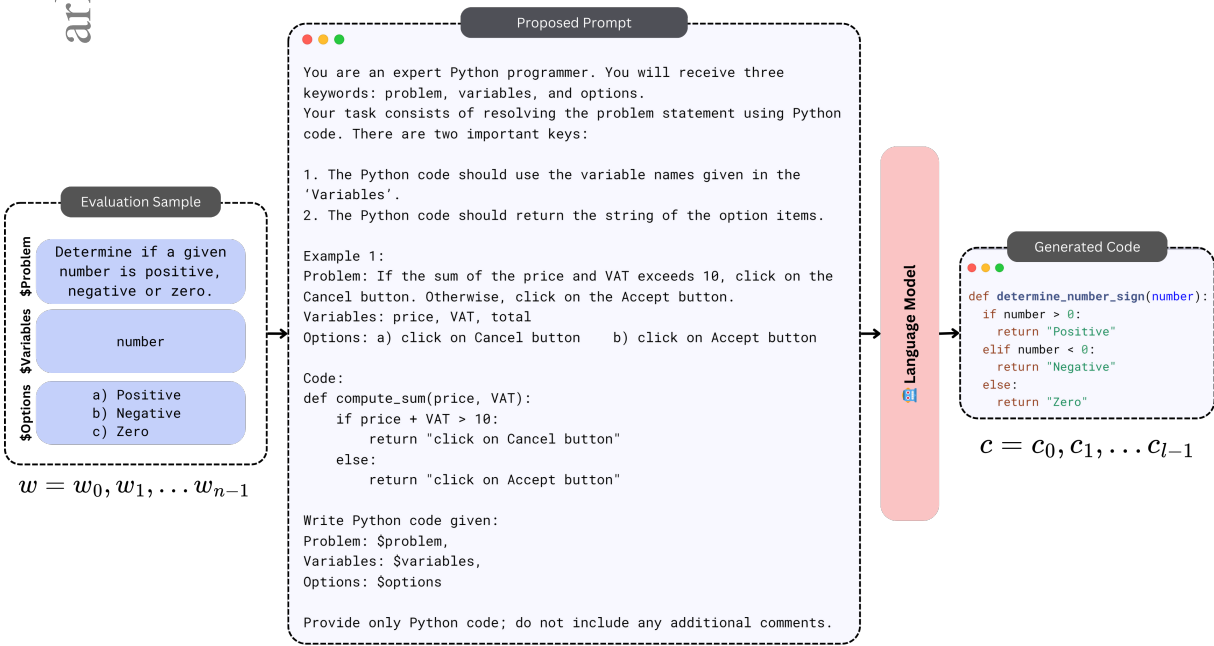


Graphical Abstract

Low-Cost Language Models: Survey and Performance Evaluation on Python Code Generation

Jessica López Espejel, Mahaman Sanoussi Yahaya Alassan, Merieme Bouhandi, Walid Dahhane, El Hassane Ettifouri



Highlights

Low-Cost Language Models: Survey and Performance Evaluation on Python Code Generation

Jessica López Espejel, Mahaman Sanoussi Yahaya Alassan, Merieme Bouhandi, Walid Dahhane, El Hassane Ettifouri

- Low-cost models achieve competitive scores compared to advanced chatbot models like ChatGPT-3.5 (OpenAI, 2022), ChatGPT-4 (OpenAI, 2023), and both Gemini 1.0 and 1.5 Pro (Team et al., 2023), which require significant computational resources.
- Some models, such as dolphin-2.6-mistral-7b (Ma et al., 2023b), excel in generating Python code but lag in meeting the required output formats. Often, when they generated correct code, these models do not return the wanted options from the input prompt.
- In contrast, models like Llama-3.1-8B-Instruct (Dubey et al., 2024) demonstrate consistent performance in understanding prompts and generating code, leading to superior results on our dataset. Similarly, they excel on HumanEval (Chen et al., 2021), and EvalPlus (Liu et al., 2023b) datasets.
- Remarkably, llama.cpp (Gerganov, 2023) project has achieved competitive outcomes using standard computers, a development that seemed very unlikely just a few months ago. This progress highlights the feasibility of running sophisticated language models on Central Processing Units (CPUs) today.

Low-Cost Language Models: Survey and Performance Evaluation on Python Code Generation

Jessica López Espejel^a, Mahaman Sanoussi Yahaya Alassan^a, Merieme Bouhandi^a, Walid Dahhane^a, El Hassane Ettifouri^a

^a*Novelis Research and Innovation Lab, 40 Av. des Terroirs de France, Paris, 75012, , France*

Abstract

Large Language Models (LLMs) have become a popular choice for many Natural Language Processing (NLP) tasks due to their versatility and ability to produce high-quality results. Specifically, they are increasingly used for automatic code generation to help developers tackle repetitive coding tasks. However, LLMs' substantial computational and memory requirements often make them inaccessible to users with limited resources. This paper focuses on very low-cost models which offer a more accessible alternative to resource-intensive LLMs. We notably: (1) propose a thorough semi-manual evaluation of their performance in generating Python code, (2) introduce a Chain-of-Thought (CoT) prompting strategy to improve model reasoning and code quality, and (3) propose a new dataset of 60 programming problems, with varied difficulty levels, designed to extend existing benchmarks like HumanEval and EvalPlus. Our findings show that some low-cost compatible models achieve competitive results compared to larger models like ChatGPT despite using significantly fewer resources. We will make our dataset and prompts publicly available to support further research.

Keywords: Python Code Generation, Natural Language Processing, Large

Preprint submitted to Engineering Applications of Artificial Intelligence August 30, 2024

1. Introduction

In recent years, automatic code generation has gained importance in the Natural Language Processing (NLP) domain because it automates the programming tasks, especially the most repetitive ones (Ahmad et al., 2021). Particularly, Python code generation has garnered extensive attention in programming language studies, being ranked among the top three most utilized languages globally (Coursera, 2023; Vailshery, 2024). In parallel, Large Language Models (LLMs) have made significant strides in NLP tasks (OpenAI, 2022, 2023; Manyika, 2023; Heidari et al., 2024), making them the go-to solution for many application domains including Python code generation. However, despite these advancements, the task of generating programming language code continues to pose greater complexity compared to standard natural language generation (López Espejel et al., 2023b). This complexity arises from the need for neural networks to accurately interpret natural language instructions to produce correct source code while also adhering to lexical, grammatical, and semantic constraints (Wang et al., 2021; Scholak et al., 2021; Liu et al., 2023c; Wong et al., 2023). Even a minor mistake, like missing a period or colon, can drastically change the meaning of the code, leading to incorrect model outputs.

While existing LLMs deliver exceptional performance, many are closed-source or face significant limitations due to their reliance on powerful Graphical Processing Units (GPUs) for inference. This reliance presents two main challenges: accessibility and cost, particularly impacting individuals and

smaller institutions. The reliance on high-performance GPUs, which are expensive and not readily available in standard computing setups, limits who can run these models efficiently. The financial burden of acquiring and maintaining such specialized hardware can be prohibitive. This diverts critical resources from other essential activities, leading to disparities in access to many Artificial Intelligence (AI) technologies. As a result, individuals and smaller organizations face significant barriers to fully participating in cutting-edge research and development, ultimately limiting the broader adoption and innovation of these powerful AI tools. In addition, even when LLMs are freely accessible, there is a risk of privacy leakage (Wu et al., 2023).

To address these issues, the research community has shifted its focus to models with fewer parameters. For instance, while large-scale models like GPT-3 (Generative Pre-trained Transformer) (Brown et al., 2020), PALM (Pathways Language Model) (Chowdhery et al., 2022), and ChatGPT-4 (OpenAI, 2023) boast 175B, 540B, and 1.7T parameters respectively, smaller models such as Mistral (Jiang et al., 2023) and LLaMa (Large Language Model Meta AI) (Touvron et al., 2023a,b; Dubey et al., 2024) demonstrate competitive or superior performance with only 7B parameters. Furthermore, there is growing interest in quantization techniques to speed up model inference. Popular techniques like GPTQ (Generative Pre-trained Transformers Quantization) (Frantar et al., 2023) and AWQ (Activation-aware Weight Quantization) (Lin et al., 2023) effectively compress LLMs while preserving much of their accuracy on downstream tasks. To leverage these techniques, libraries such as LlamaCPP (Gerganov, 2023) enable the execution of quantized LLMs on CPUs. LlamaCPP focuses on GGUF (GPT-Generated Unified Format)

files (Gerganov, 2024), an advanced binary format optimized for efficient storage and inference.

Motivated by the objective of demonstrating the potential and reliability of CPU-friendly models in natural language to Python code generation, this paper evaluates the performance of these open-source models. We then compare their performance to that of closed-source models to assess their capabilities and limitations. To do this, we propose a new dataset for Python code generation, made of 60 programming problems crafted using GPT-3.5-Turbo API (Application Programming Interface). The input text is structured around three keywords: *problem*, *variables*, and *options*. The *problem* describes the desired code function, the *variables* specify names of variables to use, and the *options* define return values for conditional statements. There are three levels of difficulty in the dataset (easy, intermediate, challenging). Problems range from basic tasks (e.g., even/odd number check) to complex ones requiring problem understanding, background knowledge (e.g., math), and code structure creation. For the sake of completeness, we also provide results on two state-of-the-art datasets: HumanEval (Chen et al., 2021), and EvalPlus (Liu et al., 2023b).

In addition to that, we improve the performance of these models by proposing a Chain-of-Thought (CoT) (Wei et al., 2022b) prompt that specifically guides the model into the problem solving. The prompt specifies the model’s role, the keywords it needs to focus on, and describes the target task. Then, a single example is provided to validate the process (one-shot inference). Finally, the prompt mentions that the goal is to have only one solving function, without any additional explanations or comments.

Our main contributions are:

- We propose a Python code generation dataset composed of 60 programming problems with three levels of difficulty. Each programming sample in our dataset is expressed as a *problem*, *variables*, and *options* (Section 4).
- We introduce a carefully-designed engineered prompt to enhance the performance of the evaluated models, drawing on insights from state-of-the-art practices (Section 5).
- We conducted an extensive semi-manual evaluation of the Python coding abilities of multiple CPU-friendly models from the state of the art. We used for evaluation our proposed dataset, HumanEval (Chen et al., 2021), and EvalPlus (Liu et al., 2023b). Here, we propose a quality-assessment prompt that is used to guide a GPT-3.5-Turbo (OpenAI, 2022) to evaluate the CPU-friendly models, then we manually correct the mistakes done by the GPT model (Sections 7 and 8).
- We will publicly share our dataset to facilitate future research.

The rest of the paper is organized as follows: in Section 2, we present the related work. Section 3 describes the problem formulation and the models used within the llama.cpp project. Section 4 presents the datasets used for the evaluation. Section 5 describes our proposed engineered prompts to enhance CPU models performance. Section 6 describes our evaluation methodology. We discuss the results in Sections 7 and 8, and finally conclude and provide some perspectives in Section 9.

2. Related Work

2.1. Python code generation

Python code generation has become a popular and challenging task in recent years. Chen et al. (2021) led a significant development in this area by introducing the Codex model and the HumanEval dataset. On the one hand, Codex is a fine-tuned GPT model containing up to 12 billion parameters. On the other hand, the HumanEval dataset was proposed as a benchmark for evaluating the proficiency of Python code generation models (Chowdhery et al., 2022; Li et al., 2023a; Rozière et al., 2024). Codex has sparked the emergence of other powerful code-generation Language Models (LMs). The study of these models can be divided into two groups: closed-source and open-source.

Closed-source models. One prominent example of closed-source models is Alphacode (Li et al., 2022), a model built upon the encoder-decoder transformer architecture (Vaswani et al., 2017), boasting a staggering 41 billion parameters. During training, the authors used the standard cross-entropy loss for next-token prediction in the decoder, as well as a Masked Language Modeling (MLM) loss (Devlin et al., 2019) for the encoder. In contrast to Alphacode, LaMDA (Thoppilan et al., 2022) relies on a decoder-only Transformer architecture (Vaswani et al., 2017), comprising 137 billion parameters that was pre-trained to predict the next token in a text corpus.

