

Algorithms for Question Answering to Factoid Question

Raihan Pambagyo Fadhila^{1*}, Detty Purnamasari²

¹Universitas Gunadarma, Jl. Margonda Raya No.100, Depok 16424, Indonesia

¹raihanpambagyo@gmail.com*; ²detty@staff.gunadarma.ac.id

* corresponding author

ARTICLE INFO

Article History:

Received February 6, 2025

Revised April 18, 2025

Accepted April 21, 2025

Keywords:

Natural Language Processing,

Question Answering,

BERT,

Sequence to Sequence,

GPT

Correspondence:

E-mail:

raihanpambagyo@gmail.com

ABSTRACT

The development of transformer-based natural language processing (NLP) has brought significant progress in question answering (QA) systems. This study compares three main models, namely BERT, Sequence-to-Sequence (S2S), and Generative Pretrained Transformer (GPT), in understanding and answering context-based questions using the SQuAD 2.0 dataset that has been translated into Indonesian. This research uses the SEMMA (Sample, Explore, Modify, Model, Assess) method to ensure the analysis process runs systematically and efficiently. The model was tested with exact match (EM), F1-score, and ROUGE evaluation metrics. Results show that BERT excels with an Exact Match score of 99.57%, an F1-score of 99.57%, ROUGE-1 of 97%, ROUGE-2 of 30%, and ROUGE-L of 97%, outperforming S2S and GPT models. This study proves that BERT is more effective in understanding and capturing Indonesian context in QA tasks. This research offers explanations for the implementation of Indonesian-based QA and can be a reference in the development of more accurate and efficient NLP systems.

1. Introduction

The advent of Transformer architectures has revolutionized Natural Language Processing (NLP) by introducing self-attention mechanisms that effectively capture both short- and long-range dependencies in text [1]. Efficiency-focused techniques such as Low-Rank Adaptation (LoRA) [2] and hardware-aware optimization strategies for GPUs [3] have further accelerated large-scale model training and inference. Additionally, quantization methods optimize model deployment on edge devices [4], while comprehensive surveys highlight challenges in complex knowledge-based Question Answering (QA), including domain adaptation and multi-hop reasoning [5].

Bidirectional Encoder Representations from Transformers (BERT) demonstrated state-of-the-art performance across language understanding tasks through bidirectional pre-training [6]. Generative Pre-trained Transformers (GPT-3) enabled zero- and few-shot learning paradigms in text generation [7]. Sequence-to-Sequence (S2S) models excelled in machine translation and summarization [8], and multilingual adaptations like mBART extended these capabilities across languages [9]. Additionally, T5-based approaches have shown efficacy in closed-book QA settings [10].

Among benchmarks, the Stanford Question Answering Dataset (SQuAD) 2.0 challenges models to answer extractive factoid questions while identifying unanswerable queries [11]. Complementary benchmarks such as Natural Questions emphasize real-world applicability [12], and HybridQA integrates tabular and textual data for multi-modal reasoning [13]. Despite these advances, QA research in low-resource languages like Indonesian is limited by the scarcity of large-scale, labeled corpora.

To address this gap, we translated SQuAD 2.0 into Indonesian using Opus-MT [14]. We then evaluate three Transformer-based approaches—BERT, GPT-3, and S2S—under both full fine-tuning and Parameter-Efficient Fine-Tuning (PEFT) with LoRA [2]. Experiments follow the SEMMA (Sample, Explore, Modify, Model, Assess) methodology for systematic data handling [15] and

leverage an NVIDIA DGX Station A100 for computational benchmarking [3]. Model performance is assessed using Exact Match and F1-score [11], sequence-level metrics (ROUGE-1, ROUGE-2, ROUGE-L) [5], and efficiency indicators (training time, GPU utilization, and energy consumption) [16].

In Indonesian NLP, resources such as IndoLEM and IndoBERT have enabled local model evaluation [9], while NusaX provides multilingual sentiment annotations [17]. Adaptations of BERT for Indonesian QA achieved significant improvements [18], and hybrid BERT-GPT systems show promise for dialog tasks [12]. Our work extends these contributions by offering a unified benchmark of three architectures, providing insights into accuracy-efficiency trade-offs in a low-resource context.

