

Large language model number handling in the Finance Domain

Anant Raj



Master of Science
School of Informatics
University of Edinburgh
2024

Abstract

This research project focuses on enhancing the performance of large language models (LLMs) in finance by improving their capabilities in numerical data handling and financial reasoning. By employing a mixture of financial datasets our approach encompasses specialized instruction tuning methods alongside the strategic use of specialised external computational tools for specific tasks. The project rigorously fine-tuned open-source LLMs like Llama2,Llama3 and evaluate their performance in comparison with commercial models such as GPT-4. The principal aim is to forge a robust framework for the processing and analysis of financial data, which minimizes human intervention while enhancing accuracy and efficiency, thus redefining the benchmark for AI applications in financial analysis.

Research Ethics Approval

This project was planned in accordance with the Informatics Research Ethics policy. It did not involve any aspects that required approval from the Informatics Research Ethics committee.

Declaration

I declare that this thesis was composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification except as specified.

(Anant Raj)

Acknowledgements

I am profoundly grateful to my supervisor, Alexandra Birch, whose expertise, understanding, and patience added substantially to my graduate experience. I deeply appreciate the countless hours she spent discussing experimental results with me and her insightful comments that continually inspired me throughout this journey. Her comprehensive guidance and unwavering support were instrumental in the successful completion of this project.

I extend my heartfelt thanks to Mateusz, whose encouragement was a beacon of light during times of struggle. His motivating words helped me navigate through the challenging moments of my research.

My sincere gratitude also goes out to my friends and family, whose unending support and love provided me with strength and encouragement. Their belief in my abilities sustained me throughout this academic endeavor.

This dissertation stands as a testament not only to my efforts but also to the valuable contributions of all those who supported me along the way. Thank you.

Table of Contents

1	Introduction	1
2	Background	4
2.1	The Transformer Architecture: A Paradigm Shift	4
2.2	Evolution of Large Language Models	5
2.3	Advancements in LLM Capabilities	6
3	Related Work	8
4	Methodology	11
4.1	Dataset Selection and Preparation	11
4.1.1	Merged Dataset Creation	11
4.1.2	Dataset Preprocessing	12
4.2	Baseline Model Selection and Evaluation	12
4.3	Fine-tuning Approaches	14
4.3.1	Training Procedure	14
4.3.2	Post Processing Algorithm	18
4.4	Evaluation and Iteration	18
4.4.1	Permissive Accuracy	18
4.4.2	Exact Match Accuracy	20
4.4.3	Inference Phase Metrics	22
4.5	External Tool Integration & Few-Shot Code Generation	22
5	Experiments	27
5.1	Research Questions	27
5.2	Experimental Setup	28
5.2.1	Dataset	28
5.2.2	Models	28

5.2.3	Evaluation Metrics	29
5.3	Results and Analysis	29
5.3.1	Model Architecture Comparison	29
5.3.2	Impact of Fine-tuning	30
5.3.3	Few-shot Prompting Strategies	30
5.3.4	Optimal Few-shot Prompting	30
5.3.5	Example Type Distributions	31
5.3.6	Tool-based Approach	31
5.4	Code Generation Approach	32
5.4.1	Experimental Setup	32
5.4.2	Results and Analysis	33
5.5	Discussion	35
6	Conclusions	37
6.1	Limitations	37
6.2	Future Work	37
	Bibliography	38
A	First appendix	42
A.1	Prompt template	42
A.2	Examples	46
A.2.1	TAT-QA Dataset	46
A.2.2	CONVFINQA Dataset	46

Chapter 1

Introduction

The intersection of Large Language Models (LLMs) and the financial industry represents a frontier of considerable promise and complexity. As documented by recent advancements, LLMs such as ChatGPT and GPT-4 have showcased remarkable proficiencies across a spectrum of Natural Language Processing (NLP) tasks, driven by sophisticated training methodologies including reinforcement learning from human feedback (RLHF) and masked language model objectives [16]. However, recent advancements have highlighted their limitations, such as the inability to access current information, difficulties with arithmetic and mathematical computations, and a propensity to generate inaccurate responses due to hallucinations [23, 14]. Despite being trained on diverse datasets covering multiple genres and subjects, their effectiveness in specialized domains like finance still necessitates further scrutiny.

In the financial arena, LLMs have begun to play an indispensable role, aiding in tasks such as investment sentiment analysis, financial named entity recognition, question-answering systems, and stock market prediction, which assist financial analysts in navigating complex datasets and predictive models. Methods such as multi-field LLMs and instruction fine-tuned LLMs have been explored to enhance the efficacy of generated results in the financial context [12]. Nonetheless, the nuanced understanding and processing of financial data by LLMs — critical for applications ranging from market analysis to investment strategy formulation — remain at an early stage for empirical study and innovation.

This project aims to investigate various number-handling strategies to fine-tune a large language model using instruction data, with the goal of achieving optimal performance in finance-related downstream tasks, particularly those involving currency or numerical questions. A diverse set of financial datasets has been selected, encompassing

tabular data, contextual information, relation extraction, and Conversational Finance Question Answering. These complex numerical reasoning datasets provide a robust foundation for generalizing across a wide range of financial tasks.

Initially, baseline performance was assessed on open-source models such as Llama2 and the newly launched Llama3 using these datasets. It is well known that large language models (LLMs) often struggle with numerical data, particularly in tasks requiring numerical reasoning, and our results were consistent with this observation. Several fine-tuning approaches were explored, starting with the Llama2-chat 7B model and subsequently extending experiments to the latest Llama3 8B model. Recognizing that a single prompting technique does not suffice for all scenarios, various techniques were experimented with, yielding notable results detailed in the experiment section.

In our study, we venture beyond the traditional confines of language model applications, drawing upon the groundbreaking work of Schick et al. (2023) who utilized external tools to elevate the performance of models handling numerical values over those processing textual data. This approach aligns with cutting-edge developments like LangChain agents and function calling, which empower models to leverage external capabilities effectively.

Our research work explores the potential of LLMs in processing and analyzing intricate financial datasets, addressing the current limitations in computational power and the scarcity of open-source models capable of handling extensive contextual inputs typical of financial data. Our study operates under the hypothesis that, with optimal resource allocation and fine-tuning, LLMs can significantly enhance their script generation capabilities, potentially surpassing existing benchmarks in numerical accuracy and analytical depth. To investigate this hypothesis, we conducted a two-phase pilot study. Initially, we utilized a compact, curated dataset sample with the FACEBOOK/OPT-1.3B model, integrated with LangChain's Python agent (*python_repl*), which yielded promising results. Building on this success, we scaled up our experiment, employing an LLM to generate Python scripts for more extensive financial datasets, executed through Python's *exec*¹ function in a pipeline-like structure. While this expanded experiment did not achieve the same level of accuracy as the smaller, handcrafted sample, it produced intriguing results that highlight the critical role of model size and sophistication in code generation based on complex inputs.

We posit that leveraging larger, more advanced language models with proper fine tuning and prompting could potentially bridge this performance gap. This research

¹<https://docs.python.org/3/library/functions.html#exec>

not only contributes to the growing body of knowledge at the intersection of natural language processing and financial analytics but also raises compelling questions about the scalability of LLMs in financial analysis, the efficacy of fine-tuning strategies, and the optimal balance between model size, computational efficiency, and analytical accuracy. As we navigate these challenges, our study lays a robust foundation for future investigations, potentially paving the way for significant advancements in how we approach and execute complex financial analyses using open source LLMs.

Chapter 2

Background

The field of natural language processing (NLP) has experienced a profound transformation in recent years, primarily propelled by the introduction of transformer-based models and the subsequent development of large language models (LLMs). This chapter presents an overview of the evolution of these models. Section 2.1 introduces the architecture of transformers, detailing their fundamental components and operational principles. Section 2.2 delineates the progression of Large Language Models, outlining key milestones and innovations. Section 2.3 examines the advancements in LLM capabilities, highlighting their impact on the field of NLP.

2.1 The Transformer Architecture: A Paradigm Shift

The introduction of the Transformer architecture by Vaswani et al. (2017) marked a pivotal moment in NLP, fundamentally changing the approach to sequence transduction problems. Unlike previous recurrent or convolutional neural networks, the Transformer relies solely on attention mechanisms, enabling more efficient processing of sequential data and facilitating the training of much larger models [26].

Key components of the Transformer architecture include:

- **Self-Attention Mechanism** This allows the model to weigh the importance of different parts of the input sequence when processing each element, enabling capture of long-range dependencies more effectively than previous architectures.
- **Multi-Head Attention** By applying multiple attention operations in parallel, the model can capture different types of relationships within the data simultaneously.

- **Positional Encoding:** To compensate for the lack of inherent sequential processing, positional encodings are added to input embeddings, allowing the model to leverage sequence order.
- **Feed-Forward Networks** These process the outputs of the attention layers, adding non-linearity and increasing the model's capacity to learn complex functions.
- **Layer Normalization and Residual Connections** These components facilitate training of deep networks by stabilizing the learning process and mitigating the vanishing gradient problem.