Gopher (Rae et al., 2021) and PALM (Pathways Language Model) (Chowdhery et al., 2022) share LaMDA’s architecture but with 280 billion and 540 billion parameters, respectively. Although Gopher used autoregressive Transformer architecture (Radford et al., 2018b), the authors made two significant

changes: they used RMSNorm (Zhang and Sennrich, 2019) instead of LayerNorm (Ba et al., 2016), and they opted for relative positional encodings (Dai et al., 2019) instead of absolute ones. On the other hand, PALM incorporates modifications such as SwiGLU activations (Shazeer, 2020) over Rectified Linear Units (Agarap, 2019), GeLU (Gated Linear Units) (Dauphin et al., 2017), or Swish activations (Ramachandran et al., 2017), and parallel layers (Wang and Komatsuzaki, 2021) within each Transformer block rather than the traditional serialized approach.

Hoffmann et al. (2022) found that for compute-optimal training, the model size and the number of training tokens should be scaled equally: for every doubling of model size, the number of training tokens should also be doubled. The study also revealed that for compute-efficient training, the model size should scale proportionally with the number of training tokens—doubling the model size necessitates doubling the training tokens. As a result, Chinchilla used the same compute budget as Gopher but with only 70B parameters, and four times more data. Chinchilla surpassed its predecessors like Gopher with 280 billion parameters, GPT-3 (Brown et al., 2020) with 175 billion parameters, and Megatron-Turing NLG (Smith et al., 2022) with 530 billion parameters, across various downstream evaluation tasks.

Chatbot models have exhibited outstanding performance in tasks involving Python code generation. Particularly, ChatGPT-3.5 (OpenAI, 2022) is built upon the GPT-3.5 model (Radford et al., 2018a,b; Ouyang et al., 2022) and marks a watershed moment in the evolution of AI, excelling in NLP tasks like code generation (Dong et al., 2023; Liu et al., 2023a,c), and text generation (Mulia et al., 2023; Lancaster, 2023). Afterward, OpenAI introduced

ChatGPT-4 (OpenAI, 2023), an upgraded version of the chatbot built upon the GPT-4 model. Although the exact number of parameters in ChatGPT-3.5 remains unknown, it was recently revealed that ChatGPT-4 contains a total of 1.7 trillion parameters (Leswing, 2024).

Following the release of ChatGPT-3.5 (OpenAI, 2022) and before ChatGPT-4 (OpenAI, 2023), Google introduced BARD (Manyika, 2023) as an enhanced iteration of the LaMDA model (Thoppilan et al., 2022). BARD enhances predictive abilities through the use of Reinforcement Learning from Human Feedback (RLHF) (Christiano et al., 2017). According to Ahmed et al. (2023) and López Espejel et al. (2023a), ChatGPT-4 demonstrates superior reasoning capabilities compared to BARD. More recently, Google unveiled Gemini (Team et al., 2023), an advanced series of models that further enhances the capabilities of BARD. Built upon Transformer decoders (Vaswani et al., 2017), Gemini models support a context length of 32,000. Notably, Gemini Ultra has become the first model to surpass human experts in the MMLU (Massive Multitask Language Understanding) benchmark (Hendrycks et al., 2021). Lastly, the Claude series represents a notable group of performance models. The latest version, Claude 3.5 Sonnet (Anthropic, 2024), surpasses Claude 3 Opus and other state-of-the-art models in reasoning, coding, question answering, and vision tasks.

Open-source models. Concurrently, open-source models emerged as an alternative to the private models. These open-source models have also showcased remarkable performance in Python code generation tasks. For instance, InCoder (Fried et al., 2023) and Starcoder (Li et al., 2023a) are models specifically trained on programming code. Diverging from other coding mod-

els (Brown et al., 2020; Chen et al., 2021), InCoder employs a causal masking objective during its training phase (Aghajanyan et al., 2022). This approach combines the strengths of both causal and masked language models (Devlin et al., 2019), enhancing its learning capabilities. Conversely, StarCoder adopts the FlashAttention mechanism (Dao et al., 2022) in its training process, which accelerates attention computations and reduces memory usage, leading to optimized performance.

Similar to InCoder, Code Llama-Python (Rozière et al., 2024) employs the causal masking technique to train infilling models. This model builds upon the LLaMA-2 architecture (Touvron et al., 2023b), incorporating fine-tuning adjustments. The introduced a Long Context Fine-Tuning (LCFT) stage (Rozière et al., 2024) allows models to handle sequences up to 16,384 tokens, which is a significant improvement from the typical 4,096 tokens in earlier LLaMA-2 stages. Despite having only 7B parameters, Code Llama-Python surpasses the 70B parameter Llama-2 in performance on the HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) datasets, showcasing its efficiency and effectiveness in handling large-scale coding challenges. Nevertheless, the performance of LLaMA-2 has been surpassed by both LLaMA-3 (AI@Meta, 2024) and the more recent LLaMA-3.1 (Dubey et al., 2024). Notably, LLaMA-3.1 is even more powerful than LLaMA-3 due to several enhancements: it is a multilingual model, benefits from improved data quality and quantity during both pre-training and post-training, and boasts a model size of 405B parameters compared to LLaMA-3’s maximum of 70B parameters.

In addition to the efforts to move towards open source language models,

the introduction of the Phi model family (Gunasekar et al., 2023; Li et al., 2023c) highlights a growing interest in creating models that compete in scale as well. The initial model, Phi-1 (Gunasekar et al., 2023), is a 1.3 billion parameter Transformer-based model designed for basic Python coding tasks, with an emphasis on textbook-quality training data. An enhanced version, Phi-1.5 (Li et al., 2023c), builds upon Phi-1 by incorporating a next-word prediction objective and utilizing a dataset of 30 billion tokens and training on 150 billion tokens. This version is an extension of Phi-2, a 2.7 billion parameter model trained with the combined data from Phi-1.5 and a new source of synthetic NLP texts and filtered websites. Therefore, in the next subsection, we will delve exclusively into small models.

2.2. Small Language Models

Since the creation of ChatGPT (OpenAI, 2022), numerous chatbots have emerged aiming to handle similar tasks while prioritizing efficient resource usage and delivering competitive performance. In this subsection, we will only discuss models with up to 13 billion parameters. One model enabling the development of small yet powerful models is LLaMA (versions 1 through 3.1) (Touvron et al., 2023a,b; AI@Meta, 2024; Dubey et al., 2024). Across its iterations, the LLaMA family includes variations of 7B, 8B, and 13B parameters, demonstrating comparable or superior performance to GPT-3 (175B parameters) (Brown et al., 2020) on most benchmarks. Inspired by LLaMA-1 and the success of instruction-following models like ChatGPT ¹, Claude ²,

¹<https://chat.openai.com/>

²<https://claude.ai>

and Microsoft BingChat ³, Stanford University launched Alpaca (Taori et al., 2023). Alpaca is a 7-billion-parameter model fine-tuned from LLaMa-1-7B using 52K instruction-following demonstrations. Despite being significantly smaller, Alpaca exhibits similar behavior to OpenAI’s text-davinci-003 . It’s worth highlighting that Vicuna (Chiang et al., 2023) surged by leveraging the strengths of both LLaMa and Alpaca. However, it slightly exceeds Alpaca in parameter count, with its 13 billion parameters.

Even though the models mentioned earlier feature up to 13 billion parameters, a groundbreaking model has emerged: Mistral (Jiang et al., 2023). Mistral is a 7-billion-parameter LM that enhances inference speed through the use of Grouped-Query Attention (GQA) (Ainslie et al., 2023), while simultaneously cost-effectively managing sequences of arbitrary length via Sliding Window Attention (SWA) (Beltagy et al., 2020). Mistral outperforms the performance of LLaMA-2 13B (Touvron et al., 2023b) on all assessed benchmarks, including code generation tasks, as well as the best among released 34B LLaMA-2 variants. Building on Mistral’s success, the improved Mixtral 8x7B (Jiang et al., 2024) utilizes a Sparse Mixture of Experts (SMoE) approach. This improvement not only enables Mixtral to outperform or match the capabilities of the 70B-parameter Llama-2 and GPT-3.5 model across all benchmarks but also to maintains the same foundational architecture as the Mistral 7B. The key difference lies in its structure: each layer in Mixtral consists of eight feedforward blocks, or “experts”, which enhance its processing efficiency and model response capability.

Zephyr (Tunstall et al., 2023) is an upgraded version of Mistral-7B, re-

³<https://www.bing.com/chat>

fined through fine-tuning with the $\sim 200,000$ samples UltraChat dataset (Ding et al., 2023; Tunstall et al., 2023). It has been further aligned using the $\sim 64,000$ samples UltraFeedback dataset (Cui et al., 2023) and optimized through the Direct Preference Optimization (DPO) algorithm (Rafailov et al., 2023). This model outperforms LLaMA2 70B on the MT-Bench benchmark (Zheng et al., 2023). However, the MiniCPM model (Inc. and TsinghuaNLP, 2024) notably surpassed Zephyr-7B-alpha in MT-Bench performance. Just like Zephyr-7B, it employs supervised fine-tuning and the DPO algorithm, but is more compact, with only 2 billion parameters.

Beyond Zephyr, several significant models were refined using Mistral-7B as their foundation, including mistral-7b-openorca (Mukherjee et al., 2023), dolphin-2.6-mistral (Ma et al., 2023b), and openhermes-2.5-mistral (Teknium, 2023). The mistral-7b-openorca model leverages the OpenOrca dataset, focusing on the interaction triad of “System message, User query, LFM response”. In contrast, dolphin-2.6-mistral enriches its training with three diverse datasets: Capybara (Daniele and Suphavadeeprasit, 2023), Magicoder-Evol-Instruct-110K (Wei et al., 2023), and Dolphin (Hartford, 2023). Openhermes-2.5-mistral, while also building upon Mistral-7B set itself apart by training on a variety of code benchmarks, including TruthfulQA (Lin et al., 2021), AGIEval (Zhong et al., 2023), and the GPT4All suite (Anand et al., 2023).

2.3. Quantization for model acceleration

In light of the success of LLMs in addressing tasks such as question answering (Omar et al., 2023; Zhuang et al., 2024), and conversational chatbots (Chiang et al., 2023; Taori et al., 2023; Ramírez et al., 2024), considerable efforts have been directed towards reducing deployment costs. These

efforts involve optimizing models for CPU efficiency using compression methods like pruning (Frantar and Alistarh, 2023; Ma et al., 2023a; Dery et al., 2024) and quantization (Li et al., 2023b; Kim et al., 2024; Huang et al., 2024), with the primary aim of decreasing the number of parameters.

Regarding quantization, this technique reduces model precision by retaining only the most significant bits, resulting in lower-bit models (1-2 or n bits). Quantization comes in various forms, including post-training quantization (PTQ) (Xiao et al., 2023; Huang et al., 2024), Quantization-Aware Training (QAT) (Chen et al., 2024; Liu et al., 2023d) and weight-only (Frantar et al., 2023). This paper focuses on weight-only quantized models, which address inference time challenges due to memory constraints related to weight parameters (Park et al., 2024). One notable early work in this area is GPTQ (Frantar et al., 2023). GPTQ (Gradient Post-Training Quantization) uses second-order information for efficient weight quantization in a single step. AWQ (Activation-aware Weight Quantization) (Lin et al., 2023) focuses on activation distributions rather than weight importance and does not require backpropagation or reconstruction. The QuIP (Quantization with Incoherence Processing) technique (Chee et al., 2023) enhances quantization by promoting incoherence in weight and Hessian matrices using random orthogonal matrices. In contrast, SpQR (Sparse-Quantized Representation) (Dettmers et al., 2023) addresses quantization errors by storing outlier weights in higher precision while compressing other weights to 3-4 bits. Collectively, these methods offer sophisticated solutions for weight-only quantization, significantly improving model efficiency and accuracy.

Building on these advancements in quantization, we utilized the open-

source LlamaCPP project (Gerganov, 2023) to conduct experiments on generating Python code from natural language using quantized models. This framework supports quantization for various models, aligning with our focus. Detailed information about the models used in our experiments is provided in Subsection 3.2.

3. Python Code Generation using Low-Cost CPU Models

3.1. Problem Definition

Given a natural language input text $w = \{w_0, w_1, \dots, w_{n-1}\}$ of length $|w| = n$, a Python code generation model seeks to generate an equivalent Python source code sequence $c = \{c_0, c_1, \dots, c_{l-1}\}$ of length $|c| = l$, where w_i and c_i are input and output tokens respectively. For the sake of clarity, we consider “CPU-friendly” any model that can be run on CPU for inference (not necessarily during training).

