
LARGE LANGUAGE MODELS FOR FAILURE MODE CLASSIFICATION: AN INVESTIGATION

A PREPRINT

Michael Stewart

Department of Computer Science and Software Engineering
The University of Western Australia
Perth, Western Australia
michael.stewart@uwa.edu.au

Melinda Hodkiewicz

School of Engineering
The University of Western Australia
Perth, Western Australia
melinda.hodkiewicz@uwa.edu.au

Sirui Li

Department of Computer Science and Software Engineering
The University of Western Australia
Perth, Western Australia
sirui.li@uwa.edu.au

September 18, 2023

ABSTRACT

In this paper we present the first investigation into the effectiveness of Large Language Models (LLMs) for Failure Mode Classification (FMC). FMC, the task of automatically labelling an observation with a corresponding failure mode code, is a critical task in the maintenance domain as it reduces the need for reliability engineers to spend their time manually analysing work orders. We detail our approach to prompt engineering to enable an LLM to predict the failure mode of a given observation using a restricted code list. We demonstrate that the performance of a GPT-3.5 model ($F1=0.80$) fine-tuned on annotated data is a significant improvement over a currently available text classification model ($F1=0.60$) trained on the same annotated data set. The fine-tuned model also outperforms the out-of-the box GPT-3.5 ($F1=0.46$). This investigation reinforces the need for high quality fine-tuning data sets for domain-specific tasks using LLMs.

Keywords Technical Language Processing · Failure Mode · Large Language Models · Maintenance

1 Introduction

The maintenance of assets plays a critical role in the safety and costs of industrial organisations. One of the key tasks within maintenance is failure mode identification. This task is done by reliability engineers to capture and code failure and other undesirable events. These failure mode codes, together with data such as the cost/ production/ service impact, safety and environmental consequence of the event are used to prioritise improvement work, update maintenance strategy and can assist product/ plant engineers to improve future design by updating their failure modes and effects analysis. Consistent and reproducible failure mode code assignment is difficult as the observation of each event are captured by field technicians in natural language. For example, consider the following maintenance work order texts:

- pump runs for a while and trip
- engin does not work
- pmp spraying out slurry
- seal leaking
- leak in seal

Each of these work orders contain an observation made by the field technician, such as “does not work”, “leaking”, and so on. In any maintenance management system there are thousands of these observations and each needs a failure mode classification (FMC), such as “leaking” and “breakdown” according to an agreed list. The challenge, is that each person doing the coding, whether it be the technician generating the work order, or the reliability engineer reviewing it, comes with their own mental model of the asset and its behaviour [Sexton et al., 2019]. Further, attention to the task of coding accurately is influenced by factors such as training, managerial support, technological input control and motivation [Murphy, 2009, Unsworth et al., 2011, Molina et al., 2013]. It is too expensive to have university-trained reliability engineers review each of these codes manually given the volume. The opportunity for AI to assist in failure mode classification is therefore an active research area [Sexton et al., 2018, Akhbardeh et al., 2020, Sala et al., 2022, Stewart et al., 2022, Usuga-Cadavid et al., 2022].

There has recently been a surge of interest in Large Language Models (LLMs), predominately as the result of the popularity of chatbot interfaces such as ChatGPT¹. LLMs such as OpenAI’s GPT-3.5² have been trained on massive corpora and thus encapsulate knowledge from a wide variety of domains. It has also been shown that LLMs require little to no fine-tuning, meaning they exhibit excellent performance with barely any annotated training data [Brown et al., 2020]. Rather than focusing on developing manually-annotated datasets to train models (like with more “traditional” text classification models such as Flair [Akbik et al., 2018]), users of LLMs typically employ *prompt engineering* in order to craft their input prompt to elicit a particular response from the model. As a result of their excellent performance on a wide range of natural language processing tasks, LLMs have already been applied to a variety of domains. Examples include medicine [Singhal et al., 2022, Thirunavukarasu et al., 2023], education [Kasneci et al., 2023], and vehicle accident records [Mumtarin et al., 2023].

