic idiosyncrasies, we aim to address the gap left by models that, despite being fine-tuned for Malaysian languages, retain traces of English nuances. MaLLaM stands as a testament to our commitment to fostering linguistic authenticity and overcoming the challenges posed by the dominance of English-centric biases in current language modeling paradigms.

- **Pre-training MaLLaM:** We utilized a powerful infrastructure consisting of 10 nodes of the Standard_ND96amsr_A100_v4 Azure instance, with each node featuring 8 A100 80 GPUs. This configuration efficiently facilitated the pre-training of language models with 1.1 billion, 3 billion, and 5 billion parameters on a substantial 349GB dataset of Malaysian texts.
- **Multi-turn Instruction-Tuned MaLLaM:** To ensure a seamless and meaningful comparison, we opted to employ the exact chat instruction dataset from Malaysian Mistral [2]. This approach enables us to replicate the same experimental setup, facilitating a direct and accurate comparison of our work with the existing model.

2 Pre-Training Data

2.1 Public Data

Malay Wikipedia, we incorporated the Malay Wikipedia dump, enriching the linguistic diversity by partially converting it into Jawi script. Utilizing the [Ejawi converter](#), we ensured that the model is not only proficient in understanding the standard Malay script but also adept at comprehending the intricacies of Jawi.

Malay Language study articles, We enriched our dataset with content from the Malay dictionary and public articles from Dewan Bahasa Pustaka. This inclusion ensures our language model is well-versed in word meanings, usages, and various linguistic styles prevalent in Malaysian literature.

Malaysia Government public documents, We included public government documents from the official Malaysia government website and Google searches to enrich our dataset. This ensures that our language model is exposed to formal language and communication styles used in official government contexts, enhancing its ability to understand and generate text relevant to administrative, legal, and policy domains in Malaysia.

Malaysian public articles, We gathered a varied dataset by scraping articles from Malaysian sources. This process exposed our language model to different topics, writing styles, and language nuances found in Malaysian articles, contributing to a model well-versed in the Malaysian language across diverse domains.

Malaysian public social media, We gathered diverse data by scraping content from specific Facebook pages, filtering tweets based on location and keywords, and extracting information from platforms like c.cari.com.my, b.cari.com.my, carigold, lowyat, and transcripts of Malaysian YouTube videos. This approach broadens our dataset to include various language styles and colloquial expressions from online Malaysian communities.

Malaysia public journals, We added content from reputable Malaysian journals like <https://mjpharm.org>, <https://myjgeosc.com>, and <https://www.akademisains.gov.my> to our dataset. This ensures our language model is familiar with formal and technical language used in academic contexts, covering a variety of subjects.

Malaysian related public research papers, We refined our dataset by filtering CrossRef using keywords like 'malaysia,' 'malay,' and 'melayu.' This targeted approach ensures that our language model is exposed to scholarly literature closely tied to Malaysia.

Complete list of gathered data at [Github Project - Prepare LLM dataset](#).

2.2 Coding Data

Incorporating a coding dataset is a crucial component of our diverse training data. For this purpose, we utilized the original dataset available at [bigcode/the-stack-dedup](#). To ensure relevance and efficiency, we selectively picked specific programming languages, including Python, Julia, C, C++, HTML, CSS, JavaScript, Go, Rust, Java, SQL, Markdown, R, Dockerfile, Ruby, Typescript, and YAML. To manage the dataset size, each programming language was limited to a maximum of 10GB.

2.3 Instruction-tuned Data

To augment the comprehensiveness of our training dataset, we devised a synthetic Malay instruction dataset. This encompassed diverse linguistic aspects, including the conversion between Rumi and Jawi scripts, dependency and constituency parsing, grammatical error generation (kesalahan tatabahasa), and coding instruction datasets related to various programming paradigms. Additionally, we included instructional content relevant to educational levels such as UPSR, PT3, and SPM, providing a broad coverage of language proficiency levels. The incorporation of syntactic and semantic elements, coupled with coding instructions, contributes to a well-rounded language model capable of handling a diverse array of linguistic tasks and understanding instructions across different domains.

Complete list of Instruction-tuned Data at [instruction-dataset](#).

3 Deduplicating and Postprocessing Data

We removed duplicate entries from our public data from 2.1 using the MinHash implementation from <https://github.com/ChenghaoMou/text-dedup>.

We configured the MinHash algorithm with the following parameters:

Parameter	Value
num_perm	256
threshold	0.95
hash_func	sha1
hash_bits	64

Complete deduplicating data implementation at [here](#). All deduped dataset published at [malaysia-ai/dedup-text-dataset](#).

After removing duplicates, we employed the postprocessing technique mentioned in Malaysian Mistral section 3.3 [2].