We provide an overview of our system in Figure 1 . As mentioned in the introduction, the input text (purple) comprises three key elements: a problem statement, variable names, and a set of options. We feed this input into our Chain-of-Thought prompting module (grey) to generate an efficient prompt that guides the model in its generation. Then, we provide this prompt to the CPU-friendly Language Model (yellow), and the latter finally provides the output Python code (blue). More details about the dataset and the Chain-of-Thought prompting module are provided in Sections 4, and 5, respectively.

3.2. Evaluation Protocol within llama.cpp Project

We performed an evaluation using GGUF format model files sourced from the llama.cpp (Gerganov, 2023) project, which is a C++ library primarily

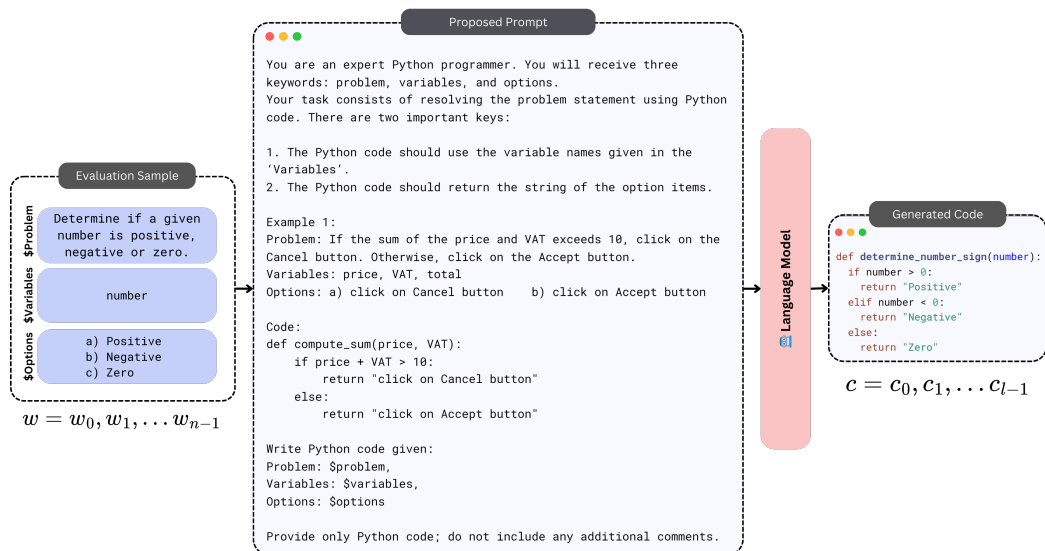


Figure 1: Overview of Python code generation with our system. The input consist of a *problem*, *variables*, and *options*. We feed it to the language model along with an engineered prompt to guide the code generation. The language model generates the equivalent Python code.

optimized for CPU usage. This library effectively implements several LLM architectures, including noteworthy models such as LLaMA (Touvron et al., 2023a,b; Dubey et al., 2024), Mistral (Jiang et al., 2023), Mixtral (Jiang et al., 2024), StarCoder (Li et al., 2023a), and Gemma (Team et al., 2024). These models are available in various quantized versions, ranging from 2 to 8 bits, and are designed to support both CPU and GPU inference, demonstrating the versatility of llama.cpp in different computing environments. The lightweight design of llama.cpp ensures portability and speed, enabling rapid responses across a wide range of devices. Its customizable low-level features empower developers to deliver effective real-time coding assistance.

In the many model variants we evaluate (the full list of models is shown

in Figure 2), the name is noted as such *model-name.Q_j-K-Size*, where:

- *model-name* is the model name and can be any of the following: llama-2-coder-7b, Llama-3.1-8B-Instruct, phi-2, mistral-7b-instruct-v0.2, mistral-7b-openorca, zephyr-7b-beta, dolphin-2.6-mistral-7b, openhermes-2.5-mistral-7b, and MiniCPM-2B-dpo-bf16.
- *Q_j-K-Size* refers to a specific type of quantization method. *j* is the bit-width used, *K* refers to the use of K-means clustering in the quantization, *Size* is the size of the model after quantization, with *S*, *M*, *L* indicating Small, Medium, and Large, respectively.

Hardware specifications. To conduct our experiments on the CPU, we used a machine equipped with the following hardware: Processor: 11th Gen Intel(R) Core(TM) i7-11850H @ 2.50GHz, RAM: 32.0 GB, Operating System: 64-bit, x64-based processor.

Software specifications. We primarily used the following libraries: Transformers (version 4.43.1), Python (version 3.10), BitsandBytes (version 0.43.2), and Accelerate (version 0.21.0). Moreover, for the experiments in GPU, we use specifically A100 on Google Colab. Additionally, all GPU-based experiments were conducted using the A100 GPU on Google Colab.

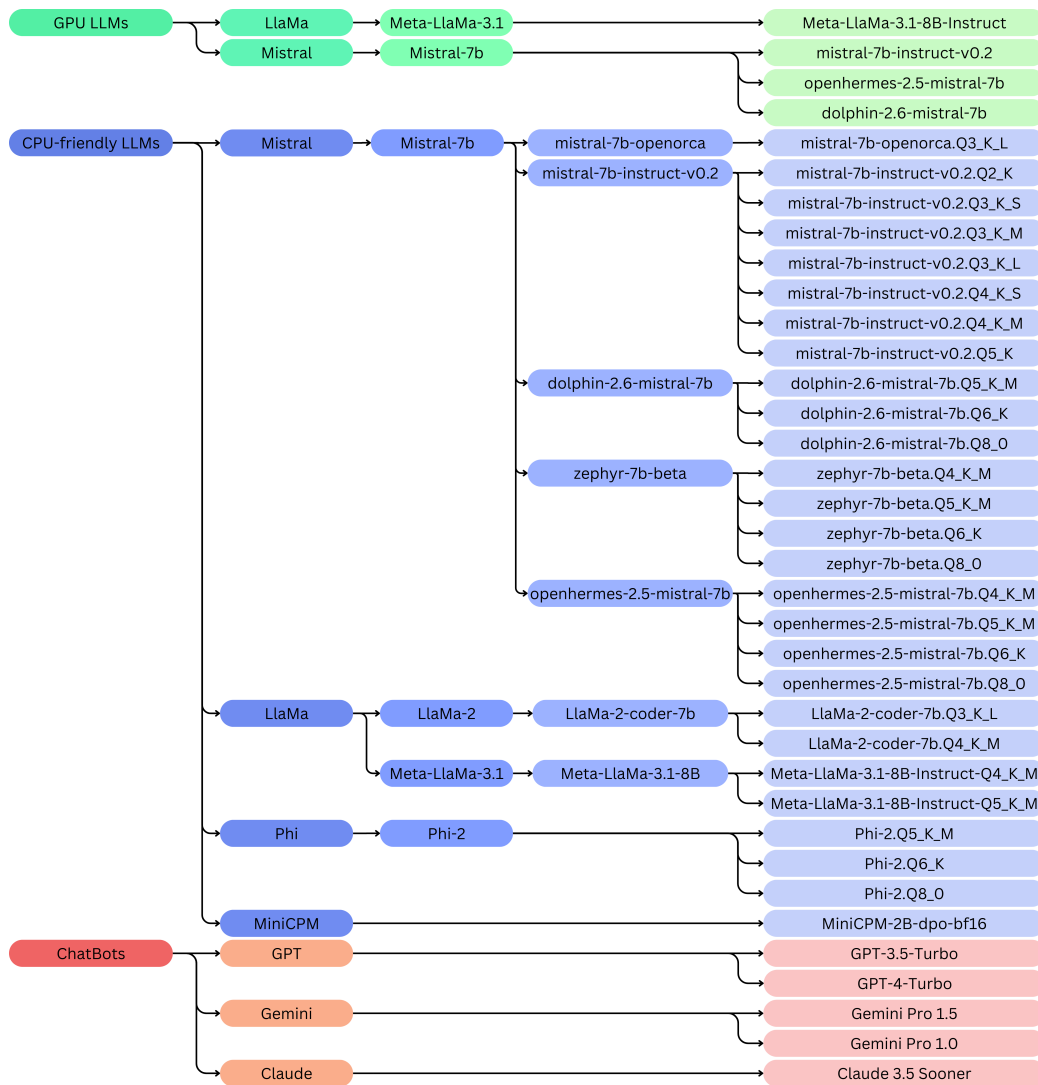


Figure 2: Hierarchical organization of all models (CPU, GPU, and Chatbots) examined in the experiments conducted in this paper.

4. Datasets

We evaluate the models using three Python datasets: our dataset, HumanEval (Chen et al., 2021), and EvalPlus (Liu et al., 2023b).

4.1. Proposed dataset

To assess the functional accuracy of code generated by various models, we designed a dataset comprising 60 crafted programming problems, each consisting of a problem statement, variable names, and a set of options. The *problem statement* describes the conditions and constraints of the user’s requirements. The *variables* keyword, as implied by their name, defines the names of variables to be utilized in the Python code. Lastly, the *options* represent the potential outcomes for each condition outlined in the problem statement. All our dataset samples are generated with the help of GPT-3.5-Turbo API, and then carefully reviewed by hand afterward.

The samples are categorized into three levels of difficulty:

- the initial 20 samples represent an easy level ;
- the subsequent 20 samples constitute an intermediate level ;
- the final 20 samples present a challenging level

The samples include mathematical, logic, and reasoning problems. We chose the topics based on popular coding platforms like LeetCode and Geeks-forGeeks, ensuring they reflect the types of challenges commonly encountered in real-world coding tasks. At the easy level, tasks involve basic operations such as identifying the type of an angle based on its degree or determining whether one string is an anagram of another. Intermediate problems typically focus on matrices, requiring a solid understanding of concepts like diagonal and orthogonal matrices. At the most challenging level, some problems that involve graphs, trees, and other complex structures. We provide an example of each difficulty level in Figure 3, and an example output for the easy example in Figure 4.

| |
|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <ul style="list-style-type: none"> • Easy: <hr/> Problem: Determine if a given number is positive, negative, or zero. Variables: number Options: a) Positive b) Negative c) Zero <hr/> |
| <ul style="list-style-type: none"> • Intermediate: <hr/> Problem: Check if a given sentence is a pangram. Variables: input_sentence Options: a) Pangram b) Not a pangram <hr/> |
| <ul style="list-style-type: none"> • Challenging: <hr/> Problem: Check if a directed graph is acyclic. Variables: directed_graph Options: a) Acyclic graph b) Cyclic graph <hr/> |

Figure 3: Examples of different programming languages.

```
def determine_number_sign(number):
    if number > 0:
        return "Positive"
    elif number < 0:
        return "Negative"
    else:
        return "Zero"
```

Figure 4: Example solution for the easy problem from Figure 3.

4.2. HumanEval and EvalPlus datasets

HumanEval (Chen et al., 2021) is a hand-written dataset, developed by OpenAI. It consists of 164 Python problems. Each prompt input contains a function signature and a docstring. Each problem in this dataset is associated with a canonical solution and an average of 7.7 unit tests. The dataset aims to evaluate the models' ability to accurately complete functions based on their

signatures and docstrings. The performance of the models is determined by how successfully they pass the given unit tests, ensuring a strict evaluation of the code’s functional correctness.

EvalPlus (Liu et al., 2023b) is an improved version of the original HumanEval dataset. The authors identified that HumanEval did not encompass all potential scenarios for testing the functional correctness of LLM-synthesized code. To address this issue, they significantly expanded the dataset by increasing the number of test cases. EvalPlus comprises an average of 764.1 tests per problem, providing a more robust and comprehensive framework for assessing code correctness.

5. Prompt Engineering to Improve Generation

To design our prompt more thoughtfully, we have integrated insights from the research of Yu et al. (2023). In their work, they recommend incorporating a Chain-of-Thought (CoT) approach (Wei et al., 2022b) to steer the model in generating suitable responses for specific tasks. They highlight that effective CoT prompts generally include either CoT demonstrations or textual instructions. To ensure the effectiveness of our prompt, we include both an example input with its corresponding output and textual instructions to guide the language model. Furthermore, drawing inspiration from the findings of Chen (2023), we opted to use a single example in our prompt. Their research demonstrates that LLMs can deliver robust performance even with a minimal, 1-shot demonstration, showcasing the efficiency of this streamlined approach. The full approach to design the prompt is as follows:

- We designate a specific role for the chatbot: *You are an expert python*

programmer.