However, to the best of our knowledge, no research has yet investigated the use of LLMs within the maintenance domain, let alone specifically for FMC. In light of this research gap, and the potential for automated FMC to enable significant time and cost benefits to industry, we present an investigation into the effectiveness of using Large Language Models for Failure Mode Classification. Our contributions are as follows:

- We investigate the most effective prompt format for performing FMC using an LLM without any fine-tuning.
- We determine whether it is necessary to fine-tune an LLM on a set of annotated data to achieve good FMC performance.
- We provide a comparison between the performance of fine-tuned LLMs and text classification models for FMC.

This paper is structured as follows. We begin by providing an outline of our models, methods and experiments, and detail the dataset that we use for fine-tuning and evaluation. We then present our results, which directly tie in to our contributions above. Finally, we present our conclusion and an outlook to future work.

The source code of this paper is open source and is available on GitHub.

2 Methods

The aim of this paper is to evaluate the applicability of Large Language Models (LLMs) to Failure Mode Classification (FMC). In this section we provide an overview of the dataset we are using for our evaluation, as well as the models that we evaluate in Section 3.

2.1 Dataset

The dataset on which we evaluate each model is an extract from the annotated maintenance work order dataset introduced by [Stewart et al., 2022] and available on PapersWithCode³. The data set consists of 502 (observation, label) pairs for training, 62 for validation, and 62 for testing. The observations, which are written in natural language, were extracted from a set of maintenance work orders using Named Entity Recognition (NER). The labels are taken from a set of 22 failure mode codes from ISO 14224⁴. Each observation was labelled by a domain expert. Some examples from this dataset are as follows:

- broken, Breakdown

¹<https://chat.openai.com/>

²<https://platform.openai.com/docs/models>

³<https://paperswithcode.com/dataset/fmc-mwo2kg>

⁴<https://www.iso.org/standard/64076.html>

- leaking fluid, Leaking
- too hot, Overheating
- tripping, Electrical
- not starting, Failure to start on demand

This open data set and the model presented in [Stewart et al., 2022] represent the state-of-the-art for FMC in the literature at this point in time and hence are used for comparative purposes.

2.2 Models

We evaluate the following models:

1. **Flair**: A Flair-based [Akbik et al., 2018] text classification model, trained on the annotated dataset.
2. **GPT-3.5**: The off-the-shelf GPT-3.5-Turbo model from OpenAI.
3. **GPT-3.5 (Fine-tuned)**: The GPT-3.5-Turbo model, fine-tuned on the annotated dataset.

The Flair model is a Bidirectional Long Short-Term Memory-based [Hochreiter and Schmidhuber, 1997] text classification model that takes a sequence of text as input, and predicts a single label. This is the same model as used in [Stewart et al., 2022], and further implementation details are available in the respective paper.

The first layer of the model, the embedding layer, was pre-trained by the Flair developers on a corpora of web, Wikipedia data, and subtitles, and thus the model has little innate knowledge of maintenance. The model was trained by [Stewart et al., 2022] on the dataset of 502 (observation, label) pairs and validated on the 62-pair validation set. In contrast to the GPT-based models, the computational requirements of training and using this model are low enough to be able to train on most desktop computers. This also means it can be used offline, and is thus appropriate for handling sensitive data.

The LLM-based models are based on OpenAI's GPT-3.5 [Brown et al., 2020]⁵, the model behind ChatGPT⁶. The GPT-3.5 model is “off-the-shelf” in that we are using the model without any form of fine-tuning. We are relying on the model’s knowledge of maintenance that it has gleaned from its massive training corpora in order to task it to perform failure mode classification. The GPT-3.5 (Fine-tuned) model, on the other hand, is fine-tuned on the annotated dataset of 502 (observation, label) pairs, and validated on the 62-pair validation set.

2.3 Data preparation

```
[{
    "role": "system",
    "content": "Determine the failure mode of the observation
               provided by the user."
},
{
    "role": "user",
    "content": "too hot"
},
{
    "role": "assistant",
    "content": "Overheating"
}]
```

Listing 1: An example prompt that is fed into the GPT-3.5 and GPT-3.5 (Fine-tuned) models. The role of the assistant is only used during fine-tuning.