The Transformer's ability to process input sequences in parallel, rather than sequentially, led to significant improvements in training efficiency and model performance. This architecture laid the groundwork for the development of increasingly large and sophisticated language models.

2.2 Evolution of Large Language Models

Building upon the Transformer architecture, researchers developed a series of increasingly powerful language models:

BERT: Bidirectional Encoders

BERT (Bidirectional Encoder Representations from Transformers), introduced by Devlin et al. (2018), represented a significant advancement in pre-training techniques [7].

Key innovations of BERT include:

- **Bidirectional Context:** Unlike previous models that processed text either left-to-right or right-to-left, BERT considers both left and right context simultaneously, enabling a more nuanced understanding of language.
- **Masked Language Model (MLM) Pre-training:** BERT is trained to predict masked words in a sentence, forcing it to learn contextual representations of words.
- **Next Sentence Prediction:** This additional pre-training task helps BERT understand relationships between sentences.

BERT's approach led to state-of-the-art performance across a wide range of NLP tasks, demonstrating the power of large-scale, unsupervised pre-training.

GPT Series: Scaling Up

The GPT (Generative Pre-trained Transformer) series, developed by OpenAI, pushed the boundaries of model size and capabilities:

- GPT : Introduced the concept of fine-tuning a large, pre-trained language model for specific downstream tasks [20].
- GPT-2 : Scaled up the model size significantly (1.5 billion parameters) and demonstrated impressive text generation capabilities [21].
- GPT-3 : With 175 billion parameters, GPT-3 marked a leap in scale and showcased strong few-shot learning abilities, reducing the need for task-specific fine-tuning [3].

The GPT series highlighted the benefits of scale in language models, showing that larger models could exhibit more flexible and generalizable language understanding and generation capabilities.

Other Notable Models

- T5 : Unified various NLP tasks into a single text-to-text format, demonstrating the versatility of the transformer architecture [22].
- BART : Combined the bidirectional encoder of BERT with the autoregressive decoder of GPT, showing strong performance on both text comprehension and generation tasks [13].

These models further expanded the capabilities of LLMs, showing their potential to handle a diverse range of NLP tasks within a single architectural framework.

2.3 Advancements in LLM Capabilities

As LLMs evolved, several key advancements emerged:

Scaling Laws

Kaplan et al. (2020) established important relationships between model size, dataset size, and computational budget [11]. Their findings showed that model performance

scales predictably with these factors, guiding the development of more efficient and powerful models. This work provided a theoretical foundation for the "bigger is better" approach in language model development.

Few-shot and Zero-shot Learning

Brown et al. (2020) demonstrated with GPT-3 that sufficiently large models could perform tasks with minimal or no task-specific training examples [3]. This capability, known as few-shot and zero-shot learning, opened new possibilities for creating more flexible and adaptable AI systems.

Instruction Tuning

Wei et al. (2022) showed that fine-tuning models on diverse sets of instructions could significantly improve their ability to follow natural language prompts [29]. This technique, known as instruction tuning, enhanced the versatility of LLMs across various tasks and made them more aligned with human intentions.

Chain-of-Thought Prompting

Another significant advancement came from Wei et al. (2022), who demonstrated that prompting models to generate step-by-step reasoning could substantially improve their performance on complex tasks, including those involving numerical reasoning [30]. This technique, called chain-of-thought prompting, showed that LLMs could be guided to break down complex problems into more manageable steps, mirroring human problem-solving processes.

Reinforcement Learning from Human Feedback (RLHF)

Ouyang et al. (2022) introduced techniques to align language models with human preferences using reinforcement learning [19]. This approach, known as RLHF, led to the development of models that could produce more helpful, safe, and contextually appropriate outputs.

Chapter 3

Related Work

The application of Large Language Models (LLMs) in the financial sector has been a subject of intense research in recent years, with studies exploring their potential and limitations across various financial tasks. This chapter reviews key works in this field, highlighting methodologies, findings, and persistent research gaps. General-purpose LLMs like GPT-3, ChatGPT, and GPT-4 have demonstrated broad capabilities across various financial tasks. Li et al. (2023) conducted a comprehensive evaluation of these models using multiple benchmark datasets [17]. Their methodology involved fine-tuning these models on financial data and comparing their performance against specialized financial models. While the general-purpose LLMs often outperformed specialized models in tasks like sentiment analysis and named entity recognition, they struggled with complex financial reasoning tasks. This highlights a significant research gap: the need for LLMs that can perform deep, domain-specific financial analysis while maintaining their general language understanding capabilities. Recognizing these limitations, researchers have begun developing models specifically tailored for the financial domain. Wu et al. (2023) introduced BloombergGPT, a 50-billion parameter language model trained on a vast corpus of financial data [31]. Their approach involved pretraining the model on a mix of general and financial text data, followed by fine-tuning on specific financial tasks. BloombergGPT demonstrated significant improvements over general-purpose models in tasks such as financial sentiment analysis and question answering. However, the authors noted that the model still struggled with complex numerical reasoning, highlighting a persistent gap in LLMs' ability to perform accurate financial calculations consistently. Wang et al. (2023) proposed FinGPT, leveraging instructional tuning to adapt LLMs for financial applications [27]. Their methodology involved fine-tuning LLMs on a diverse set of financial instructions and tasks. The

authors emphasized the importance of systematic evaluation across various financial competencies, from basic tasks to complex multitasking scenarios. While FinGPT showed improvements on text-based financial data, it overlooked other critical financial data forms such as numerical and tabular data, which are pivotal for comprehensive financial analysis.

Efforts to enhance LLM capabilities for finance have taken various forms. Chen et al. (2022) introduced the CONVFINQA dataset, focusing on complex numerical reasoning in conversational finance [5]. Their work involved creating a dataset of multi-turn financial conversations that require chained reasoning over numbers. They experimented with both neural-symbolic approaches and prompting-based methods, finding that even state-of-the-art LLMs struggled with complex, multi-step financial calculations. This work highlights a critical research gap: the need for LLMs that can perform reliable, multi-step numerical reasoning in financial contexts. Addressing the challenge of improving LLMs' instruction-following capabilities, Wang et al. (2023) proposed Self-Instruct, a method for enhancing LLMs through self-generated instructions [28]. While not specifically focused on finance, their approach of using LLMs to generate their own training data could potentially be applied to financial domains. The authors demonstrated a 33 percent improvement in GPT-3's general instruction-following capabilities. However, they noted limitations in the quality and diversity of self-generated instructions, pointing to a research gap in developing more sophisticated self-improvement mechanisms for LLMs. Araci (2023) introduced FinBERT, a BERT-based model specifically designed for financial sentiment analysis [1]. The author's methodology involved pretraining BERT on a large corpus of financial texts and fine-tuning it on labeled financial sentiment data. While FinBERT showed improvements over general-purpose models in sentiment classification tasks, it was limited to sentiment analysis and did not address more complex financial reasoning tasks, highlighting the need for more versatile financial LLMs. The challenge of context modeling and reasoning in LLMs, crucial for many financial applications, was addressed by Zhao et al. (2023) in their work on natural language-based context modeling and reasoning with LLMs [32]. They outlined various strategies for improving LLMs' ability to understand and utilize context, including prompt engineering, few-shot learning, and retrieval-augmented generation. While their work was not finance-specific, their findings have significant implications for financial applications. For instance, improved context modeling could enhance LLMs' ability to understand complex financial narratives or multi-turn financial conversations. Addressing the computational

challenges of fine-tuning large models for specific tasks, Dettmers et al. (2023) proposed QLoRA (Quantized Low-Rank Adaptation), a parameter-efficient fine-tuning method [6]. Their approach involves quantizing the pretrained language model to 4 bits, adding trainable low-rank adapters, and using a novel preconditioner which they term "Paged Optimizers" to handle GPU memory constraints. This method could potentially allow for more efficient adaptation of large, general-purpose LLMs to specific financial tasks without the need for extensive computational resources. Recent advancements in LLM capabilities have opened new avenues for their application in finance. Schick et al. (2023) introduced the Toolformer approach, which enables language models to learn to use external tools through self-supervised learning [23]. Their methodology involved augmenting a pretraining dataset with tool-use examples, allowing the model to learn when and how to call external tools. While not specifically focused on finance, this approach has significant implications for financial applications. For instance, a Toolformer-like model could potentially learn to access real-time market data, financial calculators, or regulatory databases, addressing the research gap of integrating external financial tools and data sources with LLMs. However, the authors noted challenges in tool selection and result interpretation, highlighting a persistent gap in LLMs' ability to reason about tool outputs more so in a complex domains like finance.

Chapter 4

Methodology

This chapter describes the methodology applied in this project.

4.1 Dataset Selection and Preparation

4.1.1 Merged Dataset Creation

We have used a mixture of four different financial datasets. The brief summary of these dataset is provided below. The structured sample for each of these dataset after preprocessing is provided in the appendix section.

TAT-QA contains 16,552 questions associated with 2,757 hybrid contexts from real-world financial reports. The questions typically require a range of data extraction and numerical reasoning skills, including multiplication, comparison, sorting, and their various combinations [33].