- We identify the keywords the chatbot should pay attention to: *You will receive three keywords: problem, variables, and options.*
- We explicitly describe the task that needs to be performed: *Your task consists of resolving the problem statement using python code.*
- We used a single example to show the model how to effectively execute the task. Chen (2023) suggests that providing just one example is enough for the models, and also helps in keeping the prompt short and concise.
- At the end of the prompt, we reiterate the specific format in which we expect the model to return its output. In our case, we specify that we require only a function code to be returned.

Figure 1 displays the complete prompt we propose as the input for each evaluated model. Although this is a general prompt, it was adapted for each model using the corresponding tokenizer via the `encode_chat_completion` function from Huggingface ⁴ library. This function ensures the correct template with the appropriate special tokens and roles for each model. For instance, Mistral-7B-Instruct-v0.2 utilizes the roles “user” and “assistant”, whereas Meta-Llama-3.1-8B-Instruct employs three roles, “system”, “user” and “assistant”. Consequently, the prompt was accurately adapted to fit the requirements of each model. For further details, please refer to Appendix A.

⁴<https://huggingface.co/>

6. Evaluation methodology

In this section, we provide a comprehensive breakdown of the metrics and criteria used to evaluate our dataset, HumanEval (Chen et al., 2021), and EvalPlus (Liu et al., 2023b).

6.1. Proposed dataset

To assess the performance of CPU-friendly models in generating Python code, we employ semi-manual evaluation. We provide the evaluation prompt (Figure 5) augmented with the answer generated by the model to GPT-3.5-Turbo (OpenAI, 2022) through its API. The latter returns a correctness score of 0 (wrong answer), 0.5 (passable answer), or 1 (correct answer).

```
Please compare both codes, referred to as Code 1 and Code 2.

They aim to yield identical results, disregarding the implementation
↪ details. If both codes successfully produce the same results, return
↪ 1. If the second code maintains the logic but falls short of
↪ replicating the results of Code 1, return 0.5. If the codes are
↪ entirely different, return 0.

The codes are supposed to describe the code's objective or
intention.
```

Figure 5: Proposed Prompt to evaluate our results

We proposed an approach focuses on thoroughly analyzing code errors without being overly rigid in the assessment. Instead of categorizing outputs as simply correct or incorrect (1 or 0), it accounts for the nuances and subtleties in the generated code. Specifically, it assigns a score of 0.5 to code

that, while not yielding the exact results, still adheres to the underlying logic of the problem. For a clearer illustration, please refer to Table 1.

Table 1 presents examples of how GPT-3.5-Turbo scores generated code. In the first example, the code is correct, earning a score of 1. In the second, GPT-3.5-Turbo assigns a 0.5 because the quantized model failed to convert strings to lowercase, potentially leading to incorrect results. The third example receives a score of 0 when the code incorrectly attempts to determine if a number is even instead of checking for primality.

| Problem Statement | Canonical Solution | Response from the Quantized LLM | GPT Score |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------|-----------|
| <p>Problem: Identify if a given number is a multiple of 3</p> <p>Variables: n</p> <p>Options: a) Multiple of 3 b) Not a multiple c) Other</p> | <pre>def multiple(n): assert → isinstance(n, → int), "c) → Other" if n % 3 == 0: return "a) → Multiple of → 3" else: return "b) Not a → multiple"</pre> | <pre>def multiple(n): assert → isinstance(n, → int), "Other" if n % 3 == 0: return "a)" else: return "b)"</pre> | 1 |

| Problem Statement | Canonical Solution | Response from the Quantized LLM | GPT Score |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------|-----------|
| <p>Problem: Check if a given string is an anagram of another string</p> <p>Variables: s1, s2</p> | <pre>def anagram(s1, ↪ s2): s1 = s1.lower() s2 = s2.lower() if sorted(s1) == ↪ sorted(s2):</pre> | <pre>def anagram(s1, ↪ s2): if sorted(s1) == ↪ sorted(s2): return "Anagram"</pre> | 0.5 |
| <p>Options: a) Anagram b) Not an anagram c) Other</p> <p>Problem: Determine if a given number is prime.</p> <p>Variables: n</p> <p>Options: a) Prime number b) Not a prime number c) Other</p> | <pre>def is_prime(n): if n <= 1: return "Not a ↪ prime ↪ number" for i in range(2, ↪ int(n**0.5) + ↪ 1): if n % i == 0: return "Not a ↪ prime ↪ number" return "Other"</pre> | <pre>def is_prime(n): if n <= 1: return "Not a ↪ prime ↪ number" if n % 2 == 0: return "Not a ↪ prime ↪ number"</pre> | 0 |

Table 1: Examples of scores assigned by GPT-3.5-Turbo using the prompt detailed in Table 5

We intervened manually twice. First, we double-checked the scores assigned by GPT-3.5-Turbo to code generated by the quantized models using test cases from our dataset. As described in the dataset section, these tests ensure the code uses the specified variable names, returns the expected options (e.g., "a", "b", "c" or the full text), and passes all performance criteria. For example, Table 2 shows a case where GPT-3.5-Turbo initially scored the code as 1. However, after applying our test cases, the score was reduced to 0.5 due to incorrect variable names (as required in Table 1).

| Canonical Solution | GPT Score | Response from the Quantized LLM | Test Score |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------|
| <pre>def anagram(s1, s2): s1 = s1.lower() s2 = s2.lower() if sorted(s1) == ↪ sorted(s2): return "Anagram" else: return "Not an ↪ anagram"</pre> | 1 | <pre>def anagram(str1, str2): str1 = str1.lower() str2 = str2.lower() if sorted(str1) == ↪ sorted(str2): return "Anagram" else: return "Not an ↪ anagram"</pre> | 0.5 |

Table 2: An example of penalization in the code: GPT-3.5 initially scored it as 1 (correct code), but after applying our test cases, the score was reduced to 0.5 due to incorrect use of variable names.

At this stage of our methodology, we ensure that correct answers are accurately identified. However, we manually review the code assigned scores of 0 and 0.5. Our experience shows that the model occasionally assigns a score of 0 when a 0.5 is more appropriate. This occurs because the model lacks the nuanced understanding that some code, despite failing tests, still maintains correct logic and warrants a score of 0.5. Table 3 illustrates an

example of this manual adjustment to ensure fair scoring of each generated code.

| Canonical Solution | Response from the Quantized Model | GPT Score | Manual Score |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------|--------------|
| <pre>def is_isogram(word): if len(set(word.lower() ↪)) == ↪ len(word.lower()): return "a) Isogram" else: return "b) Not an ↪ isogram"</pre> | <pre>def isogram(word): word low = word.lower() if len(set(word low)) ↪ == len(word): return "a" else: return "Other"</pre> | 0 | 0.5 |

Table 3: This example demonstrates a case where GPT-3.5-Turbo assigned a score of 0 to code generated by the quantized model. Although the code fails to produce the same results as the canonical solution and returns different responses ("Other" instead of "b" or "b) Not an isogram"), it still maintains the logic of the problem. The code attempts to solve the problem but ultimately fails in its execution.

Below, we outline the criteria for each score in more depth.

- **Score = 0 (wrong)** - A score of zero is assigned when the generated code demonstrates no comprehension of the input prompt, indicated by a lack of correlation between the function name, its algorithm, and the requirements specified in the prompt. A score of zero is also given when the function body is empty.
- **Score = 0.5 (passable)** - A score of 0.5 is awarded when the generated code understands the input prompt and attempts to address the problem but fails to deliver the expected response due to syntax issues or logical oversights. For instance, omitting a base case in recursive problems can lead to inefficiency or infinite loops, although the remainder of the algorithm suggests an understanding of the task.

Additionally, the code might technically be correct but returns an unexpected output. For example, when asked whether a number is prime with options “a) Prime number, or b) Not a prime number”, the model might simply return True or False, indicating a lack of deep understanding of the prompt.

- **Score = 1 (correct)** - A score of 1 is assigned to the code generation when the model’s output matches the gold standard code in functionality, despite possible syntactical differences. It’s worth noting that the evaluation criteria are overall not rigid. For instance, if the correct answer is “a) Prime number”, but the evaluated model omits the “a)” and simply returns “Prime number” (or vice versa), the generated code is still considered correct.

Our evaluation approach aims to be more flexible. This is achieved by focusing not only on whether the code passes unit tests, but also on understanding the challenges that different models face when generating code across various levels of difficulty.

6.2. *HumanEval and EvalPlus datasets*

To evaluate the performance of the quantized models in the Python code completion task, we used the HumanEval and EvalPlus datasets. In these datasets, the task involves completing a Python function based on the provided function header and corresponding docstring. The prompt used for this task is shown in Figure 6.

As detailed in Section 5, we adapt the template based on the model’s specific requirements. For models with both system and user prompts, we use

```

You are an expert Python programmer. You will be provided with Python
↪ code that requires your completion.

$function_header_docstring

```

Figure 6: Prompt used to generate Python code for the databases: HumanEval and EvalPlus

the following approach: the system prompt is “You are an expert Python programmer. You will be provided with Python code that requires your completion,” while the user prompt consists of the function header and corresponding docstring.

After receiving the LM’s response, we perform postprocessing to extract only the generated Python code, as some LMs may include unnecessary text. We then evaluate this code using the $\text{pass}@k$ metric (Kulal et al., 2019). Initially, it was introduced to assess functional correctness, where k denotes the number of code samples generated for each problem. A problem is considered solved if any of these k samples successfully pass the unit tests. The resulting score indicates the proportion of problems solved correctly.

A slight modification to this metric was proposed by Chen et al. (2021). They suggest generating n total samples per task, then counting the correct samples $c \leq n$ that pass the unit tests, and computing the unbiased estimator using Equation 1.

$$\text{pass}@k = \mathbb{E}_{\text{Problems}} \left[1 - \frac{\binom{n-c}{k}}{\binom{n}{k}} \right] \quad (1)$$

However, this estimator often yields large values. Thus, in practice, the score is computed using $1 - (1 - \hat{p})^k$, where \hat{p} represents the empirical estimate

of pass@1.

Therefore, we evaluated the results using the approach proposed by Chen et al. (2021). Specifically, we employ the pass@1 metric, which involves generating a single sample per problem, using the implementation from the Hugging Face library.

7. Results and Discussion: Our dataset

Table 4 presents the results of our experiments on our dataset. It is worth noting that, for comparison, we included the performance of quantized models alongside non-quantized models that require a GPU, as well as the scores for renowned chatbots such as ChatGPT-3.5, ChatGPT-4, and both ChatGemini 1.0 and 1.5-Pro. Note that, in Table 4, $Correct + Passable + Wrong = 100\%$.

ChatGPT-3.5 vs. ChatGPT-4 vs. Gemini 1.0-1.5 Pro. Our findings show that for our specific task and input problem formulation, Gemini 1.5 Pro achieved the highest score in the number of correct responses, outperforming ChatGPT-3.5, ChatGPT-4, and Gemini 1.0 by 1.67%, 8.33%, and 15%, respectively. Furthermore, consistently generated optimized and well-documented algorithms, while also passing unit tests.