4 Pre-Training Tokenizer

To ensure an efficient and versatile tokenizer, we conducted pretraining on a BPE (Byte Pair Encoding) tokenizer using diverse datasets, including Malay Wikipedia, synthetic Jawi, public articles, translated code instructions, Google-translated Tamil and Google-translated Punjabi. The objective was to develop a tokenizer capable of handling longer subwords for languages such as Malay, Mandarin, Tamil, Jawi, English, and Arabic. The decision to use BPE was motivated by certain limitations observed in SentencePiece, where newline characters caused issues, some Tamil and Jawi characters were missing, and the processing speed for very long texts was suboptimal. By opting for BPE, we addressed these challenges, ensuring a robust and efficient tokenizer for our language model training.

The BPE tokenizer was trained on a deduplicated text dataset of 85GB, with a vocabulary size of 32,000.

The decision to train our own BPE tokenizer was driven by the goal of minimizing token sizes during both input and output. To illustrate the efficiency gained, we conducted a comparison using the [Malaysian Ultrachat AstroAwani dataset](#). Our pre-trained BPE tokenizer achieved a notable reduction of up to 43%

in token size when compared to tokenizers employed by Llama2 and Mistral. This reduction in token size enhances the efficiency and resource optimization of our language model, offering advantages in terms of computational performance and memory usage during processing.

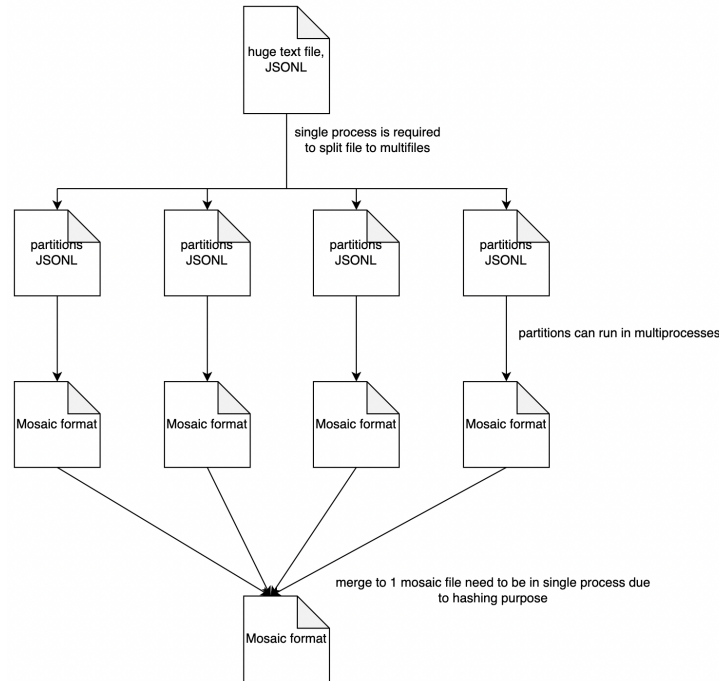
Complete pre-training tokenizer implementation at [here](#)

5 Tokenizing Data

Our dataset comprises a substantial 349GB of text in JSONL format, amounting to 90 billion tokens. The following table breaks down the distribution of token sizes within the dataset:

Distribution	Tokens (B)
deduped text Dataset	31.7
Filtered StarCoder [3]	40.98
Unfiltered MS Madlad 400 [4]	14.98
Instruction-tuned Dataset	1.58
Malaysia journals and research papers	1.14

To optimize data processing efficiency, we use MosaicML Streaming library [5], by using hashing technique. Recognizing the complexity of our dataset, we implemented data distribution processing to transform it into the MosaicML streaming format,



This involves splitting the original JSONL file into smaller ones, and each of these smaller files undergoes multiprocessing for conversion into the Mosaic format. Subsequently, these smaller Mosaic files are merged into a single Mosaic file. It's worth noting that MosaicML streaming accesses one folder at a time, necessitating the consolidation of smaller files into a unified format for seamless training data access.

References

- [1] Abhinand Balachandran. Tamil-llama: A new tamil language model based on llama 2, 2023.
- [2] Husein Zolkepli, Aisyah Razak, Kamarul Adha, and Ariff Nazhan. Large malaysian language model based on mistral for enhanced local language understanding, 2024.
- [3] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason Stiller, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder: may the source be with you!, 2023.
- [4] Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier Garcia, Christopher A. Choquette-Choo, Katherine Lee, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, and Orhan Firat. Madlad-400: A multilingual and document-level large audited dataset, 2023.
- [5] The Mosaic ML Team. streaming. <https://github.com/mosaicml/streaming/>, 2022.