Our evaluation shows that BERT, when fully fine-tuned on the Indonesian SQuAD 2.0 dataset, achieves an Exact Match of 99.57%, an F1-score of 99.57%, and ROUGE scores of 97% (ROUGE-1), 30% (ROUGE-2), and 97% (ROUGE-L), significantly outperforming both S2S and GPT-3 under identical conditions.

The main contributions of this paper are: (i) A publicly available Indonesian translation of SQuAD 2.0 via Opus-MT [14]. (ii) A comprehensive comparison of BERT, GPT-3, and S2S models under full fine-tuning and PEFT regimes [2]. (iii) Empirical insights into accuracy-efficiency trade-offs, supported by QA and sequence-level metrics [11], [5], and DGX Station A100 resource usage data [3]. By uniting methodological rigor with practical evaluations, this study advances QA research for Indonesian and other low-resource languages.

2. Method

The research utilizes the SEMMA (Sample, Explore, Modify, Model, Assess) method, a systematic approach designed for efficient data analysis [15]. It involves five stages: first, the Sample stage, where a representative data sample is selected. Next is the Explore stage, which examines data patterns, relationships, and characteristics. In the Modify stage, the data is processed and adjusted to meet the analysis requirements, including cleaning or transforming variables. The Model stage follows, where techniques are applied to build a model that aligns with the research goals. Finally, in the Assess stage, the model's performance is evaluated to ensure it provides valid and relevant results. This structured approach enables a reliable and thorough analysis.

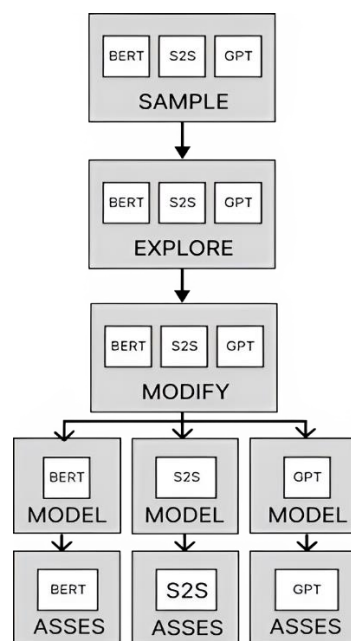


Figure 1. SEMMA Methodology

2.1. Sample

The Indonesian SQuAD 2.0 dataset [11] is split into training (80%), validation (10%), and test (10%) sets following established QA evaluation protocols [12].

2.2. Explore

Exploratory data analysis examines token distributions, answer length statistics, and translation artifacts introduced by Opus-MT [14]. Visualizations and statistical summaries leverage Transformer tokenizer utilities [6].

2.3. Modify

Preprocessing steps include text normalization (lowercasing, punctuation removal), subword tokenization (WordPiece or SentencePiece) [6], [9], and data augmentation techniques to mitigate answer class imbalance [16].

2.4. Model

Three Transformer-based models are fine-tuned:

- BERT [6], with full parameter updating.
- Sequence-to-Sequence (S2S) [8], using an encoder–decoder setup.
- GPT-3 [7], leveraging autoregressive generation.

Each model is also adapted using PEFT with LoRA [2], and we compare parameter count, memory footprint, and convergence speed..

2.5. Asses

Performance metrics include Exact Match and F1-score [11], ROUGE-1/2/L [5], and computational benchmarks (training time, GPU memory usage, energy consumption) on an NVIDIA DGX Station A100 [3]. Statistical significance is evaluated using paired t-tests [18].

3. Results and Discussion

3.1. Sample

At this stage, the collection of data is to be initiated by downloading SQuAD version 2.0 data from the official Github website. This data comprises a total of 142,000 records, which are structured in the semi-unstructured JSON data format. Following the successful download of the data, it is to be separated into two distinct sets: data intended for model training and data allocated for model testing.