It’s worth noting that ChatGPT-3.5 outperformed ChatGPT-4 by 6.66% in terms of correct answers, despite ChatGPT-4’s established superiority in tasks like reasoning (López Espejel et al., 2023a) and solving mathematical and logic problems (Plevris et al., 2023). In addition, ChatGPT-3.5 produced 0% wrong answers, compared to ChatGPT-4’s 1.67%. While ChatGPT-4’s algorithms were more optimized and concise, ChatGPT-3.5 tended to be

| | Models | Size (GB) ↓ | Required RAM (GB) ↓ | Inference Time (ms) ↓ | Correct (%) ↑ | Passable (%) ↑ | Wrong (%) ↓ | Correct+Passable (%) ↑ |
|-----|-----------------------------------|-------------|---------------------|-----------------------|---------------|----------------|--------------|------------------------|
| GPU | ChatGemini 1.5 Pro | × | × | × | 95.00 | 3.33 | 1.66 | 96.66 |
| | ChatGPT-3.5 | × | × | × | 93.33 | 6.67 | 0 | 96.67 |
| | ChatGPT-4 | × | × | × | 86.67 | 11.67 | 1.67 | 92.50 |
| | ChatGemini 1.0 | × | × | × | 80.00 | 18.33 | 1.67 | 89.17 |
| GPU | Meta-Llama-3.1-8B-Instruct | 16.5 | 29.2 | 30.12 | 91.66 | 5.00 | 3.33 | 94.16 |
| | Mistral-7B-Instruct-v0.2 | 26.98 | 27.7 | 23.08 | 90.00 | 5.00 | 5.00 | 92.50 |
| | dolphin-2.6-mistral-7b | 26.98 | 28.7 | 13.16 | 60.00 | 25.00 | 15.00 | 72.50 |
| | OpenHermes-2.5-Mistral-7B | 26.98 | 28.7 | 3.33 | 66.66 | 18.33 | 15.00 | 75.82 |
| CPU | Meta-Llama-3.1-8B-Instruct-Q4_K_M | 4.92 | 7.51 | 643.65 | 86.66 | 10.00 | 3.33 | 91.66 |
| | Meta-Llama-3.1-8B-Instruct-Q5_K_M | 5.73 | 8.14 | 662.32 | 88.33 | 8.33 | 3.33 | 92.49 |
| | llama-2-coder-7b.Q3_K_M | 3.30 | 5.80 | 215.51 | 16.67 | 16.67 | 66.67 | 25.00 |
| | llama-2-coder-7b.Q3_K_L | 3.60 | 6.10 | 268.36 | 11.67 | 20.00 | 68.33 | 21.67 |
| | phi-2.Q5_K_M | 2.07 | 4.57 | 147.48 | 28.33 | 26.67 | 45.00 | 41.67 |
| | phi-2.Q6_K | 2.29 | 4.79 | 145.41 | 26.67 | 41.67 | 31.67 | 47.50 |
| | phi-2.Q8_0 | 2.96 | 5.46 | 182.52 | 30.00 | 41.67 | 28.33 | 50.83 |
| | mistral-7b-instruct-v0.2.Q2_K | 3.08 | 5.58 | 209.85 | 33.33 | 40.00 | 26.67 | 53.33 |
| | mistral-7b-instruct-v0.2.Q3_K_S | 3.16 | 5.66 | 200.09 | 56.67 | 20.00 | 23.33 | 66.67 |
| | mistral-7b-instruct-v0.2.Q3_K_M | 3.52 | 6.02 | 250.16 | 55.00 | 35.00 | 10.00 | 72.50 |
| | mistral-7b-instruct-v0.2.Q3_K_L | 3.82 | 6.32 | 261.95 | 50.00 | 21.67 | 28.33 | 60.83 |
| | mistral-7b-instruct-v0.2.Q4_K_S | 4.14 | 6.64 | 278.04 | 73.33 | 16.67 | 10.00 | 81.67 |
| | mistral-7b-instruct-v0.2.Q4_K_M | 4.37 | 6.87 | 304.72 | 86.67 | 8.33 | 5.00 | 90.83 |
| | mistral-7b-instruct-v0.2.Q5_K_M | 5.13 | 7.63 | 330.64 | 58.33 | 20.00 | 21.67 | 68.33 |
| | mistral-7b-openorca.Q3_K_L | 3.82 | 6.32 | 252.96 | 28.33 | 40.00 | 31.67 | 48.33 |
| | zephyr-7b-beta.Q4_K_M | 4.37 | 6.87 | 240.36 | 45.00 | 38.33 | 16.67 | 64.17 |
| | zephyr-7b-beta.Q5_K_M | 5.13 | 7.63 | 281.22 | 41.67 | 33.33 | 25.00 | 58.33 |
| | zephyr-7b-beta.Q6_K | 5.94 | 8.44 | 329.11 | 36.67 | 46.67 | 16.67 | 60.00 |
| | zephyr-7b-beta.Q8_0 | 7.70 | 10.20 | 416.54 | 36.67 | 53.33 | 10.00 | 63.33 |
| | dolphin-2.6-mistral-7b.Q5_K_M | 5.13 | 7.63 | 358.30 | 48.33 | 36.67 | 15.00 | 66.67 |
| | dolphin-2.6-mistral-7b.Q6_K | 5.94 | 8.44 | 393.90 | 30.00 | 55.00 | 15.00 | 57.50 |
| | dolphin-2.6-mistral-7b.Q8_0 | 7.70 | 10.20 | 422.94 | 50.00 | 40.00 | 10.00 | 70.00 |
| | openhermes-2.5-mistral-7b.Q4_K_M | 4.37 | 6.87 | 260.72 | 53.33 | 21.67 | 25.00 | 64.17 |
| | openhermes-2.5-mistral-7b.Q5_K_M | 5.13 | 7.63 | 298.93 | 55.00 | 25.00 | 20.00 | 67.50 |
| | openhermes-2.5-mistral-7b.Q6_K | 5.94 | 8.44 | 325.59 | 58.33 | 25.00 | 16.67 | 70.83 |
| | openhermes-2.5-mistral-7b.Q8_0 | 7.7 | 10.20 | 413.84 | 58.33 | 21.67 | 20.00 | 69.17 |
| | MiniCPM-2B-dpo-bf16 | 5.45 | 10.9 | 510.25 | 50.00 | 36.67 | 13.33 | 68.33 |

Table 4: Evaluation on our dataset. Q_j : quantization using j bit width, K: the use of k-means clustering in the quantization, S , M , L : Small, Medium, and Large model size after quantization. The best results of each category are in bold.

more verbose in its outputs. The confidential nature of these models’ development processes and OpenAI’s versioning strategy makes it challenging to explain this behavior. We hypothesize that ChatGPT-3.5’s better understanding of diverse prompts may be due to its more frequent usage, facilitated by free access.

Finally, Gemini 1.0 presents the lowest performance in our evaluation, with 80.0% correct responses. This model often complicates code generation by adding unnecessary explanations to nearly every line and failing to maintain proper indentation, making the code difficult to copy and paste. Interestingly, both ChatGPT-4 and Gemini 1.0 had the same percentage of wrong answers (1.67%). However, ChatGPT-4 produced fewer passable answers (11.67%) compared to Gemini’s 18.33%.

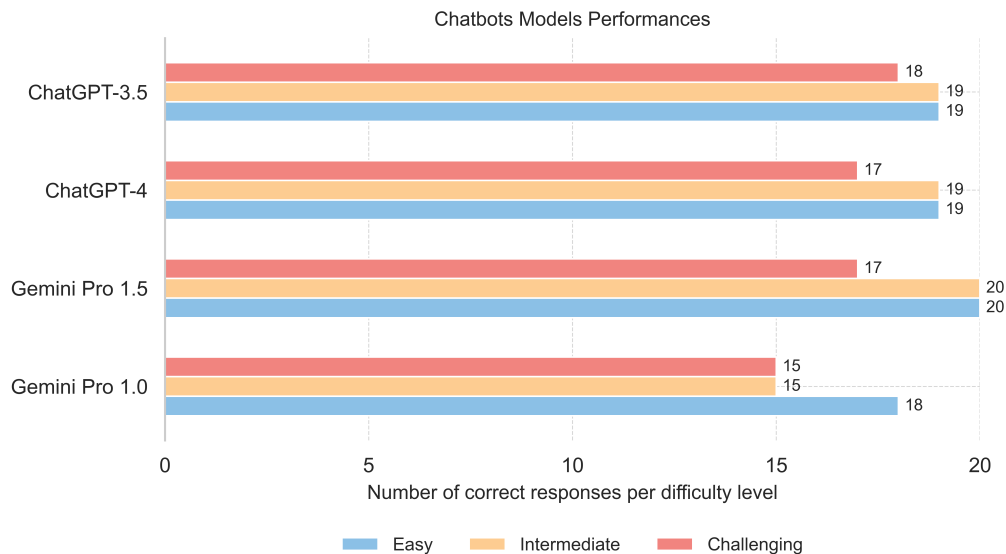


Figure 7: Number of correct answers obtained by Gemini, ChatGPT-4, and ChatGPT-3.5 across easy, intermediate, and challenging levels.

Figure 7 illustrates the number of correct answers obtained by Gemini 1.5 Pro, ChatGPT-3.5, ChatGPT-4, and Gemini 1.0 across easy, intermediate, and challenging levels. On the easy level, it shows that Gemini 1.5 leads with 20/20, and both ChatGPT-3.5 and ChatGPT-4 scored 19/20, slightly outperforming Gemini 1.0, which scored 18/20. At the intermediate level,

Gemini 1.5 Pro achieved a perfect score of 20/20. In comparison, ChatGPT-3.5 scored 19/20, while ChatGPT-4 and Gemini scored 16/20 and 15/20, respectively. In the challenging category, ChatGPT-3.5 excelled with a score of 18/20, followed by ChatGPT-4 and Gemini 1.5 Pro with 17/20 each, and Gemini 1.0 with 15/20. This comparison effectively highlights the relative strengths and challenges of each model across varying levels of difficulty.

LLaMA models. In this series of models, we evaluated two versions of LLaMA: LLaMa-2 and LLaMa-3.1. For LLaMA-2, we focused on two quantized models specialized in code: llama-2-coder-7b.Q3_K_L and llama-2-coder-7b.Q4_K_M. Contrary to the expectations set by Wei et al. (2022a) and Cobbe et al. (2021), who suggest that larger models perform better, our findings reveal that this is not always true. Specifically, the larger model, llama-2-coder-7b.Q3_K_L, performed worse than the medium-sized llama-2-coder-7b.Q4_K_M in terms of correct responses, highlighting the significance of bit quantization over model size alone. Overall, the large model scored 21.67%, while the medium-sized model achieved 25%, considering both correct and passable answers.

The primary weakness observed in the llama-2-coder-7 models is a flaw in their background knowledge. For instance, when asked whether a specific month falls within the winter season, the model generated the following code: `if (input_month \geq 12 and input_month \leq 21): return "Winter"`. This response shows a misunderstanding of the problem or a lack of knowledge that a year comprises only 12 months, making the model incorrectly extend beyond this range. Furthermore, the LLaMA models tend to generate verbose code and struggle offer a variety of range of options for answers, opting

instead to provide direct responses.

We evaluated the 8B parameter version of Llama-3.1 and for a fair comparison, we tested the non-quantized version, Meta-Llama-3.1-8B-Instruct, which achieved 91.66% correct responses. In contrast, the 4-bit and 5-bit quantized versions scored 86.66% and 88.33% correct responses, respectively. This clearly demonstrates the impact of quantization on performance. For example, in the quantize version of 4 and 5 bits, Meta-Llama-3.1-8B-Instruct overcome these models with 5% and 3.33% in the correct responses, respectively.

We also evaluated the 8B parameter version of Llama-3.1. To ensure a fair comparison, we tested the non-quantized version, Meta-Llama-3.1-8B-Instruct, which achieved a 91.66% accuracy in the correct responses. In contrast, the 4-bit and 5-bit quantized versions scored 86.66% and 88.33%, respectively, highlighting the impact of quantization on performance. Notably, the non-quantized Meta-Llama-3.1-8B-Instruct outperformed the quantized versions by 5% and 3.33% in accuracy.

Among CPU-friendly models for Python code generation, Meta-Llama-3.1-8B-Instruct excelled compared to llama-2-coder-7b, phi-2, mistral-7b-instruct-v0.2, zephyr-7b-beta, dolphin-2.6-mistral-7b, openhermes-2.5-mistral, and MiniCPM-2B-dpo-bf16, across all quantization levels (from 2 to 8 bits). This superior performance can be attributed to the extensive data used during pre-training and post-training, carefully curated through preprocessing. For a detailed breakdown of Llama’s performance across easy, medium, and challenging levels, refer to Figure 8.