The default behaviour of the GPT-based models is to act as a chatbot, and thus it will not respond with a failure mode code for a given observation unless the instruction to do so is included as part of the prompt. Structuring an input prompt to elicit a particular response from a large language model is known as *prompt engineering*.

The latest versions of the GPT-based models require a three-part prompt. The system-level prompt dictates the desired response format of the model. For example, one can use this prompt to ask the model to reply in a sarcastic tone, or to reply with a one-word answer, and so on. The user-level prompt is the input from the user. Finally, the

⁵GPT-4.0 was not available for fine-tuning as of the time of writing, hence the decision to use GPT-3.5.

⁶<https://chat.openai.com/>

Observation	LLM output
runs for a while and trip	The failure mode of the observation "runs for a while and trips" suggests an electrical failure. This could be due to an overload, short circuit, or other electrical issue that causes the equipment to shut down or trip a circuit breaker.
does not work	The failure mode of the observation "does not work" is a breakdown.
spilling	The failure mode of the observation "spillage" is a leakage.
spraying out slurry	The failure mode of the observation "spraying out slurry" is leaking.

Table 1: Some examples of predictions made by the off-the-shelf GPT-3.5-Turbo on a sample of the test data. The system-level prompt is “Determine the failure mode of the observation provided by the user.”

assistant-level prompt is the desired input from the LLM (this is used when fine-tuning to inform the model of the expected output).

To create the prompts, we wrote Python code to iterate through the annotated CSV-based dataset and convert each (observation, label) pair into a prompt as shown in Listing 1. The same system-level prompt is used for each input to the model, and describes the task to perform (failure mode classification). We use the user-level prompt to provide the model with the observation that we want it to label. During the fine-tuning of the GPT-3.5 (Fine-tuned), we include an assistant-level prompt that informs the model of the desired output for each observation (i.e. the failure mode). The design behind these prompts were based on the best practices listed in the OpenAI Documentation⁷.

In our experiments we also investigate the necessity to add the following two texts to the system-level prompt:

- In Section 3.1, we include the sentence “Your answer should contain only the failure mode and nothing else.” to instruct the language model to avoid outputting unnecessary text (e.g. “The failure mode is ...”, etc.
- In Section 3.2 we include “Valid failure modes are: ” followed by a newline-separated list of valid labels from the dataset. This is an attempt to ensure that the model does not come up with its own failure modes, but instead outputs a failure mode code from the prescribed list.

2.4 Evaluation metrics

In the same manner as [Stewart et al., 2022], we evaluate each model using Micro F1 and Macro F1 score. Micro F1 calculates an F1-Score by adding the true positives (TPs), false positives (FPs) and false negatives (FNs) from all class labels together and then calculating F1-Score:

$$\text{MicroF1} = \text{F1}_{(class1+class2+\dots+class_n)} \quad (1)$$

Macro f1, on the other hand, simply averages the F1-Score of each class. Given N is the number of class labels, it is calculated as follows:

$$\text{MacroF1} = \frac{\sum_{n \in N} \text{F1}_{class_n}}{N} \quad (2)$$

3 Results

This section aims to answer the following questions:

1. How best to use an off-the-shelf Large Language Model (LLM) to perform Failure Mode Classification (FMC)?
2. Is it necessary to fine-tune the LLM to perform FMC?
3. Are LLMs more effective at FMC than text classification models?
4. What are some barriers one may face when using LLMs for FMC?

3.1 How best to use an off-the-shelf LLM to perform Failure Mode Classification?

To address the first research question we begin by investigating the use of a simple system-level prompt of “Determine the failure mode of the observation provided by the user.”. Upon feeding this prompt into the model, along with

⁷<https://platform.openai.com/docs/guides/fine-tuning/fine-tuning-examples>

Observation	LLM output
runs for a while and trip	Overheating
very stiff to operate	Stiff operation
requires rebuild	Noisy operation
has no equipment earth	N/A
high earth reading	No failure mode can be determined from the given observation.
failed electrical	Failure mode: Electrical failure

Table 2: Some examples of predictions made by the off-the-shelf GPT-3.5-Turbo on a sample of the test data. The system-level prompt is “Determine the failure mode of the observation provided by the user. Your answer should contain only the failure mode and nothing else.”

the user-level prompt (the observation, e.g. “runs for a while and trip”), the LLM produces outputs as shown in Table 1. These outputs, which are conversational in nature, are not machine-readable and are therefore not applicable to downstream analysis. A more specific prompt is needed to perform FMC.