3.2. Explore

This stage is characterized by an in-depth examination of the dataset, encompassing three distinct tasks: data analysis, visualization, and data quality validation. The primary objective of the data analysis stage is to comprehend the structural elements and characteristics of the data that will be employed in the natural language processing (NLP) question answering (QA) model. The SQuAD v2.0 dataset serves as the training and testing foundation for the QA model, comprising a substantial collection of question and answer pairs accompanied by a context or reference text. Each entry in the dataset under consideration contains several elements of significance.

Table 1. Data Description

Column Name	Description
title	Title or topic of the article or text used for context.
paragraphs	A collection of paragraphs containing context and related questions.
qas	A collection of questions and answers related to the context of a particular paragraph.
question	The text of the question asked based on context.
id	Unique identity for each question.
answers	A collection of answers to questions that can be answered.

Column Name	Description
title	Title or topic of the article or text used for context.
text	Text the correct answer to a question.
answer_start	Index the starting position of the answer text in the context column.
is_impossible	A binary label that indicates whether the question has answer in context or not.
context	A paragraph of text that contains relevant information to answer the question.
id	Unique identity for each paragraph.

The dataset under consideration in this study has undergone extensive statistical analysis of the primary components in the SQuAD version 2.0 dataset, namely context, question, and answer, to ascertain the characteristics of the dataset.

Table 2. Statistical Analysis Results

Elements	Mean	Median	Max
Context	106.84 Words	99 Words	653 Words
Questions	8.78 Words	8 Words	35 Words
Answers	3.46 Words	2 Words	223 Words

The visualization stage is used to see the distribution picture of the data, see the distribution of the context, questions, answers, and is_impossible labels in a bar graph. The results at this stage will display a distribution data graph.

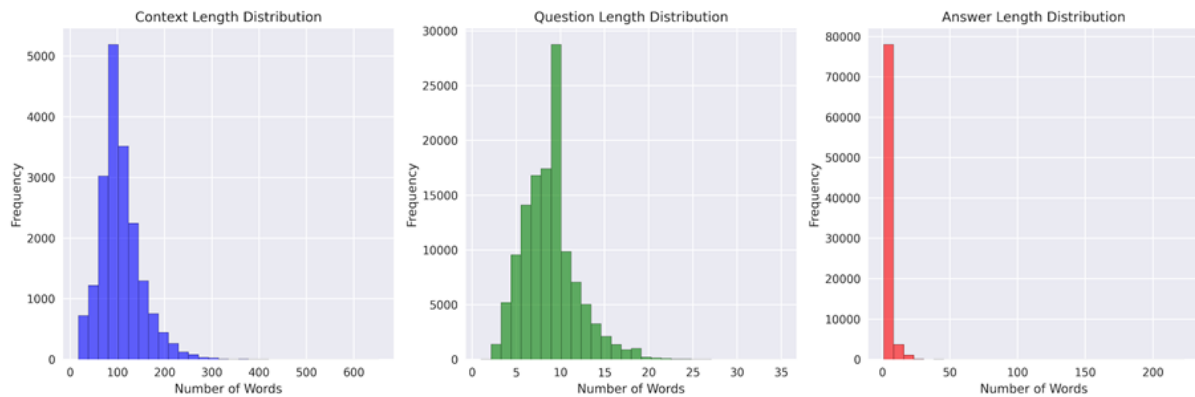


Figure 2. Data Distribution Chart

In the data quality validation stage, various checks are carried out to ensure that the data has good and consistent quality. The data quality validation check is performed by checking for the presence of empty or null values and duplicate values in the data. The output at this stage includes a summary of the total article data size, total paragraph data size, and total question data size. Additionally, it includes checking for empty or null values in the data and checking for duplicate data.

3.1. Modify

At this stage, data processing is executed to ensure its subsequent utilization. This process involves two stages: initial data construction, which involves transforming the data into a format suitable for the subsequent stage, and subsequent integration, which integrates the data with the next stage. In the data construction stage, the data is processed into a form that suits the needs of the model. Some of the steps taken include raw data extraction, tokenization, combining tokens, convert to token IDs, attention mask & padding, calculation answer position, tensor creation.

3.1.1. Raw Data Extraction

This stage displays the raw data from the JSON format file to be processed and converts it into a data structure that is easier to process for the model. This is useful in order to understand the original structure of the data, such as the question text and its context, as well as the location of the answer in context.