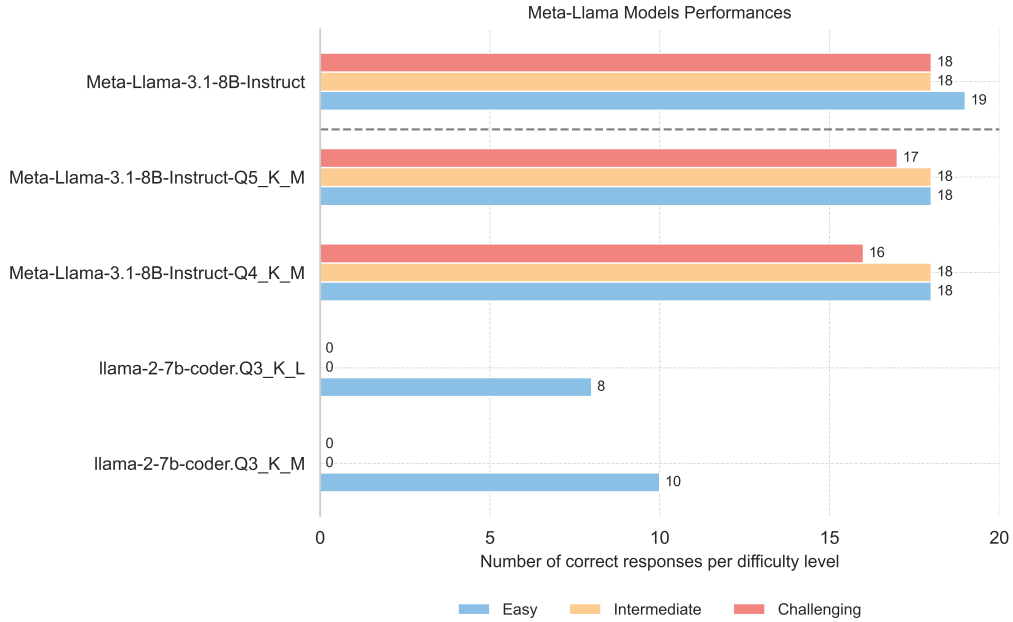


Figure 8: Number of correct answers obtained by LLaMa models across easy, intermediate, and challenging levels.

Phi-2 models. Our evaluation of Phi-2 models, each with different quantization levels (5, 6, and 8 bits), revealed that Phi-2.Q_8 achieved the highest accuracy (30%) in correctly generating Python code from prompts, as shown in Table 4. This surpassed the performance of Phi-2 models quantized to 5 and 6 bits, which obtained 28.33% and 26.67% correct answers, Interestingly, both Phi-2.Q6_K and Phi-2.Q_8 models achieved an identical success rate of 41.67% in providing passable responses. Overall, the Phi-2.Q_8 model demonstrates superior performance, with a combined correct and passable accuracy rate of 50.83%.

However, on our dataset, Phi-2 models generally performed at the lower end compared to other models in terms of correct answers at the easy level,

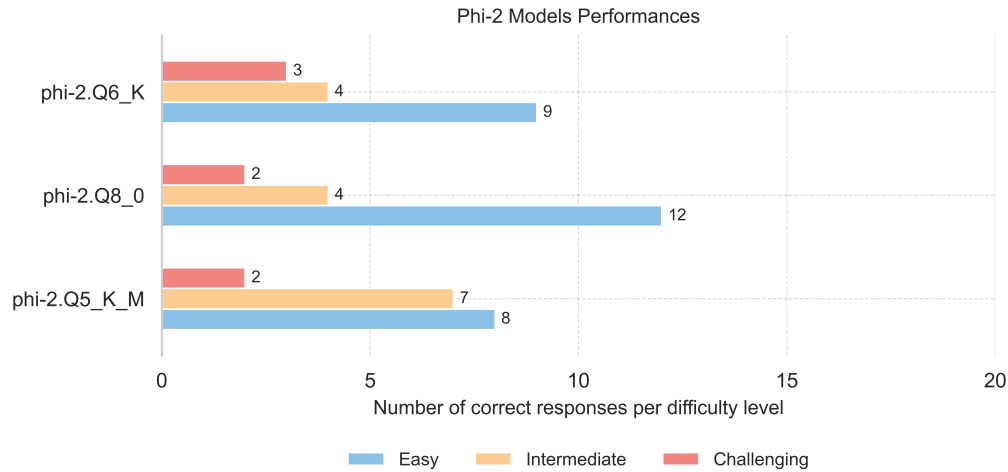


Figure 9: Number of correct answers obtained by Phi-2 models across easy, intermediate, and challenging levels.

as seen in Figure 9. The best-performing Phi-2 model, Phi-2.Q8_0, answered only 12 out of 20 questions correctly. For harder prompts, it yielded less than 20% accuracy, with only 4 out of 20 and 2 out of 20 correct responses for intermediate and challenging problems, respectively.

The primary limitations of the Phi-2 models are their verbosity, including unnecessary comments and often incomplete code, where the models generate the function header and docstring but leave the body empty. Additionally, there are numerous irrelevant comments unrelated to the prompt that indicate the models' tendency to hallucinate. Specifically, the Phi-2.Q5_K_M model struggles with proper Python code indentation, making it difficult to read and interpret it.

Mistral models. In our study, we assessed the performance of the mistral-7b-instruct-v0.2 model and its quantized variants. We explored quantization

levels from 2 to 5 bits across three model sizes: small, medium, and large. Additionally, we evaluate a large 3-bit quantized variant of the Openorca .

According to Table 4, the non-quantized mistral-7b-instruct-v0.2 model unsurprisingly achieved the highest performance among the mistral-7b-instruct family, with 90% correct answers. Among the quantized versions, mistral-7b-instruct-v0.2.Q4_K_M achieved the best accuracy (86.67%), followed by the mistral-7b-instruct-v0.2.Q4_K_S model (73.33%). An interesting observation is that both top-performing models were quantized to 4 bits. This finding contradicts the hypothesis that models quantized with more bits would necessarily provide higher accuracy. In fact, the mistral-7b-instruct-v0.2.Q5_K_M model, quantized with 5 bits, achieved a 28.34% lower accuracy than the mistral-7b-instruct-v0.2.Q4_K_M model. Furthermore, the accuracy of the remaining evaluated models (2 and 3 bits) consistently fell below 55.00%.

The last column of Table 4 provides the total performance of the model considering both the correct and the passable answers. The non-quantized mistral-7b-instruct-v0.2 achieved the highest score of 92.5%, followed by mistral-7b-instruct-v0.2.Q4_K_M with 90.83% and mistral-7b-instruct-v0.2.Q4_K_S with 81.67%. Surprisingly, the fourth-best performer was not the model quantized with 5 bits, instead, it was the mistral-7b-instruct-v0.2.Q3_K_M model that achieved 72.50% of accuracy. This result confirms that the performance of the models in Python generation is not solely determined by the number of bits used to quantize the model.

Further analysis of Figure 10 revealed interesting trends in correct responses by difficulty level. While mistral-7b-instruct-v0.2.Q4_K_S performed exceptionally well with easy samples, mistral-7b-instruct-v0.2.Q4_K_M sur-

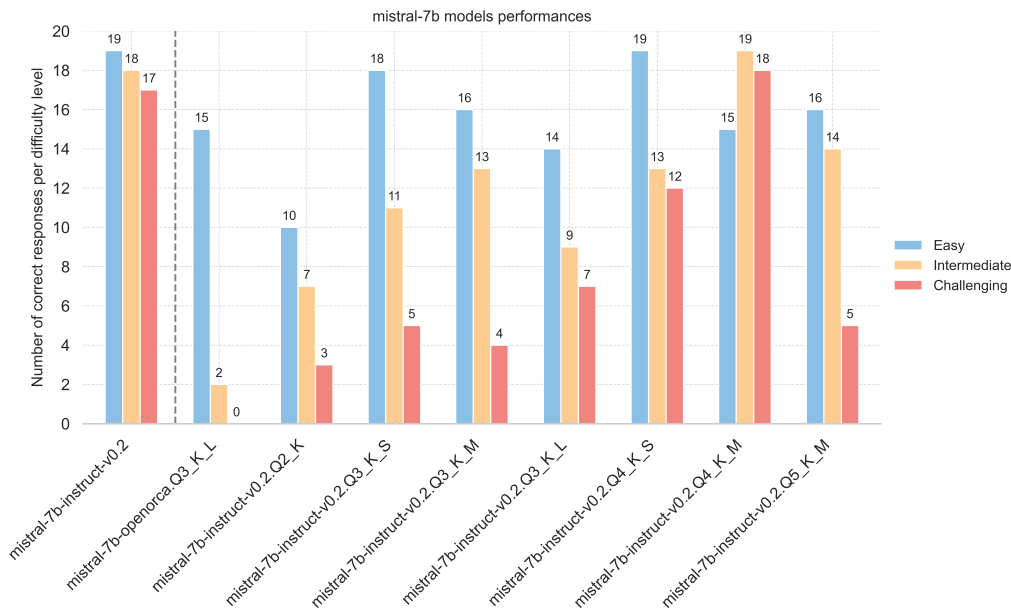


Figure 10: Number of correct answers obtained by mistral models across easy, intermediate, and challenging levels.

passed it in generating correct responses for intermediate and challenging samples. This suggests that mistral-7b-instruct-v0.2.Q4_K_M may be better suited for handling complex Python code generation.

After finding that mistral-7b-instruct-v0.2.Q4_K_M has comparable or superior performance with powerful close-models such as Gemini 1.0 and ChatGPT-4, we explored the performance of other mistral-based models. Notably, we looked into the capabilities of Mistral-7b-openorca quantized to 3-bits and found that it underperformed compared to the original Mistral model. It’s important to note that this model was fine-tuned on the OpenOrca dataset, which is not specifically designed for coding tasks. Nonetheless, it achieved 28.33% correct responses and 40% passable answers.

In addition, we delved deeper into other models that were fine-tuned on

top of Mistral, such as zephyr-7b-beta, dolphin-2.6-mistral, and openhermes-2.5-mistral. We discuss these models in the following sections, and compare their performance at different quantization levels of 4 bits and higher, allowing us to assess how such adjustments impact their efficiency and output quality. For models with high accuracy in their quantized versions, such as dolphin-2.6-mistral and openhermes-2.5-mistral-7b, we also provided the scores of their non-quantized counterparts.

zephyr-7b-beta models. Our tests on Zephyr models (4, 5, 6, and 8 bits), as seen in Figure 11, showed lower performance compared to mistral-7b-instruct-v0.2.Q4_K_M and some other quantized mistral variants. While zephyr-7b-beta.Q4_K_M achieved the highest score (45%) within the Zephyr family, it lagged behind mistral-7b-instruct-v0.2.Q4_K_M and ChatGPT-3.5 by 41.67% and 48.33%, respectively.

Zephyr-7b-beta.Q5_K_M achieved the second-best performance of 41.67% accuracy. Interestingly, both 6-bit and 8-bit Zephyr models obtained the same score (36.67%). However, considering overall performance within the Zephyr family, zephyr-7b-beta.Q4_K_M is the best performer, followed by zephyr-7b-beta.Q8_0. This is because the 8-bit Zephyr model achieved the highest percentage of passable responses, demonstrating an understanding of the prompt and code logic, albeit with syntactic inaccuracies. It is noteworthy to mention that zephyr-7b-beta.Q8_0 showcases distinct advantages, such as generating optimized and compact code by leveraging existing libraries like NumPy and utilizing list comprehensions. Contrary to its 5- and 6-bit-width counterparts, zephyr-7b-beta.Q8_0 seldom produces empty functions.

Zephyr is fine-tuned on top of Mistral, but worsens its results in Python

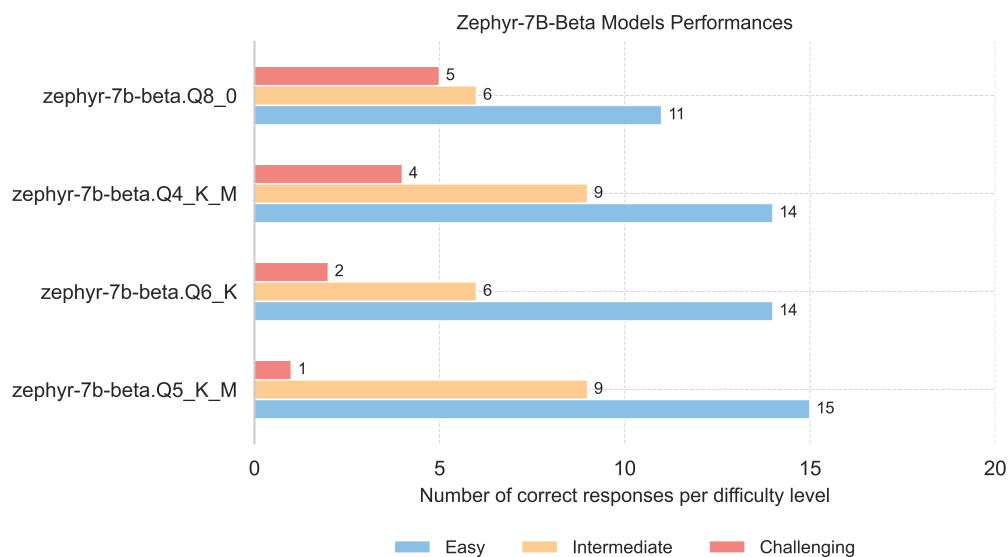


Figure 11: Number of correct answers obtained by zephyr-7b-beta models across easy, intermediate, and challenging levels.