ConvFinQA contains 3,892 conversations consisting 14,115 questions from real-world scenario of conversational question answering over financial reports. The dataset is formulated using both textual content and structured table [5].

FINQA is an expert annotated dataset that contains 8,281 financial QA pairs, along with their numerical reasoning processes. The reasoning processes answering these questions are made of many common calculations in financial analysis, such as addition, comparison, and table aggregation [4].

FinGPT FinRED-RE dataset available on Hugging Face is designed for financial relationship extraction tasks. It contains 13.5k rows of data, divided into 11.4k training rows and 2.14k test rows. The dataset features text inputs with associated financial entities and their relationships. The task contains instructions for extracting financial

relationships from textual data, utilizing specific relations like employer, industry, and product/material produced ¹.

Merged dataset: The datasets are merged to form a comprehensive unified dataset, balancing various types of financial data to ensure robust coverage of multiple scenarios and data formats encountered in the financial domain. This integration includes tabular data, text-based financial reports, conversational question-answer pairs, and extracted relationships, enhancing the dataset's versatility and applicability for a wide range of financial data analysis tasks. This comprehensive approach improves the ability to analyze diverse financial situations and generalize across tasks, making it valuable for both financial domain applications and general numerical reasoning scenarios.

4.1.2 Dataset Preprocessing

- **Algorithm Development:** A preprocessing algorithm 1 was created to prepare the merged dataset for model training. This algorithm ensures that the data follows a unified structure comprising context, questions, and answers.
- **Prompt Template:** The data was formatted into a specific template prompt format required by both Llama2 7b and Llama3 8b models that are used in this project. This involved:
 - **Data Cleaning:** Removing inconsistencies and irrelevant information.
 - **Normalization:** Standardizing numerical values and textual content to maintain uniformity.
 - **Prompt Structuring:** Organizing the data into the structured prompts that the models were originally trained on, ensuring compatibility and effectiveness during fine-tuning. This include use of special tokens and prompt structure of the specific model ².

4.2 Baseline Model Selection and Evaluation

- **Baseline Models Used:**
 - meta-llama/Llama-2-7b-hf ³

¹<https://huggingface.co/datasets/FinGPT/fingpt-finred-re>

²<https://llama.meta.com/docs/model-cards-and-prompt-formats/meta-llama-3>

³<https://huggingface.co/meta-llama/Llama-2-7b-hf>

Algorithm 1 Preprocessing Merged Dataset

```

1: procedure PREPROCESSDATASETS
2:   Load TAT-QA, ConvFinQA, FinQA, and RelEx datasets
3:   combined_data  $\leftarrow \emptyset$ 
4:   for dataset  $\in \{TAT-QA, ConvFinQA, FinQA, RelEx\}$  do
5:     processed  $\leftarrow$  PROCESSDATASET(dataset)
6:     combined_data  $\leftarrow$  combined_data  $\cup$  processed
7:   end for
8:   combined_data  $\leftarrow$  RANDOMSAMPLE(combined_data,  $\lfloor |combined\_data|/3 \rfloor$ )
9:   Split combined_data into train_data (90%) and test_data (10%)
10:  return train_data, test_data
11: end procedure
12: procedure PROCESSDATASET(dataset)
13:   processed  $\leftarrow \emptyset$ 
14:   for each sample  $\in$  dataset do
15:     (context, question, answer)  $\leftarrow$  EXTRACTINFO(dataset, sample)
16:     processed  $\leftarrow$  processed  $\cup \{(\textit{context}, \textit{question}, \textit{answer})\}$ 
17:   end for
18:   return processed
19: end procedure
20: procedure EXTRACTINFO(dataset, sample)
21:   if dataset = TAT-QA then
22:     context  $\leftarrow$  Combine paragraphs and format table
23:     for each q  $\in$  sample.questions do
24:       return (context, q.question, q.answer)
25:     end for
26:   else if dataset = ConvFinQA then
27:     return (sample.input, sample.instruction + instructions, sample.output)
28:   else if dataset = FinQA then
29:     context  $\leftarrow$  Combine pre_text, post_text, and table
30:     return (context, sample.qa.question, sample.qa.answer)
31:   else if dataset = RelEx then
32:     return (sample.input, sample.instruction, sample.output)
33:   end if
34: end procedure

```

- meta-llama/Meta-Llama-3-8B ⁴
- Performance Testing: The merged dataset was tested on these baseline models to establish initial performance metrics. This step was crucial for:
 - Benchmarking: Establishing a reference point for comparing improvements post fine-tuning.
 - Identifying Weaknesses: Understanding the initial strengths and weaknesses of the models with raw financial data.

4.3 Fine-tuning Approaches

4.3.1 Training Procedure

Model Architecture

We employed the Meta-Llama-2-7B-hf and Meta-Llama-3-8B model, a variant of the LLaMA architecture specifically adapted for few-shot learning scenarios. This choice was driven by its recent success in various NLP tasks, particularly those requiring nuanced understanding and generation based on limited context. For the code generation and execution task, we have used the recently launched Meta-Llama-3.1-8B model due to its enhanced ability for code generation tasks.

Dataset Preparation

For the first experiment of fine tuning, We curated a dataset from multiple sources as mentioned in section 4.1, aiming to encompass a wide range of financial topics to enhance the model's ability to generalize across diverse financial contexts. The data was split into 90% for training and 10% for validation, ensuring a representative distribution of topics in each subset.

For the X experiment, we exclusively used the TAT-QA dataset, following a pre-processing strategy similar to the one outlined in the algorithm 1, specifically steps 20-25.

⁴<https://huggingface.co/meta-llama/Meta-Llama-3-8B>

Training Details

The model was fine-tuned using a quantization-aware training approach, utilizing 4-bit quantization to balance performance with computational efficiency. The BitsAndBytes library facilitated the integration of low-precision arithmetic during training. Quantization is a two-step process, involves normalizing constants to scale vectors into a target range and rounding to the nearest target value. This technique can cause significant quantization loss if weights have outliers. To address this, *bitsandbytes* employs vector-wise quantization and mixed precision decomposition, maintaining performance akin to non-quantized states [9]. Although LLM int8 does not impair performance, it increases inference time due to quantization overhead but significantly reduces memory usage by 71%, which enabled us to fine tune the models on NVIDIA GPUs. Parameter-Efficient Fine-Tuning (PEFT) enhances the performance of pre-trained language models for specific tasks in Natural Language Processing. By adjusting only a subset of the model's parameters on smaller datasets, PEFT conserves computational resources and time. This method typically involves freezing several layers of the model and fine-tuning only the final layers that directly pertain to the target application, thereby achieving greater efficiency [9]. Low-Rank Adaptation (LoRA)⁵ is an efficient training technique for large language models that significantly reduces the number of trainable parameters by inserting a smaller set of new weights, which are the only parts trained. This method speeds up training, enhances memory efficiency, and results in much smaller model weights (only a few hundred megabytes), making the models easier to manage and distribute. We employed the LoraConfig from the PEFT library, setting a rank of 16 and an alpha of 64 to finely tune the attention mechanism to our dataset while minimizing hardware needs without extra inference latency. In our experiment fine-tuning the Llama-3-8b model, the original parameter count was 4,582,543,360. With LoRA, we reduced it to 41,943,040 trainable parameters, constituting just 0.915% of the total parameters. This reduction underscores the efficiency of LoRA in managing computational resources.

In our training setup, we employed a distributed system featuring eight NVIDIA GeForce RTX 3090 GPUs, leveraging PyTorch's DistributedDataParallel (DDP) framework. This choice was informed by insights from the work by Shen Li et. al (2020) which demonstrated DDP's superiority in synchronizing gradients efficiently across multiple GPUs. They noted that DDP minimizes communication overhead and optimizes computation by overlapping gradient reduction with backpropagation [15]. This

⁵<https://huggingface.co/docs/diffusers/en/training/lora>

results in enhanced training speed and scalability, crucial for handling large datasets and complex models in a distributed environment. This approach was particularly necessary for fine-tuning Llama models in our resource-constrained environment, where utilizing multiple GPUs for data-parallel training was essential.

Gradient accumulation ⁶ allows for the use of larger batch sizes than those limited by hardware memory by summing up gradients over multiple mini-batches and updating the model only after a predefined number of these batches. In our training setup, we managed a batch size of one per device, accumulating gradients over 12 steps. This strategy effectively simulates training with larger batch sizes, optimizing the trade-off between memory usage and training convergence speed.

AdamW is an optimization algorithm that modifies the classic Adam optimizer by decoupling weight decay from the gradient updates [18]. This adjustment allows for more effective and theoretically sound management of weight decay, improving generalization compared to standard weight decay in optimizers like Adam. The modification ensures that the weight decay is applied directly to the weights themselves rather than as part of the gradient descent, which helps in better preserving the training stability and often leads to better performance on validation and test datasets. We utilized the AdamW optimizer with a learning rate of 2E-4, chosen based on preliminary testing to ensure rapid convergence without compromising stability. The model underwent training over five epochs, incorporating early stopping based on validation loss to mitigate overfitting, thus enhancing model generalizability and performance efficiency.