Table 3. Raw Data Extraction

Column Name	Description
Id	56be85543acaaa14008c9063
Question	Kapan Beyonce mulai populer?
Context	Beyoncé Giselle Knowles-Carter (/biːljɒnsel/ bee-YON-say) (lahir 4 September 1981) adalah seorang penyanyi, penulis lagu, produser rekaman dan aktris Amerika. Lahir dan dibesarkan di Houston, Texas , ia tampil di berbagai kompetisi menyanyi dan menari sebagai seorang anak, dan menjadi terkenal di akhir 1990-an sebagai penyanyi utama dari grup R&B Destiny's Child . Dikelola oleh ayahnya, Mathew Knowles , grup ini menjadi salah satu grup wanita terlaris di dunia sepanjang masa. Hiatus mereka melihat perilisan album debut Beyoncé, Dangerously in Love (2003), yang memantapkannya sebagai artis solo di seluruh dunia, memperoleh lima Grammy Awards dan menampilkan singel nomor satu Billboard Hot 100 "Crazy in Love" dan "Baby Boy".
Answer_text	di akhir 1990-an
Start_char	295
End_char	311
Is_impossible	False

3.1.2. Tokenization

In the subsequent stage of the process, the question and context data that have been extracted are processed into tokens. These tokens can then be utilized by the Transformers-based model. Each question and context is broken down into tokens using a tokenizer.

Table 4. Tokenization

Question	Question
['ka', '##pan', 'beyonce', 'mu', '##lai', 'pop', '##ule', '##r', '']	['beyonce', 'gi', '##selle', 'knowles', '-', 'carter', '(', '/', 'bi', '##:']...

3.1.3. Combining Tokens

In the combining tokens stage, the tokens from the previous process are combined to form a sequence that matches the input format in the Transformer model. After the tokenization process, the tokens from the question and the tokens from the context are combined into one token sequence. Special tokens such as [CLS] are used as start tokens to mark the beginning of input, while [SEP] tokens are used as separation tokens between question and context, and to mark the end of input.

Table 5. Combining Tokens

Combining Tokens
['[CLS]', 'ka', '##pan', 'beyonce', 'mu', '##lai', 'pop', '##ule', '##r', '?', '[SEP]', 'beyonce', 'gi', '##selle', 'knowles']...

3.1.4. Convert To Token IDs

Then the next step is to convert the merged tokens into a numerical representation using the tokenizer function. This representation is needed because the model does not understand the text directly, but processes the data in the form of numbers (IDs) that match the vocabulary.

Table 6. Convert To Token IDs

Convert To Token IDs
[101, 10556, 9739, 20773, 14163, 19771, 3769, 9307, 2099, 1029, 102, 20773, 21025, 19358, 22815]...

3.1.5. Attention Mask and Padding

Following the conversion of data into a numerical representation, the next stage is the attention mask stage. During this stage, the model is informed about the part of the input that is relevant for processing and the part that is padding. The attention mask is used to optimize the model's focus on important tokens, such as context and questions, and to ignore additional tokens, such as padding. This stage creates a list of ones to initialize the mask with the number of relevant input tokens, and adds zeros to fill the maximum required length (padding).

Table 7. Convert To Token IDs

Attention Mask (Initial mask length)	Padding (After padding length)
257	384

3.1.6. Calculation Answer Position

The primary step at this stage is to calculate the position of the start token and end token of the answer in the given context. This process is essential for language comprehension, as it requires the precise position of the token to determine the location of the answer in the context. If the data labeled is_impossible is true, then the answer position is set to zero for the start position and zero for the end position. Conversely, if the data labeled is_impossible is false, the token position of the answer will be calculated based on the character position in the context. If the start or end position exceeds the maximum length, the start and end positions will be set to zero.

Table 8. Calculation Answer Position

Start Position	End Position
114	119

3.1.7. Tensor Creation

In the tensor creation stage, the previously processed data will be converted into a tensor format. This will allow the model to use the data and optimize the computational utility of GPU processing. This stage aims to convert the processed data into a tensor format using the Pytorch library. The data elements that are converted into tensors are input_ids, attention_mask, start_position, and end_position.