Code generation. Designed to align with user intent and handle natural prompts effectively, it was fine-tuned on the UltraChat 200k dataset (Ding et al., 2023; Tunstall et al., 2023). UltraChat focuses on three broad areas: understanding real-world concepts, generating various text formats, and working with existing text (rewriting, translation, etc.). Since code generation isn’t a core focus in UltraChat, the dataset likely contains minimal code content, which could explain Zephyr’s weaker performance.

dolphin-2.6-mistral. We conducted tests on three quantized versions of the dolphin-2.6-mistral model, using 5, 6, and 8 bits. , The dolphin-2.6-mistral-7b.Q8_0 model, with the highest bit width, achieved the best performance, with 50% correct responses, closely followed by dolphin-2.6-mistral-7b.Q5_K_M at 48.33%. While the non-quantized version, dolphin-2.6-mistral-

7b, scored 60%, the performance gap with the quantized models is modest. However, a significant drop is observed with the dolphin-2.6-mistral-7b.Q8_0 model, which only achieved 30% accuracy.

When considering only correct responses, Figure 12 highlights the superior performance of the dolphin-2.6-mistral-7b model, followed by the dolphin-2.6-mistral-7b.Q8_0 variant, across difficult and intermediate levels. Notably, dolphin-2.6-mistral-7b.Q5_K_M achieved the highest number of correct answers at the easy level (16/20), while the other two models closely followed with 15/20 each. Both the 8-bit and 5-bit dolphin-2.6-mistral models maintained consistent performance in the intermediate (8/20 and 7/20, respectively) and difficult levels (7/20 and 6/20, respectively). In contrast, the dolphin-2.6-mistral-7b.Q6_K_M model experienced a sharp decline, failing to generate any correct answers at the difficult level (0/20).

The enhanced performance of the Dolphin-2.6-Mistral models over Zephyr models can be explained by the datasets used for fine-tuning. Zephyr models were fine-tuned using the UltraChat 200k dataset (Ding et al., 2023; Tunstall et al., 2023), which has a limited portion of source code. In contrast, Dolphin-2.6-Mistral models benefited from a fine-tuning with Magicoder-Evol-Instruct-110K (Wei et al., 2023) and Magicoder-OSS-Instruct-75K datasets (Team, 2024b), among others. This selection of datasets is a key factor in Dolphin-2.6-Mistral being better than Zephyr models.

openhermes-2.5-mistral. We evaluated four quantized versions of the openhermes-2.5-mistral model, utilizing 4, 5, 6, and 8 bits. The analysis revealed that the 4-bit version and the 5-bit version yielded accuracy rates of 53.33% and 55% in the correct responses, respectively. In contrast,

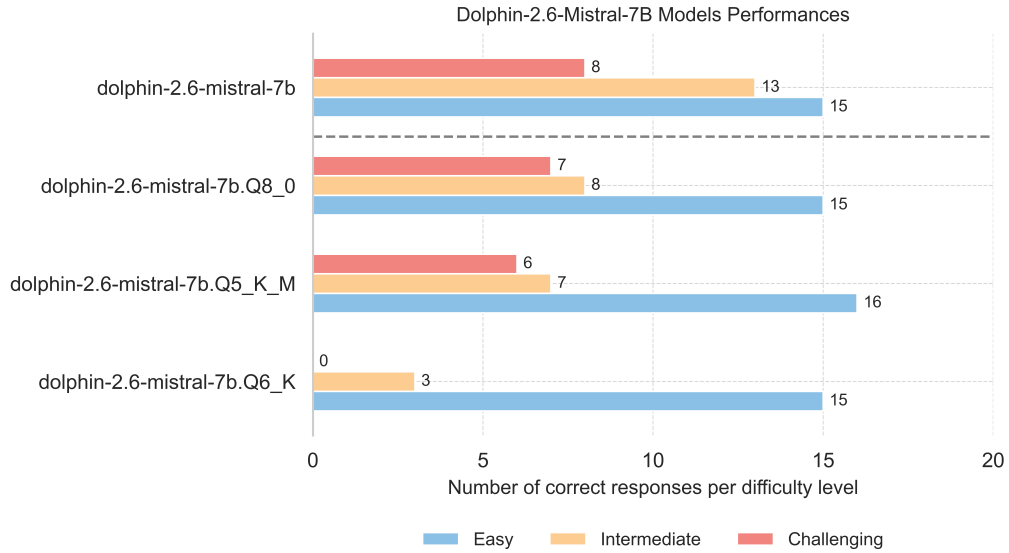


Figure 12: Number of correct answers obtained by dolphin-2.6-mistral models across easy, intermediate, and challenging levels.

both the 6-bit and 8-bit versions demonstrated superior accuracy, achieving a rate of 58.33%. While these scores are close to the 66.66% accuracy of the non-quantized OpenHermes-2.5-Mistral-7B, they highlight the effectiveness of CPU-friendly models for generating Python code from natural language.

Moreover, as shown in Table 4, the 6-bit model (openhermes-2.5-mistral-7b.Q6_K) outperforms the 8-bit model (openhermes-2.5-mistral-7b.Q8_0) in terms of overall quality. This is because the 6-bit model produced more passable responses and fewer incorrect answers. While the 8-bit model achieved a combined accuracy of 69.17% (correct + passable), the 6-bit model reached a higher accuracy of 70.83%.

Examining correct answers across different difficulty levels (Figure 13), the findings indicate that both the openhermes-2.5-mistral-7b.Q4_K_M and

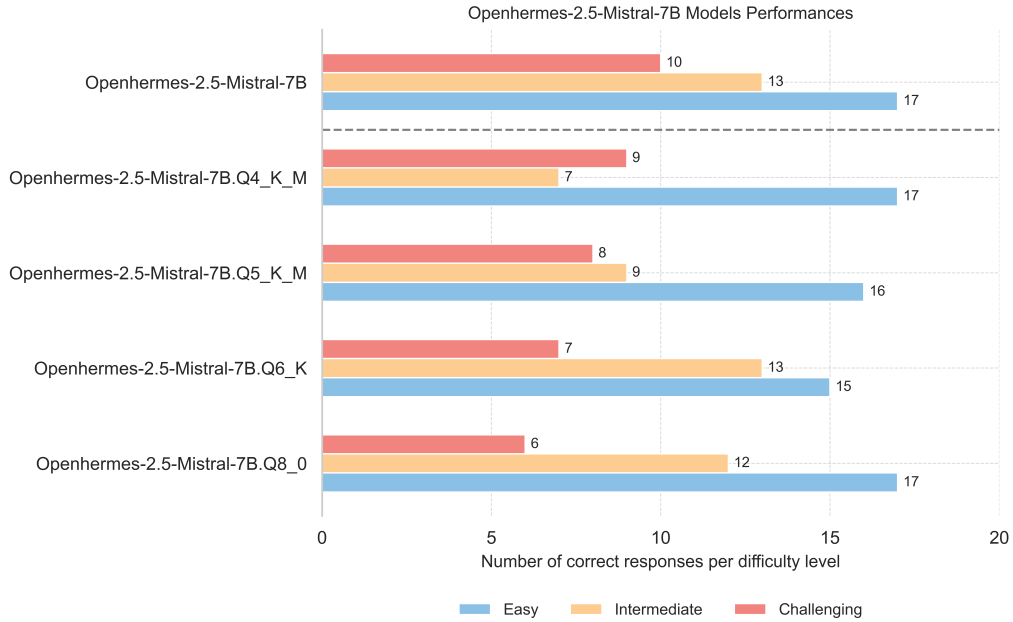


Figure 13: Number of correct answers obtained by openhermes-2.5-mistral models across easy, intermediate, and challenging levels.

openhermes-2.5-mistral-7b.Q8_0 models excelled in the easy category, each securing 17 correct responses. They were closely followed by the 5-bit model with 16 correct responses, while the 6-bit version recorded the lowest success rate in this category. Furthermore, the openhermes-2.5-mistral-7b.Q4_K_M and openhermes-2.5-mistral-7b.Q5_K_M models secured the highest success rate in the challenging category, each achieving 8 out of 20 correct responses. While not the best in either easy or challenging categories, the 6-bit openhermes-2.5-mistral-7b.Q6_K_M model achieves consistent performance across all difficulty levels. This stability makes it the overall best choice, considering both correct and passable responses.

The main weaknesses of the openhermes-2.5-mistral family are halluci-

nations and lack of code conciseness. For instance, the model references nonexistent NumPy functions or incorrectly attributes SciPy functions to the NumPy library. Moreover, the generated code lacks compactness; for example, rather than employing the more concise condition: $18.5 \leq \text{bmi} \leq 24.9$, it opts for the lengthier $\text{bmi} \geq 18.5$ and $\text{bmi} \leq 24.9$. Additionally, it's noteworthy that the openhermes-2.5-mistral 4-bit version tends to leverage native structures more effectively for algorithm development, in contrast to the 6 and 8-bit versions, which more frequently utilize library functions.

MiniCPM-2B-dpo-bf16. Lastly, we assessed the performance of the MiniCPM model, which outperformed models like Llama2-7B, Mistral-7B, and Llama2-13B in code and mathematical reasoning (Inc. and TsinghuaNLP, 2024), including on the HumanEval dataset. While MiniCPM also surpassed Llama2-7B in our dataset, it did not consistently outperform all Mistral variants. Specifically, MiniCPM scored higher in terms of correct responses compared to Mistral models up to 3 bits. However, Mistral models from 4 bits and above achieved better results. It is important to note that MiniCPM evaluation was impacted by its non-compliance to the expected output format. Despite this penalization, MiniCPM achieved 50% of correct responses, 36.67% of passable responses, and 13.33% of incorrect answers, as shown in Figure 14.

Size, RAM and Inference Time. The analysis presented in Table 4 offers an overview of the performances and system requirements of various models. Notably, the models that need GPU to run the inference, are the heaviest exceeding the 15GB in space. With respect to quantize models, the heaviest models are zephyr-7b-beta.Q8_0, dolphin-2.6-mistral-7b.Q8_0,

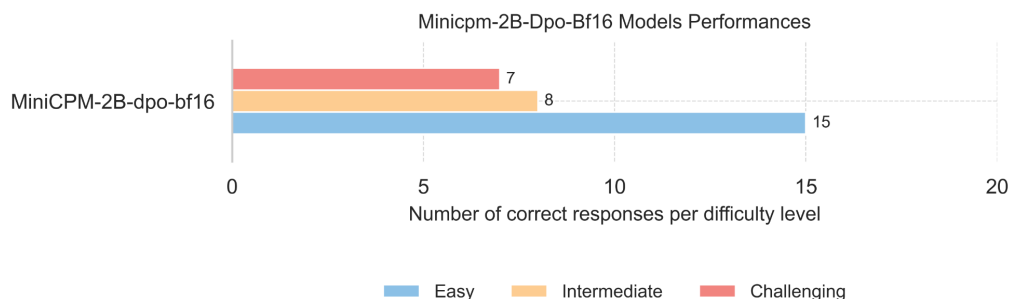


Figure 14: Number of correct answers obtained by the MiniCPM model across easy, intermediate, and challenging levels.

and openhermes-2.5-mistral-7b.Q8_0, with a size of 7.70 GB, which is still manageable on a conventional machine. In terms of RAM usage, MiniCPM is the most demanding, requiring 10.9 GB. This moderate use of resources allows for the possibility of performing other tasks in parallel, thus optimizing machine utilization.