In the code generation and execution experiment conducted on a Google Colab platform with a single NVIDIA A100 40GB GPU, the recently introduced META-LLAMA/META-LLAMA-3.1-8B model was deployed to generate Python code for financial analysis. This experiment was distinctive as it did not involve model fine-tuning. Instead, few-shot learning techniques were utilized, which included an enhanced prompt template featuring context, a specific question, and Python code examples, formatted with actual newline characters. To identify the most contextually relevant examples for the model, TF-IDF vectorization [24] and cosine similarity metrics ⁷ were applied. This methodological choice was influenced by the findings of Omid et al. (2019), who observed that despite the prevalent expectation favoring advanced topic models (e.g., Latent Semantic Indexing) and neural models (e.g., Paragraph Vectors) for superior performance across all measures of textual similarity, TF-IDF demonstrated

⁶https://huggingface.co/docs/accelerate/en/usage_guides/gradient_accumulation

⁷https://scikit-learn.org/stable/modules/generated/sklearn.metrics.pairwise.cosine_similarity.html

remarkable efficacy. Specifically, TF-IDF excelled in scenarios involving longer and more technical texts and in making more nuanced distinctions among nearest neighbors. This attribute proved particularly advantageous for the current study, which involves complex financial datasets, thereby necessitating a method capable of handling the intricacies and technicalities embedded within such texts.

Evaluation Strategy

Model performance was periodically evaluated on the validation set at the end of a predefined step in the training setup.

Few-Shot Example Configuration

In Section 5.3.3, we detail the outcomes of the experiment involving few-shot prompting techniques, where we evaluated various few-shot learning methods. Few-shot learning (FSL) involves training models to recognize categories from very limited examples. One-shot learning (OSL) refers to learning tasks where only a single example per rare category is available. Zero-shot learning (ZSL), meanwhile, deals with categories for which no examples are provided [2]. Suvarna et al. (2019) discusses the challenges and techniques of generalization in few-shot learning, emphasizing that while machines need numerous examples to learn like humans, they lack certain cognitive functions. In contrast to earlier methods requiring human intervention for learning representations, modern deep learning autonomously learns these representations. However, learning effective representations from few examples remains difficult. Deep models require diverse, *representative examples* to generalize well but often struggle in few-shot scenarios due to the scarcity of such examples [10]. To tackle the challenges of few-shot learning, we developed a strategy that utilizes both relevant and random examples. The selection of relevant examples is based on cosine similarity⁸ measurements between the embeddings of training samples and a predefined set of few-shot examples. These embeddings are generated using the SentenceTransformer model, specifically *all-MiniLM-L6-v2*⁹ which provides a dense representation of text ideal for similarity assessments. The selection process is detailed in Algorithm 2, where TOP_K indicates the number of relevant samples calculated. In our experiments, TOP_K is set to 2 due to memory constraints and the context length limitations of the Llama-2-7b-hf and Llama-3-8b models, which are 4096 and 8192 tokens, respectively. The parameters

⁸https://scikit-learn.org/stable/modules/generated/sklearn.metrics.pairwise.cosine_similarity.html

⁹<https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

context_weight and *question_weight* are both set to 0.5, reflecting the equal importance of context and questions in datasets such as ConvFinQA, where the history and the final question are crucial. This approach ensures that no input is truncated during training, although it limits our ability to test performance with a larger number of examples, i.e., a much higher TOP_K.

In the experiment for code generation detailed in section 5.4.1, we used few-shot learning techniques, leveraging an improved prompt template that included context, a question, and Python code examples, with a focus on proper formatting using actual newline characters. To select the most relevant examples for the model, we employed TF-IDF vectorization and cosine similarity measures. This approach ensured that the examples used in the prompt were highly pertinent to the task at hand, enhancing the model's ability to generate accurate and executable code.

4.3.2 Post Processing Algorithm

During inference, the model generates responses based on varying prompt templates, which differ across experiments and models. Even with the same model, the prompts vary depending on whether we use few-shot prompts or additional chain-of-thought prompting to guide the model to generate explanations. Consequently, a robust post-processing algorithm is necessary. Algorithm 3 presents a generic post-processing approach used during inference and evaluation, with slight modifications based on the model's output response. The code for all post-processing steps is available in my Github Repo¹⁰.

4.4 Evaluation and Iteration

4.4.1 Permissive Accuracy

We have designed a permissive accuracy measure that accommodates minor precision differences by allowing slight deviations in decimal points within an alpha threshold. This measure is particularly useful in financial contexts, where exact numerical precision may not always be critical. The method involves extracting numerical values and comparing them within a tunable alpha difference.

For threshold setting, we define acceptable deviation ranges, ensuring flexibility in

¹⁰<https://github.com/rjanant/dissertation>

Algorithm 2 Precompute Relevant Examples for Few-Shot Learning

```

1: function GET_EMBEDDINGS(text)
2:   inputs  $\leftarrow$  embedding_tokenizer(text, return_tensors='pt', padding=True, truncation=True)
3:   embedding  $\leftarrow$  embedding_model(**inputs).last_hidden_state.mean(dim=1).cpu().numpy()
4:   return embedding[0]
5: end function
6: function PRECOMPUTE_RELEVANT_EXAMPLES(dataset, few_shot_examples, top_k,
   context_weight, question_weight)
7:   few_shot_context_embeddings  $\leftarrow$  [GetEmbeddings(ex["context"]) for ex in
   few_shot_examples]
8:   few_shot_question_embeddings  $\leftarrow$  [GetEmbeddings(ex["question"]) for ex in
   few_shot_examples]
9:   relevant_examples_map  $\leftarrow$  {}
10:  for idx, item in enumerate(dataset) do
11:    context_embedding  $\leftarrow$  GetEmbeddings(item['context'])
12:    question_embedding  $\leftarrow$  GetEmbeddings(item['question'])
13:    context_similarities  $\leftarrow$  cosine_similarity([context_embedding],
   few_shot_context_embeddings)[0]
14:    question_similarities  $\leftarrow$  cosine_similarity([question_embedding],
   few_shot_question_embeddings)[0]
15:    combined_similarities  $\leftarrow$  context_weight * context_similarities + ques-
   tion_weight * question_similarities
16:    most_relevant_indices  $\leftarrow$  argsort(combined_similarities)[-top_k:][::-1]
17:    relevant_examples_map[str(idx)]  $\leftarrow$  most_relevant_indices.tolist()
18:  end for
19:  return relevant_examples_map
20: end function

```

Algorithm 3 Post-Processing Model Output

```

1: function PREPROCESSOUTPUT(output)
2:   Input: Model's output string
3:   Pattern Matching: Use regex to find the answer section
4:   answer_pattern  $\leftarrow$  r"Answer (including calculation steps and final answer,
5:   use 'n' for line breaks):(.*)?"
6:   (?:(nContext:—$))"
10:  match  $\leftarrow$  re.search(answer_pattern, output, re.DOTALL)
11:  if match is not None then
12:    answer  $\leftarrow$  match.group(1).strip()
13:    lines  $\leftarrow$  [line.strip() for line in answer.split('\n') if line.strip()]
14:    return lines[0] if lines else ""
15:  else
16:    return ""
17:  end if
18: end function

```

various applications. This measure is especially relevant in scenarios like relational extraction, where generated outputs are textual values. In such cases, the Sequence-Matcher is used for string comparison, returning 1 if the similarity exceeds β ¹¹ or if key parts of the prediction and reference match. This approach is beneficial when the model generates the correct answer but includes additional values, which can be ignored. Additionally, we have an exact match function to report precise match accuracy, ensuring comprehensive evaluation of model performance.

4.4.2 Exact Match Accuracy

In the context of financial datasets, exact match accuracy is vital for ensuring the precision of numerical output generation. Given the importance of accuracy in financial analysis, even a minor deviation, such as a difference in decimal points, is unacceptable and considered an error. This metric enforces a strict criterion by directly comparing the generated output with the expected answer, ensuring only perfectly matching results are considered correct.

¹¹ $\beta = 0.9$ & $\alpha = 0.001$

Algorithm 4 Permissive Accuracy

```

1: function PERMISSIVE_ACCURACY(predicted, actual)
2:   Input: predicted, actual
3:   predicted  $\leftarrow$  str(predicted).lower().strip()
4:   actual  $\leftarrow$  str(actual).lower().strip()
5:   pred_num  $\leftarrow$  ExtractNumericalValue(predicted)
6:   actual_num  $\leftarrow$  ExtractNumericalValue(actual)
7:   if pred_num is not None and actual_num is not None then
8:     Return int(abs(pred_num - actual_num) <  $\alpha$ ) else
9:     pred_words  $\leftarrow$  predicted.split()
10:    actual_words  $\leftarrow$  actual.split()
11:    similarity  $\leftarrow$  SequenceMatcher(None, pred_words,
12:    actual_words).ratio()
13:    key_parts_match  $\leftarrow$  all(part in predicted for part in
14:    actual.split(':')[0].split(','))
15:    Return int(similarity > 0.9 or key_parts_match)
16:   end if
17: end function

```

4.4.3 Inference Phase Metrics

4.4.3.1 Comprehensive Evaluation: Using a suite of metrics to assess various aspects of model performance

- ROUGE: Measures overlap of n-grams between generated and reference texts.
- METEOR: Considers precision, recall, and synonymy for evaluating translation quality.
- BLEU Score: Evaluates the precision of n-grams in generated text against reference.
- BERT Precision, Recall, F1 Score: Uses contextual embeddings to measure semantic similarity.
- GLUE Benchmark: It is a collection of resources for training, evaluating, and analyzing natural language understanding systems across diverse tasks.