Table 9. Tensor Creation

Input ids	Attention Mask	Start Position	End Position
tensor([101, 10556, 9739, 20773, 14163, 19771, 3769, 9307, 2099, 1029, 102, 20773, 21025, 19358, 22815, ...])	tensor([1, 1, 1, 1, 1, 1, 1, ...])	tensor(114)	tensor(119)

The next step is to ensure that the data processed through various previous stages can be combined and integrated structurally so that it is ready to be used in the model training process. The data integrated into a dictionary-shaped data structure are `input_ids`, `attention_mask`, `start_position`, and `end_position`.

3.2. Model

At this stage of the process, the model is being built based on the modified dataset from the previous stage. This stage is divided into two parts: the first part is building test scenarios, and the second part is building models.

3.2.1. Building Test Scenarios

The objective of constructing this test scenario is to identify the most effective model outcomes during the training process. In this study, the test scenario will be executed five times on the `learning_rate` parameter, four times on the `batch_size` parameter, and five times on the `epochs` parameter. These executions will vary several parameters believed to influence model performance. The parameters that will be altered and utilized as test scenarios are `batch_size`, `learning_rate`, and `epochs` parameters.

The learning rate is a critical factor in determining the step size of the optimization algorithm, which is used to update the model weights during the training process. To ensure the model achieves optimal performance, experiments were conducted with various values of the `learning_rate` parameter.

Table 10. Learning Rate Testing Scenario

Learning_rate	Batch Size	Epochs
1e-5	8	5
2e-5		
3e-5		
4e-5		
5e-5		

The batch size is the number of samples processed at once before the model performs parameter updates. In the test scenario, to ensure the model achieves optimal performance, tests are conducted using various batch size values.

Table 11. Batch Size Testing Scenario

Batch Size	Learning Rate	Epochs
8	Using the learning rate that has the highest result in the previous test	5
16		
32		
64		

Epochs refer to the number of full iterations of the model through the training dataset. Insufficient epochs may result in underfitting, while excessive epochs may lead to overfitting. To ascertain the optimal model in the test scenario, experiments were conducted using various epochs values.

Table 12. Epochs Testing Scenario

Epochs	Learning Rate	Batch Size
5	Using the learning rate that has the highest result in the previous test	Using the batch size that has the highest result in the previous test
10		
20		
30		
50		

3.2.2. Infrastructure Preparation

The model was trained using NVIDIA DGX Station A100 hardware.

Table 13. Infrastructure Specifications

Component	Specifications
Server	NVIDIA DGX Station A100
GPU	1x NVIDIA A100 40GB Tensor Core GPUs
CPU	AMD EPYC 7742 (128 Cores, 2.25 GHz)
RAM	512 GB DDR4
Memory	7.68 TB NVMe SSD
OS	Ubuntu 20.04 LTS

3.2.3. BERT Testing Results

Parameter testing on the BERT algorithm model using the test scenario that has been done, with a learning rate parameter of $2e-5$ and a batch size parameter of 32.

Table 14. Epochs Testing Scenario BERT

Epochs	Exact Match	F1-Score	Execution Time
5	90.67%	90.67%	198.74 Minutes
10	97.17%	97.17%	397.07 Minutes
20	98.30%	98.30%	872.95 Minutes
30	99.33%	99.33%	734.22 Minutes
50	99.57%	99.57%	1478.49 Minutes

3.2.4. Sequence to Sequence Testing Results

Parameter testing on the Sequence to Sequence algorithm model using the test scenario that has been done, with a learning rate parameter of $5e-5$ and a batch size parameter of 8.

Table 15. Epochs Testing Scenario Sequence to Sequence

Epochs	Exact Match	F1-Score	Execution Time
5	51.11%	57.00%	668.80 Minutes
10	58.40%	63.96%	788.86 Minutes
20	67.77%	73.42%	1153.50 Minutes
30	73.32%	78.62%	1520.63 Minutes
50	83.82%	88.18%	2340.76 Minutes

3.2.5. GPT Testing Results

Parameter testing on the Sequence to Sequence algorithm model using the test scenario that has been done, with a learning rate parameter of $3e-5$ and a batch size parameter of 8.