In light of this, we next add the phrase “Your answer should contain only the failure mode and nothing else.” to the system-level prompt. Adding this sentence to the prompt results in the model predicting a single failure mode for each observation, as shown in Table 2. However, there are several notable issues with the outputs of the model after adding this phrase. Firstly, despite the addition of the phrase in the prompt, the model still occasionally adds additional text to its response. One such example is its response for the phrase “failed electrical”, to which it also adds “Failure mode:” prior to the actual classification. It also occasionally disregards the instruction when it was not capable of recognising a particular failure mode, for example in its classification of “high earth reading”.

While the LLM is capable of predicting failure modes using this prompt, they are not aligned with any particular failure mode ontology. Downstream analysis using these failure modes is thus not possible, due to the sheer number of possible failure modes and inconsistency between them. For example, the model predicts both “Leakage” and “Leaking”, which are the same failure mode written two different ways. One can liken the LLM’s predicted failure modes to that which might be produced by a layperson, i.e. not a domain expert.

The non fine-tuned model also has difficulties producing consistent failure mode labels when dealing with uncertainty. When the model is unable to classify the observation, it responds in a variety of different ways, for example “Insufficient information”, “N/A”, “None”, “No failure mode detected.”, “No failure mode provided.”, and so on. Attempting to resolve all possible variations of these phrases into a single classification (such as “Unknown” or “Other”) is a non-trivial task, and thus the outputs of this model are not readily applicable to downstream tasks.

In an attempt to solve this issue we add a final phrase to the prompt: “Valid failure modes include: ” followed by a newline-separated list of the failure mode labels appearing across the entire dataset. We found that this addition generally causes the model to behave as expected. However, it occasionally hallucinates labels: for example, it predicts the label “Fail to open” for “sticking shu”, and “Fail to adjust” for “cant be adjusted”. It also has issues with label consistency - for example, it predicts both “Fail to function” and “Failure to function”. Similarly to the previous attempt without constraining the label space, this attempt at using the LLM directly without fine-tuning is not directly applicable to failure mode analysis as a result of these issues.

In summary we have demonstrated that it is possible to engineer the prompt to enable the LLM to predict failure mode codes without any fine-tuning. However, these outputs are not grounded in any particular ontology and are inconsistent.

3.2 Is it necessary to fine-tune the LLM to perform Failure Mode Classification?

We now aim to determine whether fine-tuning the LLM on a purpose-built dataset is necessary, or whether similar performance can be achieved without fine-tuning. We focus our attention on a comparison between the GPT-3.5 model, and GPT-3.5 (Fine-tuned). The former model has been fed with the prompt discussed at the end of 3.1, i.e. it constrains the model to predict only the failure mode and nothing else, and also provides it with a list of the valid failure modes from the dataset. The latter model has been fine-tuned on the 500 (observation, label) pairs in the training dataset, and the prompt does not contain the aforementioned constraints (as they are not necessary due to the fine-tuning).

Table 3 shows the results of each model on the test dataset. It is clear that fine-tuning has a significant impact on performance, as the Micro-F1 score jumps from 0.46 to 0.81 between the non fine-tuned and fine-tuned models respectively. The results of the non fine-tuned model indicate that it does possess knowledge of maintenance, though, as it was capable of getting nearly half of all predictions correct without any form of fine-tuning.