4.5 External Tool Integration & Few-Shot Code Generation

[Ask potential research question and try to answer those questions in this section as well as in the experiments - Results and Discussion section — Also add diagram of this code generation approach in a well defined manner like the one on langgraph for code gen page — In Experiments talk about use of atmax 5 rel examples but couldn't due to model limit context length and on truncation losing useful context for specific dataset i.e. convfinqa - add sample in appendix from each dataset]

Recent advancements in tackling complex reasoning tasks, have opened new avenues for applying large language models to specialized domains like financial analysis. Notably, the study by Suzgun et al.(2022) provides compelling evidence for the efficacy of step-by-step approaches in solving intricate problems. The study focused on BIG-Bench, a collaborative benchmark comprising over 200 diverse text-based tasks, including traditional NLP, mathematics, commonsense reasoning, and question-answering [25]. Particularly relevant to our research is their curation of BIG-Bench Hard (BBH), a subset of 23 especially challenging tasks where previous models had not surpassed average human-rater performance.

The study's results are particularly illuminating: using various *chain-of-thought* (CoT) approaches, they found that the Codex model with CoT prompting outperformed the average human rater score on 17 out of 23 chosen tasks. This remarkable performance on diverse, complex tasks underscores the potential of step-by-step reasoning approaches in tackling challenging problems. Drawing parallels to our research, we recognized that the complex nature of financial data analysis - involving intricate calculations, temporal considerations, and the integration of tabular and textual data - could benefit from a similar methodical approach.

Inspired by these findings, our research leverages a step-by-step code generation methodology to address the unique challenges posed by financial datasets. By breaking down complex financial analyses into discrete, programmable steps, we aim to enhance the accuracy and reliability of our model's outputs. This approach allows us to handle the multifaceted nature of financial data, including numerical precision in calculations, temporal dependencies in time-series data, and the integration of qualitative information from textual sources.

Our methodology extends the concept of chain-of-thought reasoning to the domain of financial analysis, where each step in the chain is represented by a generated code snippet. This not only facilitates more accurate computations but also provides transparency and interpretability in the analysis process - crucial factors in the financial sector where decisions often require clear justification and auditable processes.

Our experimental design commenced with the careful selection of a small, hand-crafted sample derived from a diverse array of financial question-answer pairs. These samples were meticulously chosen from the comprehensive dataset detailed in Section 4.1.1, ensuring a representative cross-section of financial queries and their corresponding responses.

To enhance the analytical capabilities of our model, we implemented a novel approach: replacing the original answers in the dataset with Python scripts. These scripts were generated with the assistance of GPT-4, a state-of-the-art language model, and subsequently underwent rigorous manual verification to ensure accuracy and functionality. This transformation allowed for a more dynamic and computationally rich representation of financial data analysis. The sample dataset was intentionally constructed to encompass a wide spectrum of answer types, including Numerical responses, Textual answers and Temporal data. This diversity in response types was crucial for assessing the model's versatility in handling various financial data representations. For the execution of these Python scripts embedded within the answer fields, we leveraged the

LangChain LLM framework integrated into our pipeline. The execution environment was facilitated by a Python agent interfacing with a Python REPL (Read-Eval-Print Loop) tool. The REPL tool was instrumental in enabling real-time execution and evaluation of Python code, a feature paramount for the dynamic and precise computation of financial metrics and analyses. The language model underpinning this integration was the FACEBOOK/OPT-1.3B MODEL, selected for its robust natural language processing capabilities and its ability to understand and generate context-aware responses. As demonstrated by [8], Nucleus Sampling is highly effective for neural generation. By sampling from the dynamic nucleus of the probability distribution, it balances diversity and reliability, resulting in text that mirrors human quality while maintaining fluency and coherence. In our pipeline, we use parameters: temperature=0.7, top_p=0.95, and repetition_penalty=1.15. Using these parameters means controlling the randomness of the generation (temperature=0.7), focusing on the top 95% of the probability mass (top_p=0.95), and discouraging repetitive sequences (repetition_penalty=1.15). The results are presented and analyzed in Chapter 5.

In light of the promising results obtained from the initial experiment, we proceeded to evaluate the latest META-LLAMA/META-LLAMA-3.1-8B-INSTRUCT model, a cutting-edge language model. This model demonstrated superior performance in code generation on the HumanEval benchmark, surpassing the GPT-3.5 Turbo model, and also achieved comparable results on reasoning tasks using the ARC Challenge benchmark ¹², thereby establishing its suitability for our study.

The process begins with a dynamic few-shot prompting mechanism designed to optimize the model's responses to specific financial queries. To ensure contextual alignment, we utilize TF-IDF vectorization and cosine similarity to identify and select the most relevant examples from our dataset for each query. These selected examples are then integrated into an enhanced few-shot prompt template contained in ??, along with the corresponding context and question, to guide the model in generating accurate Python code.

The code generation process is fine-tuned with parameters such as temperature (0.9) and top_p (0.95) to strike a balance between creativity and coherence. The generated code is subsequently extracted using regex pattern matching and executed within a controlled environment. This environment is equipped with the necessary libraries and contextual data, facilitating comprehensive financial data manipulation and analysis.

To evaluate the effectiveness of our approach, we processed a randomly selected

¹²<https://ai.meta.com/blog/meta-llama-3-1/>

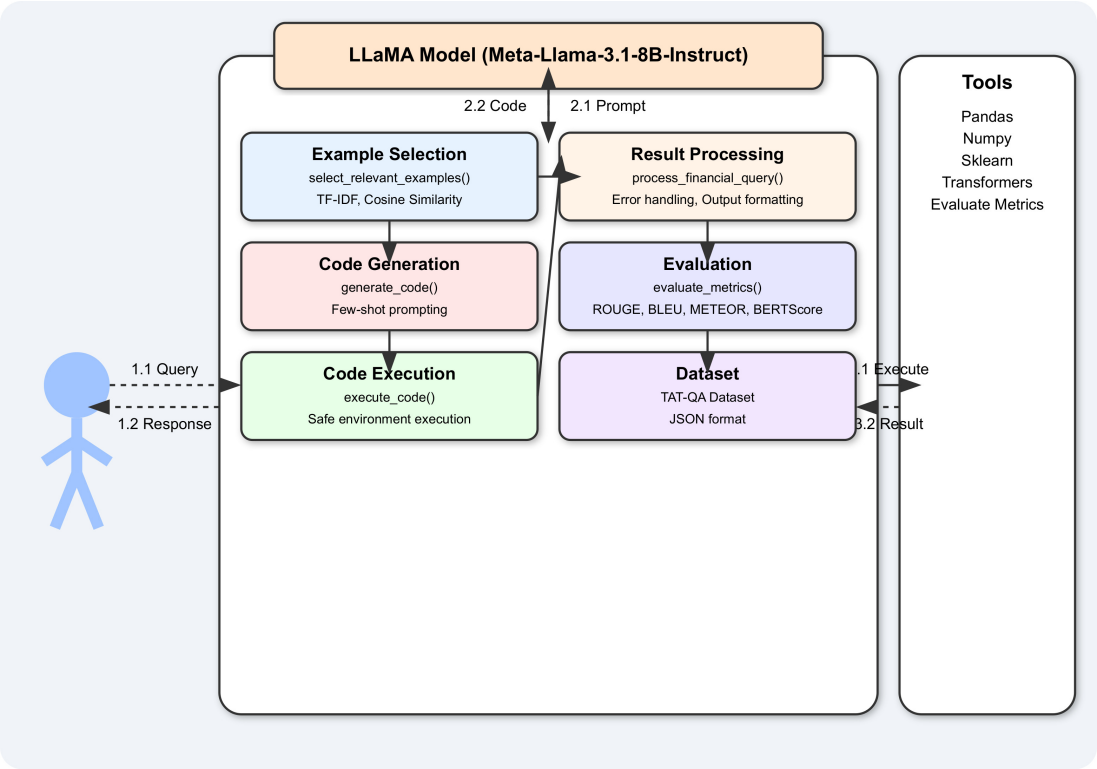


Figure 4.1: LLaMA-3.1-8B Workflow: From Query to Code Execution

subset of 1,000 samples from the TAT-QA dataset, ensuring the representativeness of the sample. For each sample, Python code was generated, executed, and the output was compared against a reference answer. The evaluation utilized the same suite of metrics detailed in Section 4.4.3, ensuring a rigorous assessment of the code generation approach within the context of financial tasks. The flowchart is explained in Figure 4.2