Table 16. Epochs Testing Scenario Sequence to Sequence

Epochs	Exact Match	F1-Score	Execution Time
5	93.27%	95.53%	1599.58 Minutes
10	Out of Memory	Out of Memory	Out of Memory
20	Out of Memory	Out of Memory	Out of Memory
30	Out of Memory	Out of Memory	Out of Memory
50	Out of Memory	Out of Memory	Out of Memory

3.3. Asses

At this stage, an evaluation of the model built using the SEMMA (Sample, Explore, Modify, Model, Assess) method is conducted. The primary objective is to assess the model's performance by testing it with appropriate evaluation metrics for the research. The evaluation metrics employed include exact match (EM), F1-Score, and ROUGE Score, which are designed to gauge the similarity between the model's outputs and the reference answers. At this stage, two key tasks are identified: evaluating the modeling results and reviewing the modeling process. Once the model training is complete, testing is carried out using validation and testing datasets. In the ROUGE evaluation, the model is tested using ten factoid questions.

Table 17. ROUGE Evaluation Questions

Context	No	Question	Reference
Indonesia adalah negara kepulauan yang terletak di Asia Tenggara dan Oseania, yang terdiri dari lebih dari 17.000 pulau. Negara ini berbatasan dengan Malaysia di utara, Papua Nugini di timur, dan Australia di selatan. Indonesia memiliki populasi lebih dari 270 juta jiwa, menjadikannya negara dengan populasi terbesar keempat di dunia. Jakarta adalah ibu kota dan kota terbesar di Indonesia. Bahasa Indonesia adalah bahasa resmi negara ini. Indonesia dikenal dengan keanekaragaman budaya, bahasa, dan etnis yang sangat kaya.	1.	Apa ibu kota Indonesia ?	Jakarta
	2.	Indonesia terdiri dari berapa pulau ?	lebih dari 17.000
	3.	Indonesia memiliki populasi lebih dari berapa jiwa ?	270 juta
	4.	Australia berbatasan dengan Indonesia di sebelah mana ?	selatan
	5.	Papua Nugini berbatasan dengan Indonesia di sebelah mana ?	timur
	6.	Indonesia adalah negara kepulauan yang terletak di mana ?	Asia Tenggara dan Oseania
	7.	Jakarta adalah ibu kota negara apa ?	Indonesia
	8.	Apa nama kota terbesar di Indonesia ?	Jakarta
	9.	Malaysia berbatasan dengan Indonesia di sebelah mana ?	utara
	10.	Apa bahasa resmi negara Indonesia ?	bahasa Indonesia

The following are the Model Evaluation Results using the ROUGE method.

Table 18. Model Evaluation Results

Model	Exact Match	F1-Score	ROUGE-1	ROUGE-2	ROUGE-L
BERT	99.57%	99.57%	97%	30%	97%
S2S	83.82%	88.18%	63%	25%	63%
GPT	93.27%	95.53%	60%	10%	60%

4. Conclusion

Based on the results of the research conducted, a comparative analysis of BERT, S2S, and GPT algorithms on question answering for factoid questions using the SQuAD v2.0 dataset that has been translated into Indonesian has been carried out. Based on the test results, it can be concluded that the BERT (Bidirectional Encoder Representations from Transformers) algorithm model obtained the best results in the question answering task for factoid questions using the translated SQuAD v2.0 dataset compared to the S2S (Sequence to Sequence) and GPT (Generative Pretrained Transformers) models. This is evidenced by the exact match value and F1-score which has a value of 99.57%, and has a fairly high ROUGE evaluation, with a ROUGE-1 value of 97%, ROUGE-2 value of 30%, and

ROUGE-L value of 97%. This indicates that the BERT model is more accurate in understanding the context and providing answers that match the questions given. This advantage can be attributed to the BERT architecture which is based on bidirectional attention, so that the model can capture the relationship between words in a sentence better than other models. Meanwhile, the Sequence-to-Sequence (S2S) and GPT models show relatively lower performance, especially in maintaining the accuracy of answers to factoid questions. The S2S model, which is commonly used in translation or text sequencing tasks, is less optimal in handling context-specific understanding. On the other hand, the GPT model, although superior in generating more natural text, has a weakness in maintaining answer precision due to its autoregressive nature, which leads to the possibility of generating answers that are less in line with the facts of the given context. The results of this study reinforce the understanding that the choice of model in a question answering task is highly dependent on the type of question being asked. For factoid-based questions, the BERT model proved superior to other approaches due to its ability to capture the relation between words more accurately.