Failure mode	Support	Flair	GPT-3.5	GPT-3.5 (FT)
Abnormal instrument reading	1	1.00	1.00	0.00
Breakdown	7	0.37	0.44	1.00
Contamination	1	1.00	1.00	1.00
Electrical	6	0.67	0.50	0.67
Erratic output	1	0.00	0.00	0.00
Fail to function	3	0.50	0.00	0.00
Failure to start on demand	1	0.40	0.33	1.00
Failure to stop on demand	1	0.00	1.00	1.00
High output	1	0.00	1.00	1.00
Leaking	3	0.67	0.86	1.00
Low output	2	0.00	0.00	0.00
Minor in-service problems	17	0.73	0.11	1.00
Other	2	0.67	0.40	0.00
Overheating	4	1.00	1.00	1.00
Plugged / choked	6	0.67	0.25	1.00
Spurious stop	1	0.00	0.00	0.00
Structural deficiency	3	0.60	0.57	1.00
Vibration	2	0.67	1.00	1.00
Micro-F1		0.60	0.46	0.81
Macro-F1		0.46	0.53	0.62

Table 3: A comparison of the Flair model [Stewart et al., 2022] and the GPT-3.5 LLMs (non-fine-tuned and fine-tuned) on the test dataset. Support is the number of times the label appears in the test dataset. The results of the top-performing model (when there are no ties) are in **bold**.

We also tested the effectiveness of “few-shot learning”, i.e. providing a list of example (observation: failure mode) pairs to the model as part of the system-level prompt as opposed to a list of only the valid failure modes. We found that the results were near identical to the non fine-tuned model, and thus did not include these results in the table for brevity. Overall, the results show that fine-tuning is necessary to achieve strong performance. This demonstrates the importance of high quality annotated data when applying LLMs to maintenance work orders.

3.3 Are LLMs more effective at failure mode classification than text classification models?

To answer this final research question we focus our attention to a comparison between the Flair text classification model from [Stewart et al., 2022] and the GPT-3.5 models. As shown in Table 3, the LLM significantly outperforms Flair, but only after fine-tuning. Without fine-tuning, Flair exhibits much stronger performance, indicating the necessity of annotated training data to be able to perform this particular task.

After fine-tuning on the annotated data, the LLM performs significantly better than Flair. It also tends to fair better on the minority classes, such as “Failure to start on demand”, “Failure to stop on demand”, etc, which we argue can be attributed to the underlying knowledge made available as part of the LLM’s lengthy training process on a large corpora.

In summary, our results show this LLM is more effective at FMC than the text classification model, but only when the LLM is fine-tuned to perform this task.

3.4 What are some barriers one may face when using LLMs for FMC?

Overall we found the process of using and fine-tuning GPT-3.5 fairly straightforward, though we experienced a couple of issues that are worth noting. Firstly, the non-deterministic nature of LLMs mean that they can produce different output given the same input. There is a built-in temperature parameter which can be set to 0 to reduce the likelihood of this occurring, but in our experience we were still receiving slightly different results each time we ran our experiments. This effect is most noticeable in the non fine-tuned model with no prompt engineering (i.e. from Section 3.1, and has less of an effect when the model is informed of the list of valid labels).

We also noticed that during inference, the OpenAI API would occasionally refuse our requests due to being overloaded, causing us to have to start the inference process again. This was not a significant problem for our small 62-record test set, but it would be more problematic when running inference over a large dataset.

Finally, we note that the overall fine-tuning and inference process was fairly inexpensive, costing approximately \$1 USD for each of our experiments. This shows that cost is not a barrier for achieving an acceptable level of performance on failure mode classification using LLMs.

4 Conclusion

In this paper we have demonstrated the use of Large Language Models (LLMs) to perform Failure Mode Classification (FMC). We have investigated the use of prompt engineering to determine the best prompt to feed in to an LLM, such as GPT-3.5, in order to perform FMC without any fine-tuning. However, we have also found that fine-tuning an LLM is necessary to obtain significantly better performance on FMC when compared to text classification models such as Flair. The fine tuning is performed using a relatively small, high quality, annotated data set.

The annotated data set we used for fine-tuning is publicly available. It maps observations to failure modes based on ISO 14224 classes. For the benefit of industry users wishing to use this fine-tuned data set on their own data, we note they will need to preprocess their maintenance work orders to extract observations. An example of a code pipeline to do this is in [Stewart et al., 2022].