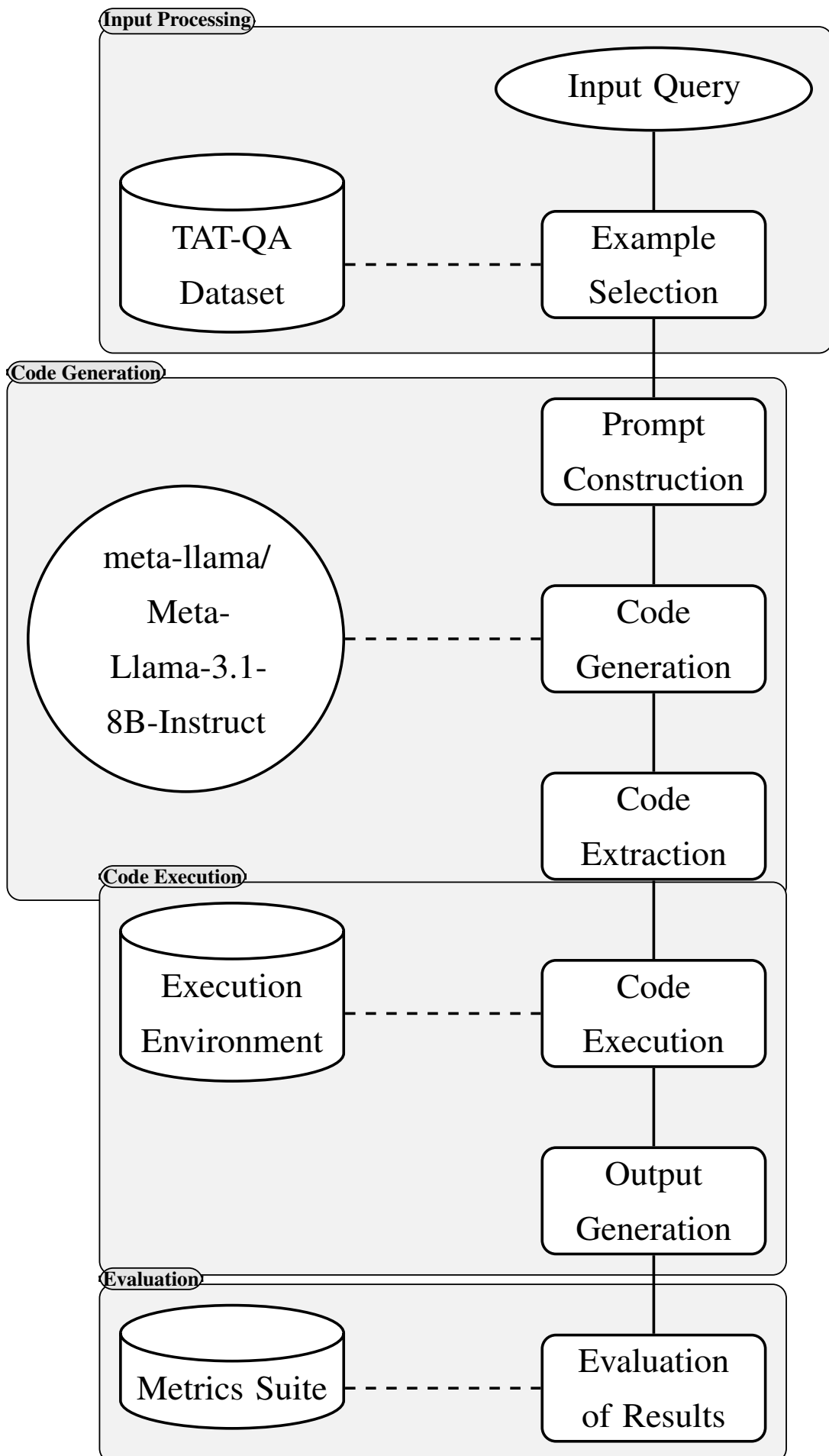


Figure 4.2: Flowchart of the Code Generation Process for Financial Analysis

Chapter 5

Experiments

This chapter presents the experimental setup, methodologies, and results of our research aimed at enhancing the number handling capabilities of language models in financial contexts. We explore various approaches, including few-shot prompting and fine-tuning, to outperform existing models across diverse financial tasks. We introduced a novel approach - code generation from input and running it through python exec in a pipeline like structure to enhance the numerical capability of the model. These are detailed in section 5.3.7.

5.1 Research Questions

Our experiments were designed to address the following research questions:

- How do different model architectures (Llama2-7B vs. Llama3-8B) compare in their baseline performance on financial tasks?
- What is the impact of fine-tuning on the models' performance across various metrics?
- How effective are different few-shot prompting strategies (e.g., selective, random, chain-of-thought) compared to fine-tuning?
- What is the optimal number and type of examples for few-shot prompting in financial tasks?
- How does the performance of language models compare to specialized tool-based approaches in financial data processing?

5.2 Experimental Setup

5.2.1 Dataset

We utilized a composite dataset combining four financial datasets which are explained in section 4.1 in detail. This diverse dataset allows us to evaluate model performance across various financial data types and tasks.

5.2.2 Models

We experimented with the following models:

- Llama2-7B chat (baseline and fine-tuned)
- Llama3-8B (baseline and fine-tuned)
- Llama3.1-8b (baseline)
- Llama3.1-8b-instruct (baseline)

5.2.2.1 Approaches

We implemented and compared several approaches:

1. Baseline performance (zero-shot)
2. Fine-tuning
3. Few-shot prompting with chain-of-thought (FSP + COT)
4. Selective few-shot prompting
5. Random few-shot prompting
6. Tool-based approach using Langchain (on a subset of 150 samples)
7. Code generation and execution (on TAT-QA dataset)

5.2.3 Evaluation Metrics

We employed a comprehensive set of metrics to evaluate model performance which are all detailed in section 4.4.3:

- Average Permissive Accuracy (custom metric allowing for slight numerical differences and partial text matches)
- Average Exact Accuracy
- ROUGE Score
- BLEU Score
- METEOR Score
- BERTScore (Precision, Recall, F1)
- GLEU Score

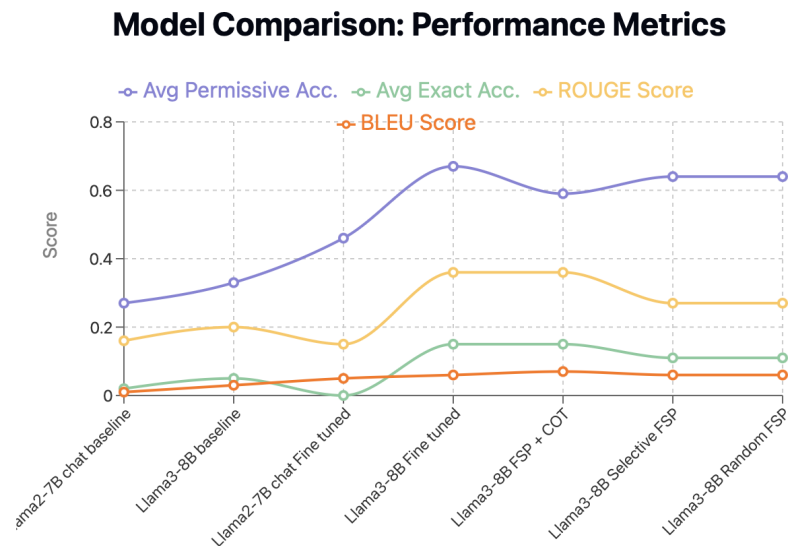


Figure 5.1: Performance Metrics of various models

5.3 Results and Analysis

5.3.1 Model Architecture Comparison

Table 5.1 presents the evaluation metrics for various models and approaches.

Key observations:

Table 5.1: Evaluation Metrics for Various Models

Model	Avg Permissive Acc.	Avg Exact Acc.	ROUGE Score	BLEU Score	METEOR Score	BERTScore P	BERTScore R	BERTScore F1	GLEU Score
Llama2-7B chat baseline	0.27	0.02	0.16	0.01	0.19	0.814	0.860	0.836	0.02
Llama3-8B baseline	0.33	0.05	0.20	0.03	0.19	0.836	0.876	0.855	0.04
Llama2-7B chat Fine tuned	0.46	0.00	0.15	0.05	0.21	0.778	0.892	0.830	0.07
Llama3-8B Fine tuned	0.67	0.15	0.36	0.06	0.31	0.850	0.919	0.881	0.07
Llama3-8B FSP + COT	0.59	0.15	0.36	0.07	0.31	0.854	0.917	0.883	0.07
Llama3-8B Selective FSP	0.64	0.11	0.27	0.06	0.28	0.809	0.913	0.855	0.05
Llama3-8B Random FSP	0.64	0.11	0.27	0.06	0.28	0.807	0.912	0.854	0.05

- Llama3-8B baseline outperforms Llama2-7B chat baseline across all metrics, indicating that the newer architecture has inherently better number handling capabilities.
- The performance gap is particularly noticeable in Average Permissive Accuracy (0.33 vs. 0.27) and Average Exact Accuracy (0.05 vs. 0.02).