References:

- [1] A. Vaswani et al., "Attention Is All You Need," 2017.
- [2] E. Hu et al., "LoRA: Low-Rank Adaptation of Large Language Models," 2021.
- [3] K.-H. Kim and C.-S. Jeong, "Optimizing Single DGX-A100 System: Overcoming GPU Limitations via Efficient Parallelism and Scheduling for Large Language Models," *Applied Sciences*, vol. 13, no. 16, p. 9306, Mar. 2023, doi: 10.3390/app13169306.
- [4] M. Nagel et al., "Quantization Techniques for Transformer Inference on Edge Devices," in *Proc. 2024 Int. Conf. Edge AI*, 2024.
- [5] Y. Lan et al., "Complex Knowledge Base Question Answering: A Survey," *IEEE Trans. Knowl. Data Eng.*, vol. 35, no. 11, pp. 11196–11215, Mar. 2023, doi: 10.1109/TKDE.2022.3223858.
- [6] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," Mar. 2018.
- [7] T. B. Brown et al., "Language Models are Few-Shot Learners," Mar. 2020.
- [8] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to Sequence Learning with Neural Networks," Mar. 2014.
- [9] F. Koto, A. Rahimi, J. H. Lau, and T. Baldwin, "IndoLEM and IndoBERT: A Benchmark Dataset and Pre-trained Language Model for Indonesian NLP," in *Proc. 28th Int. Conf. Comput. Linguist.*, 2020, pp. 757–770. doi: 10.18653/v1/2020.coling-main.66.
- [10] T. Alkhaldi, C. Chu, and S. Kurohashi, "A Peek Into the Memory of T5: Investigating the Factual Knowledge Memory in a Closed-Book QA Setting and Finding Responsible Parts," *J. Nat. Lang. Process.*, vol. 29, no. 3, pp. 762–784, 2022, doi: 10.5715/jnlp.29.762.
- [11] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, "SQuAD: 100,000+ Questions for Machine Comprehension of Text," Mar. 2016.
- [12] T. Kwiatkowski et al., "Natural Questions: A Benchmark for Question Answering Research," *Trans. Assoc. Comput. Linguist.*, vol. 7, pp. 453–466, Mar. 2019, doi: 10.1162/tacl_a_00276.
- [13] E. Qin et al., "HybridQA: A Benchmark for Multi-Modal Question Answering," in *Proc. 2021 Conf. Empirical Methods Natural Lang. Process.*, 2021, pp. 4540–4551. doi: 10.18653/v1/2021.emnlp-main.365.
- [14] P. Budzianowski and I. Vulić, "Hello, It's GPT-2 - How Can I Help You? Towards the Use of Pretrained Language Models for Task-Oriented Dialogue Systems," in *Proc. 3rd Workshop Neural Generation and Translation, ACL*, 2019, pp. 15–22. doi: 10.18653/v1/D19-5602.
- [15] SAS Institute, "SEMMA: Data Mining Methodology," SAS Technical Report, 2001.
- [16] S. Ruder et al., "Transfer Learning in Natural Language Processing," in *Proc. NAACL-HLT*, 2019, pp. 15–18. doi: 10.18653/v1/N19-5004.
- [17] G. I. Winata et al., "NusaX: Multilingual Parallel Sentiment Dataset for 10 Indonesian Local Languages," Mar. 2022.
- [18] K. Duan et al., "Enhancement of Question Answering System Accuracy via Transfer Learning and BERT," *Applied Sciences*, vol. 12, no. 22, p. 11522, Mar. 2022, doi: 10.3390/app122211522.