One of the key drawbacks of OpenAI’s LLMs is that to be able to fine-tune the models, one must upload potentially sensitive data to OpenAI’s servers. This is a non-issue for companies with the capability to run and fine tune LLMs in their own secure environments, but presents complications for others. In light of this, in the future we aim to investigate the performance of offline large language models, such as LLaMA [Touvron et al., 2023], on failure mode classification. We also plan to explore how well the Flair-based model performs on this task when it is fed with GPT-based embeddings. Finally, we also plan to release a larger annotated dataset than the one proposed by [Stewart et al., 2022], which will enable further fine-tuning and improved evaluation quality.

Acknowledgments

This research is supported by the Australian Research Council through the Centre for Transforming Maintenance through Data Science (grant number IC180100030), funded by the Australian Government.

References

- Thurston Sexton, Melinda Hodkiewicz, Michael P Brundage, et al. Categorization errors for data entry in maintenance work-orders. In *Proceedings of the Annual Conference of the PHM Society*, volume 11, 2019.
- Glen D Murphy. Improving the quality of manually acquired data: Applying the theory of planned behaviour to data quality. *Reliability Engineering & System Safety*, 94(12):1881–1886, 2009.
- Kerrie Unsworth, Elisa Adriasola, Amber Johnston-Billings, Alina Dmitrieva, and Melinda Hodkiewicz. Goal hierarchy: Improving asset data quality by improving motivation. *Reliability Engineering & System Safety*, 96(11):1474–1481, 2011.
- Roger Molina, Kerrie Unsworth, Melinda Hodkiewicz, and Elisa Adriasola. Are managerial pressure, technological control and intrinsic motivation effective in improving data quality? *Reliability Engineering & System Safety*, 119: 26–34, 2013.
- Thurston Sexton, Melinda Hodkiewicz, Michael P Brundage, and Thomas Smoker. Benchmarking for keyword extraction methodologies in maintenance work orders. In *Proceedings of the Annual Conference of the PHM Society*, volume 10, 2018.
- Farhad Akbardeh, Travis Desell, and Marcos Zampieri. NLP tools for predictive maintenance records in maintnet. In *Proceedings of the 1st conference of the Asia-Pacific chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing: System demonstrations*, pages 26–32, 2020.
- Roberto Sala, Fabiana Pirola, Giuditta Pezzotta, and Sergio Cavalieri. NLP-based insights discovery for industrial asset and service improvement: an analysis of maintenance reports. *IFAC-PapersOnLine*, 55(2):522–527, 2022.
- Michael Stewart, Melinda Hodkiewicz, Wei Liu, and Tim French. Mwo2kg and echidna: Constructing and exploring knowledge graphs from maintenance data. *Proceedings of the Institution of Mechanical Engineers, Part O: Journal of Risk and Reliability*, page 1748006X221131128, 2022.
- Juan Pablo Usuga-Cadavid, Samir Lamouri, Bernard Grabot, and Arnaud Fortin. Using deep learning to value free-form text data for predictive maintenance. *International Journal of Production Research*, 60(14):4548–4575, 2022.

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33:1877–1901, 2020.
- Alan Akbik, Duncan Blythe, and Roland Vollgraf. Contextual string embeddings for sequence labeling. In *COLING 2018, 27th International Conference on Computational Linguistics*, pages 1638–1649, 2018.
- Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfahl, et al. Large language models encode clinical knowledge. *arXiv preprint arXiv:2212.13138*, 2022.
- Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature Medicine*, pages 1–11, 2023.
- Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günemann, Eyke Hüllermeier, et al. Chatgpt for good? on opportunities and challenges of large language models for education. *Learning and Individual Differences*, 103:102274, 2023.
- Maroa Mumtarin, Md Samiullah Chowdhury, and Jonathan Wood. Large language models in analyzing crash narratives—a comparative study of chatGPT, BARD and GPT-4. *arXiv preprint arXiv:2308.13563*, 2023.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. LLAMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.