5.3.2 Impact of Fine-tuning

Fine-tuning significantly improved the performance of both model architectures:

- Llama2-7B chat: Average Permissive Accuracy increased from 0.27 to 0.46
- Llama3-8B: Average Permissive Accuracy increased from 0.33 to 0.67

The fine-tuned Llama3-8B model showed the best overall performance, with notable improvements in ROUGE Score (0.36) and METEOR Score (0.31).

5.3.3 Few-shot Prompting Strategies

Different few-shot prompting strategies yielded varying results:

- FSP + COT: Achieved high performance (0.59 Average Permissive Accuracy) but did not surpass fine-tuned models.
- Selective FSP and Random FSP: Both achieved 0.64 Average Permissive Accuracy, indicating that careful example selection may not always be necessary.

5.3.4 Optimal Few-shot Prompting

Table 5.2 shows the performance metrics for different few-shot prompting approaches.

Key findings:

Table 5.2: Performance Metrics for Different Few Shot Prompting

Fine Tuning Type	Avg Permissive Acc.	Avg Exact Acc.	ROUGE Score	BLEU Score	METEOR Score	BERTScore P	BERTScore R	BERTScore F1	GLEU Score
1 Relevant Example	0.62	0.12	0.32	0.06	0.28	0.841	0.914	0.874	0.06
1 Random Example	0.63	0.13	0.33	0.06	0.29	0.840	0.915	0.874	0.06
Results 1	0.61	0.14	0.34	0.06	0.28	0.852	0.916	0.881	0.07
Results 2	0.61	0.14	0.34	0.06	0.28	0.853	0.915	0.881	0.07
Results 3	0.61	0.14	0.35	0.06	0.28	0.855	0.916	0.883	0.07

- Using a single relevant example or a single random example yielded similar performance (0.62 vs. 0.63 Average Permissive Accuracy).
- The results suggest that the number of examples may be more important than their specific relevance in this context.

5.3.5 Example Type Distributions

Table 5.3 presents the distribution and weights of different example types used in the experiments.

The varied distributions across different result sets allow us to analyze the impact of example selection strategies on model performance.

Table 5.3: Example Type Distributions and Weights

Results	Combined	Relevant	Randomized	Baseline
Results 1	341 (0.25)	304 (0.25)	325 (0.25)	317 (0.25)
Results 2	147 (0.10)	324 (0.25)	499 (0.40)	317 (0.25)
Results 3	338 (0.25)	616 (0.50)	197 (0.15)	136 (0.10)

5.3.6 Tool-based Approach

The tool-based approach using Langchain on a subset of 150 samples showed promising results:

- Exact Accuracy: 0.75
- BLEU Score: 0.23
- METEOR Score: 0.33

The findings from this study underscore the efficacy of specialized tools in performing precise financial tasks, with a notable requirement for exact matches in data handling. It is pertinent to acknowledge that these results are derived from scenarios utilizing manually generated code, specifically using GPT-4 ¹, which represents a potentially idealized set of conditions.

An interesting continuation of this research could involve leveraging OpenAI's advanced code generation capabilities in conjunction with Langchain's Python agent REPL tool to automate and refine this process further. However, such an approach would entail additional costs, which were not within the purview of this current study. Nevertheless, the potential of commercial language models to adeptly manage the complexity and variability inherent in our financial dataset remains a compelling aspect of our ongoing investigation. This line of inquiry suggests substantial promise for enhancing automated financial analysis tools through the integration of sophisticated AI-driven technologies.

5.4 Code Generation Approach

In this series of experiments, we explore a novel approach to enhancing numerical precision in language models for financial tasks. Rather than relying on traditional fine-tuning methods, we leverage the models' inherent code generation capabilities through strategic few-shot prompting. This approach is particularly relevant in the context of financial data processing, where precise numerical handling is crucial. Our motivation stems from the observation that while large language models have demonstrated impressive natural language understanding and generation capabilities, their performance in tasks requiring precise numerical computations often falls short. By guiding these models to generate executable code, we aim to bridge this gap and achieve higher accuracy in financial calculations.

5.4.1 Experimental Setup

Dataset: We utilized the TAT-QA dataset [33], a benchmark for text-and-table question answering. To ensure a manageable yet representative sample, we randomly selected 1000 samples from the dataset. This selection was kept constant across all models using a fixed seed value for reproducibility.

¹<https://chatgpt.com>

Models: We experimented with several models from the Llama family:

1. Llama3-8b
2. Llama3.1-8b
3. Llama3.1-8b-instruct

Additionally, we included results from a FACEBOOK/OPT-1.3B model using Langchain for comparison, though this was tested on a smaller sample of 150 instances.

Our approach can be summarized in the following steps:

1. Few-shot prompting: We provided the models with example samples that demonstrate the desired code structure.
2. Code generation: The models were then prompted to generate similar code for new inputs.
3. Code execution: We used Python’s `exec()` function to execute the generated code and produce final results.
4. Evaluation: The results were evaluated using various metrics, including permissive accuracy, exact accuracy, ROUGE, BLEU, METEOR, BERTScore, and GLEU.

It’s worth noting that the Llama models used in this study, particularly the 8b variants, are not primarily known for their code generation capabilities. This choice was deliberate, as it allows us to test the limits of our approach and potentially demonstrate its efficacy even with models not specialized for this task.

5.4.2 Results and Analysis

Table 5.4 presents the evaluation metrics for various models and approaches.

Table 5.4: Evaluation Metrics for Code Generation

Model	Avg Permissive Acc.	Avg Exact Acc.	ROUGE Score	BLEU Score	METEOR Score	BERTScore P	BERTScore R	BERTScore F1	GLEU Score
facebook/opt-1.3b (langchain - Test samples 150)	-	0.75	0.30	0.23	0.33	0.822	0.878	0.849	0.23
Llama3-8b	0.39	0.30	0.40	0.18	0.32	0.870	0.898	0.884	0.19
Llama3.1-8b	0.40	0.30	0.39	0.17	0.32	0.868	0.898	0.882	0.18
Llama3.1-8b-instruct	0.52	0.43	0.49	0.29	0.34	0.887	0.909	0.897	0.24
Llama3.1-8b (code-gen)	0.23	0.18	0.25	0.05	0.14	0.841	0.874	0.856	0.06
Llama3.1-8b-instruct (code-gen)	0.35	0.28	0.41	0.14	0.21	0.866	0.898	0.881	0.14

Key observations:

Model Comparison: Code Generation Experiment Results

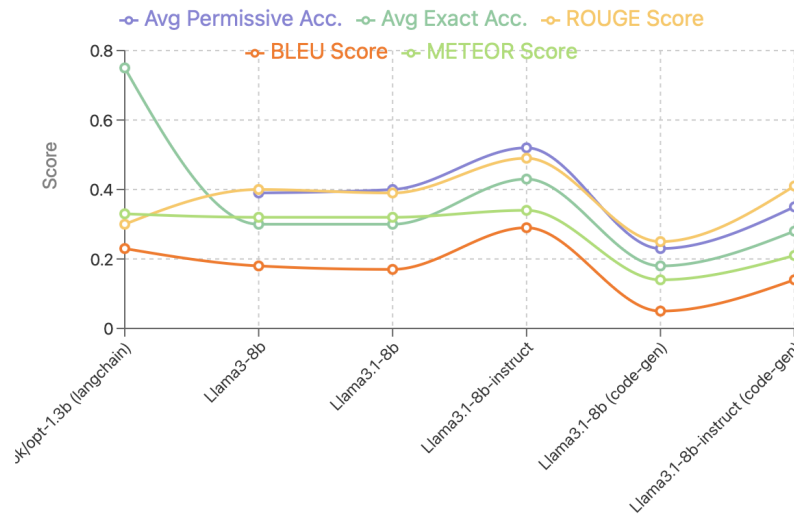


Figure 5.2: Performance Metrics of various models in Code Generation

1. Performance improvement with instruction tuning: The Llama3.1-8b-instruct model consistently outperformed its non-instructed counterpart across all metrics. For instance, in the code generation task, the instruct model achieved a 52% higher average permissive accuracy (0.35 vs 0.23) and a 55% higher average exact accuracy (0.28 vs 0.18) compared to the non-instruct model.
2. Code generation vs. direct output: Interestingly, for both Llama3.1-8b and Llama3.1-8b-instruct, the direct output outperformed the code generation approach. This suggests that while code generation shows promise, there's still room for improvement in the quality and reliability of the generated code.
3. Comparison with Langchain: The facebook/opt-1.3b model using Langchain showed impressive performance on its test set, achieving 0.75 exact accuracy. However, it's important to note that this was on a much smaller sample size (150) compared to our main experiments (1000).

To further understand the challenges in the code generation approach, we analyzed the filtered sampling results, presented in Table 5.5.

These results reveal crucial insights:

1. Sample filtration: A significant number of samples were filtered out due to various issues, primarily code generation failures or execution errors.

Table 5.5: Filtered sampling in Code Generation

Model	Total samples selected	Remaining samples	Percentage of samples with execution errors	Utilization Rate
Llama3.1-8b (code-gen)	1000	779	31.58	53.3
Llama3.1-8b-instruct (code-gen)	1000	843	7.12	78.3

2. Improvement with instruction tuning: The Llama3.1-8b-instruct model showed a marked improvement in code generation reliability. It had fewer samples filtered out (157 vs 221) and a drastically lower percentage of execution errors (7.12% vs 31.58%) compared to the non-instruct model.
3. Utilization rate: The instruct model achieved a much higher utilization rate (78.3% vs 53.3%), indicating that it successfully generated and executed code for a larger proportion of the input samples.

5.5 Discussion

Our experiments reveal several important insights into improving number handling capabilities in language models for financial tasks:

1. **Model Architecture:** The Llama3-8B architecture demonstrates superior baseline performance compared to Llama2-7B, indicating that newer model iterations have improved inherent number handling abilities.
2. **Fine-tuning Effectiveness:** Fine-tuning proves to be a highly effective strategy, significantly boosting performance across all metrics. The fine-tuned Llama3-8B model achieved the best overall results, suggesting that combining advanced architectures with task-specific fine-tuning is a powerful approach.
3. **Few-shot Prompting:** While not outperforming fine-tuned models, few-shot prompting strategies show promise, especially when computational resources for full fine-tuning are limited. The similar performance of selective and random few-shot prompting indicates that the quantity of examples may be more crucial than their specific relevance in this domain.
4. **Example Selection:** Our results suggest that careful curation of examples for few-shot prompting may not always yield significantly better results than random selection. This finding could simplify the implementation of few-shot learning in practical applications.

5. **Tool-based Approaches:** The high performance of the Langchain-based tool on a subset of tasks highlights the potential of combining language models with specialized financial tools for certain applications.
6. **Potential of code generation approach:** Despite using models not primarily designed for code generation, we achieved promising results. This suggests that with more specialized models, the performance could be even better.
7. **Impact of instruction tuning:** The consistent outperformance of the instruct model demonstrates the value of instruction tuning in improving code generation capabilities and overall task performance.
8. **Challenges in code generation:** The high filtration rate and execution errors, especially in the non-instruct model, highlight the challenges in generating reliable, executable code for complex financial calculations.
9. **Trade-off between direct output and code generation:** While code generation offers the potential for more precise calculations, our results show that it currently underperforms compared to direct output. This suggests a need for further refinement of the code generation approach.
10. **Scalability and resource limitations:** Our experiments were limited to 8b models due to resource constraints. The significant improvement seen with the instruct model suggests that scaling to larger models (e.g., 70b or 405b) could yield even more impressive results, potentially establishing a new state-of-the-art for numerical precision in financial NLP tasks.

Chapter 6

Conclusions

6.1 Limitations

6.2 Future Work

Bibliography

- [1] Dogu Araci. Finbert: Financial sentiment analysis with pre-trained language models, 2019.
- [2] Preston Billion Polak, Joseph D Prusa, and Taghi M Khoshgoftaar. Low-shot learning and class imbalance: a survey. *Journal of Big Data*, 11(1):1, 2024.
- [3] 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.
- [4] Zhiyu Chen, Wenhui Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. Finqa: A dataset of numerical reasoning over financial data. *Proceedings of EMNLP 2021*, 2021.
- [5] Zhiyu Chen, Shiyang Li, Charese Smiley, Zhiqiang Ma, Sameena Shah, and William Yang Wang. Convfinqa: Exploring the chain of numerical reasoning in conversational finance question answering, 2022.
- [6] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- [8] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751*, 2019.
- [9] Mathav Raj J, Kushala VM, Harikrishna Warriar, and Yogesh Gupta. Fine tuning llm for enterprise: Practical guidelines and recommendations, 2024.

- [10] Suvarna Kadam and Vinay Vaidya. Review and analysis of zero, one and few shot learning approaches. In *Intelligent Systems Design and Applications: 18th International Conference on Intelligent Systems Design and Applications (ISDA 2018) held in Vellore, India, December 6-8, 2018, Volume 1*, pages 100–112. Springer, 2020.
- [11] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [12] Jean Lee, Nicholas Stevens, Soyeon Caren Han, and Minseok Song. A survey of large language models in finance (finllms), 2024.
- [13] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- [14] Junyi Li, Jie Chen, Ruiyang Ren, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. The dawn after the dark: An empirical study on factuality hallucination in large language models, 2024.
- [15] Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff Smith, Brian Vaughan, Pritam Damania, and Soumith Chintala. Pytorch distributed: Experiences on accelerating data parallel training, 2020.
- [16] Xianzhi Li, Samuel Chan, Xiaodan Zhu, Yulong Pei, Zhiqiang Ma, Xiaomo Liu, and Sameena Shah. Are chatgpt and gpt-4 general-purpose solvers for financial text analytics? a study on several typical tasks, 2023.
- [17] Yinheng Li, Shaofei Wang, Han Ding, and Hang Chen. Large language models in finance: A survey. In *Proceedings of the fourth ACM international conference on AI in finance*, pages 374–382, 2023.
- [18] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.
- [19] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.

- [20] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [21] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [22] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- [23] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools, 2023.
- [24] Omid Shahmirzadi, Adam Lugowski, and Kenneth Young. Text similarity in vector space models: A comparative study. In *2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA)*, pages 659–666, 2019.
- [25] Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*, 2022.
- [26] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023.
- [27] Neng Wang, Hongyang Yang, and Christina Dan Wang. Fingpt: Instruction tuning benchmark for open-source large language models in financial datasets, 2023.
- [28] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions, 2023.
- [29] Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners, 2022.

- [30] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [31] Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kambadur, David Rosenberg, and Gideon Mann. Bloomberggpt: A large language model for finance, 2023.
- [32] Haoyi Xiong, Jiang Bian, Sijia Yang, Xiaofei Zhang, Linghe Kong, and Daqing Zhang. Natural language based context modeling and reasoning for ubiquitous computing with large language models: A tutorial, 2023.
- [33] Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. TAT-QA: A question answering benchmark on a hybrid of tabular and textual content in finance. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3277–3287, Online, August 2021. Association for Computational Linguistics.

Appendix A

First appendix

A.1 Prompt template

PROMPT TEMPLATE FOR LLAMA2-7B-CHAT-HF

Prompt_template:

```
<s>[INST]
<<SYS>>
{system_prompt}
<</SYS>>
{user_message} [/INST]
```

system_prompt: "You are a helpful AI assistant specializing in financial analysis.
Answer the following question based on the given context."

user_message:

```
f"Context:\n{context}\n\nQuestion:_{question}"
```


PROMPT TEMPLATE FOR LLAMA3-8B**Prompt template:**

Below is an instruction that describes a task,
paired with an input that provides further context.
Write a response that appropriately completes the request.

Instruction:

{}

Input:

{}

Response:

{}

PROMPT TEMPLATE FOR LLAMA3-8B WITH FSP**Prompt template:**

```
You are a financial expert assistant. Answer the following
question
based on the given context. It is CRUCIAL that you use
step-by-step reasoning and provide detailed numerical
calculations where applicable.
Break down your answer into clear, numbered steps.
If asked to extract subject and object in relation, return
only the single
most relevant and applicable extraction. Use '\n' for line
breaks in your response.

Examples:
{examples}

Now, please answer the following question using the same
step-by-step
approach as demonstrated in the Examples:

Context: {context}
Question: {question}
Answer (including calculation steps and final answer, or
single
most relevant relation extraction, use '\n' for line
breaks): """
```

IMPROVED FEW SHOT PROMPT TEMPLATE FOR CODE GENERATION AND EXECUTION

Below is an instruction that describes a task, paired with examples and an input that provides further context. Learn from the examples and write a response that appropriately completes the request.

Instruction:

You are an AI assistant specialized in financial analysis. Your task is to generate Python code that answers questions based on given financial data and context. Learn from the examples provided and generate accurate, executable Python code that solves the given problem.

Important: When writing Python code, use actual newline characters to separate lines, not 'n' string literals. Each line of code should be on its own line in your response.

Examples:

{examples}

Input:

Context:

{context}

Question: {question}

Response:

Here's the Python code to answer the question:

```
``python
```

A.2 Examples

A.2.1 TAT-QA Dataset

Context: Actuarial assumptions. The Group's scheme liabilities are measured using the projected unit credit method with principal actuarial assumptions set out below:

- Figures represent a weighted average assumption of the individual schemes.
- The rate of increases in pensions in payment and deferred revaluation are dependent on the rate of inflation.

	2019 %	2018 %	2017 %
Table: Rate of inflation	2.9	2.9	3.0
Rate of increase in salaries	2.7	2.7	2.6
Discount rate	2.3	2.5	2.6

Question: What does the Weighted average actuarial assumptions consist of?

Answer: Rate of inflation, Rate of increase in salaries, Discount rate

A.2.2 CONVFINQA Dataset