The inference times vary among the models. The Phi family of models is the fastest, followed by Mistral. Conversely, Meta-Llama-3.1-8B-Instruct quantize models exhibits the longest inference time, followed by MiniCPM and Dolphin. Overall, Mistral stands out as offering the best balance between size, RAM usage, inference speed, and overall performance on this dataset.

8. Results and Discussion: HumanEval and EvalPlus dataset

In this section, we discuss the results obtained by the evaluated models on the HumanEval and EvalPlus datasets. Our primary reference for comparison is the EvalPlus Leaderboard (Team, 2024a), as outlined by Liu et al. (2023b). However, it is noteworthy that the samples in this leaderboard were generated from scratch and underwent preprocessing with a sanitizer script.

| | | | HumanEval | EvalPlus |
|--------|-----------------------------------|-----------------------|-----------------|-----------------|
| | Models | Inference Time (ms) ↓ | Pass@1 0-shot ↑ | Pass@1 0-shot ↑ |
| Closed | Claude 3.5 Sooner | × | 92.0 | 82.3 |
| | GPT-4-Turbo (April 2024) | × | 90.2 | 86.6 |
| | Gemini 1.5 Pro | × | 84.1 | 77.25 |
| | GPT-3.5-Turbo (Nov 2023) | × | 76.8 | 70.7 |
| | Gemini 1.0 Pro | × | 63.4 | 55.5 |
| GPU | Meta-Llama-3.1-8B-Instruct | 28.26 | 60.4 | 57.35 |
| | Mistral-7B-Instruct-v0.2 | 14.96 | 42.10 | 36.00 |
| | dolphin-2.6-mistral-7b | 8.89 | 60.25 | 52.32 |
| | OpenHermes-2.5-Mistral-7B | 7.56 | 50.12 | 46.25 |
| CPU | Meta-Llama-3.1-8B-Instruct-Q4_K_M | 458.25 | 54.2 | 50.7 |
| | Meta-Llama-3.1-8B-Instruct-Q5_K_M | 478.12 | 54.32 | 51.87 |
| | llama-2-coder-7b.Q3_K_M | 209.93 | 11.52 | 10.21 |
| | llama-2-coder-7b.Q3_K_L | 224.33 | 12.52 | 10.87 |
| | llama-2-coder-7b.Q4_K_M | 279.63 | 14.12 | 11.36 |
| | mistral-7b-instruct-v0.2.Q2_K | 188.38 | 15.85 | 14.02 |
| | mistral-7b-instruct-v0.2.Q3_K_S | 207.47 | 32.92 | 28.65 |
| | mistral-7b-instruct-v0.2.Q3_K_M | 233.43 | 27.40 | 25.00 |
| | mistral-7b-instruct-v0.2.Q3_K_L | 240.58 | 23.78 | 21.34 |
| | mistral-7b-instruct-v0.2.Q4_K_S | 242.29 | 23.17 | 21.95 |
| | mistral-7b-instruct-v0.2.Q4_K_M | 245.57 | 23.17 | 24.39 |
| | mistral-7b-instruct-v0.2.Q5_K_M | 292.74 | 13.41 | 12.19 |
| | mistral-7b-openorca.Q3_K_L | 219.98 | 34.75 | 32.31 |
| | dolphin-2.6-mistral-7b.Q5_K_M | 288.57 | 51.82 | 45.12 |
| | dolphin-2.6-mistral-7b.Q6_K | 332.83 | 54.26 | 48.17 |
| | dolphin-2.6-mistral-7b.Q8_0 | 425.05 | 49.39 | 44.51 |
| | openhermes-2.5-mistral-7b.Q4_K_M | 265.82 | 43.29 | 35.97 |
| | openhermes-2.5-mistral-7b.Q5_K_M | 294.62 | 45.12 | 37.80 |
| | openhermes-2.5-mistral-7b.Q6_K | 328.12 | 43.90 | 36.58 |
| | openhermes-2.5-mistral-7b.Q8_0 | 426.11 | 39.63 | 32.31 |

Table 5: Evaluation on HumanEval and EvalPlus datasets. Q_j : quantization using j bit width, K: the use of k-means clustering in the quantization, S , M , L : Small, Medium, and Large model size after quantization. Best results of each category are in bold.

To ensure a fair comparison with the scores reported on the EvalPlus Leaderboard, we used the greedy search decoding algorithm to generate code outputs. We evaluated both quantized models and the original, missing models that require GPU for inference, as shown in Table 5. Additionally,

we assessed code correctness using the pass@1 metric and the 0-shot approach in the input prompts.

Claude 3.5 Sooner vs ChatGPT-3.5 vs ChatGPT-4 vs Gemini 1.0-1.5 Pro. According to (Dubey et al., 2024; Anthropic, 2024) and Table 5, Claude 3.5 Sonnet outperforms both ChatGPT-3.5, ChatGPT-4, and the Gemini models in the HumanEval dataset, achieving 92.0% accuracy. GPT-4-Turbo follows closely with 90.2%. Gemini 1.5 Pro shows a significant drop in performance, scoring 84.1%, which is 7.3% higher than GPT-3.5-Turbo’s 76.8%. Finally, Gemini Pro 1.0 ranks fifth with 63.4% accuracy.

In the EvalPlus evaluation, the rankings shifted. GPT-4-Turbo leads at 86.6% accuracy, followed by Claude 3.5 Sonnet at 82.3%. The other positions remain consistent: Gemini 1.5 Pro is third with 77.25%, GPT-3.5-Turbo is fourth with 70.7%, and Gemini Pro 1.0 is fifth with 55.5%.

These results highlight the rapid advancements in language model capabilities, especially in code generation. Since the release of ChatGPT-3.5 in November 2022, the community has achieved impressive progress in less than two years.

LLaMA Models. In this set of experiments, we evaluated the Meta-Llama-3.1-8B-Instruct model alongside two of its quantized variants (4-bit and 5-bit), as well as three quantized versions of the Llama-2-Coder-7B model (Q3_K_M, Q3_K_L, Q4_K_M). As shown in Table 5, the non-quantized Meta-Llama-3.1-8B-Instruct achieved the highest scores in both the HumanEval and EvalPlus datasets, with 60.4% and 57.35% accuracy, respectively. The quantized versions were close behind, with the 4-bit and 5-bit models scoring 54.2% and 50.7%, and 54.32% and 51.87%, respectively.

In contrast, the Llama-2-Coder-7B quantized models recorded the lowest scores in both datasets. In the HumanEval dataset, they achieved 11.52%, 12.52%, and 14.12% accuracy for the Q3_K_M, Q3_K_L, and Q4_K_M models, respectively. Performance on the EvalPlus dataset was even lower, with the Llama-2-Coder-7B.Q4_K_M emerging as the most effective variant at 11.36% accuracy.

These results highlight the superiority of LLaMA-3.1 in Python code generation tasks. Despite being specialized for coding, the Llama-2-Coder-7B model was outperformed by the Llama-3.1-8B-Instruct, and our findings consistently revealed that the code produced by Llama-2-Coder was frequently incomplete and riddled with syntax errors.

Mistral. In this set of experiments, we evaluate the performance of the non-quantized Mistral-7B-Instruct-v0.2 model and its quantized versions. According to the EvalPlus Leaderboard and Table 5, the non-quantized model achieved 42.1% accuracy on HumanEval and 36% on EvalPlus. The performance difference between the non-quantized and quantized models is minimal. For instance, the best quantized performance on the HumanEval dataset comes from the Mistral-7B-Instruct-v0.2.Q3_K_S model with a pass@1 score of 32.92%, followed by Mistral-7B-Instruct-v0.2.Q3_K_M at 27.40%. Surprisingly, the lowest scores were obtained by Mistral-7B-Instruct-v0.2.Q2_K and Mistral-7B-Instruct-v0.2.Q5_K_M, with pass@1 scores of 15.85% and 13.41% respectively. This highlights, once again, that quantizing with more bits does not necessarily lead to better results.

Unsurprisingly, the quantized Mistral models showed lower performance on the EvalPlus dataset, which includes a greater number of unit tests. The

models achieved an average score 1.6 points lower than on the HumanEval dataset. Nevertheless, the overall ranking of the tested models remained consistent across both datasets.

Finally, we evaluated the performance of the Mistral-7b-openorca.Q3_K_L model, where it surpassed other quantized Mistral variants with accuracies of 34.75% and 32.31%, respectively. However, as detailed in Section 7, its performance on our dataset was significantly lower than that of other quantized Mistral models due to penalties for failing to generate outputs in the correct format. This discrepancy highlights the model’s difficulties in interpreting input prompts with restrictive conditions for generating Python code.

dolphin-2.6-mistral-7b. In our evaluation, we assessed the non-quantized dolphin-2.6-mistral-7b model alongside its 5-bit, 6-bit, and 8-bit quantized versions. Unsurprisingly, the non-quantized model achieved the highest accuracy, scoring 60.25% on HumanEval and 52.32% on EvalPlus, with only minor differences compared to the quantized versions. All dolphin models outperformed the mistral models. On HumanEval, the 6-bit model led the quantized versions with 54.26%, followed by the 5-bit model at 51.82%, and the 8-bit model at 49.39%. Similarly, in EvalPlus, the 6-bit model scored highest among the quantized versions with 48.17%, while the 8-bit model had the lowest score at 44.51%.

These findings show that the dolphin-2.6-mistral-7b quantized models outperformed the Mistral quantized models in Python code generation across both the HumanEval and EvalPlus datasets. However, it’s worth noting that on our dataset, the dolphin-2.6-mistral-7b model scored lower than Mistral. It’s important to note that although dolphin-2.6-mistral-7b generates correct

code, it was penalized for not following the specific output format we provided. This highlights the importance to adhere to the correct format, even if the code is functionally correct. Thus, while Mistral excels in understanding the required output format, dolphin-2.6 demonstrates a high capability in generating Python code but often fails to deliver it in the prescribed format.

openhermes-2.5-mistral-7. The evaluation results for the non-quantized OpenHermes-2.5-Mistral-7B model were 50.12% on HumanEval and 46.25% on EvalPlus. These results are competitive compared to more complex models like Meta-Llama-3.1-8B-Instruct, which scored 60.4% and 57.35%, respectively. Additionally, the quantized OpenHermes-2.5-Mistral-7B outperformed the quantized Mistral models. However, it scores lower than the dolphin-2.6-mistral-7b models, which is expected since the latter were trained using code datasets, enhancing their performance in such tasks. Although mistral models achieved better scores in our dataset but lower scores on HumanEval and EvalPlus, it's important to note that openhermes-2.5-mistral-7b also produces correct and functional codes. However, similarly to dolphin-2.6-mistral-7b models, these were also penalized due to not adhering to the required output format, resulting in a reduced score of 0.5.

Consistent with our previous findings, the openhermes-2.5-mistral-7 model often struggles with variable name consistency when generating code. Although it adheres to the variable names specified in the prompt, it frequently introduces errors when using new variables. For example, it might refer to the same variable as 'time' initially and then as 'times' later, leading to failures in unit tests. It would be advantageous to develop strategies to address this issue, such as through post-processing or improved token generation

techniques.

Inference Time. Table 5 shows minimal differences in inference times among the evaluated models. Mistral and LLaMa are the fastest, with a small difference of 104.36 milliseconds between Mistral’s 2-bit and 8-bit quantized versions. The Meta-Llama-3.1-8B-Instruct family of models have longer inference times on average, with the 5-bit quantized version being the slowest at 478.12 milliseconds. These results highlight the efficiency and practicality of these CPU-friendly models in terms of both inference speed and performance.

9. Conclusions and Future Directions

In this paper, we examine the effectiveness of CPU-friendly models in generating Python code under various scenarios. Specifically, our dataset tasked models with generating Python code based on a given problem statement, using predetermined variable names and returning a specified set of options included in the input prompt. In contrast, in the HumanEval and EvalPlus datasets, the challenge was to generate Python code from the function’s signature and docstring, allowing us to assess the functional correctness of LLM-synthesized code.

Our evaluation of CPU-friendly models across different datasets assesses not only their ability to generate correct Python code but also how well they followed the specified variable names and output formats.

The advantages of this work are to: (1) facilitate performance comparisons of CPU-friendly models, enabling individuals or small companies with

limited computational resources to better choose a smaller model as an alternative to large, costly models; (2) propose a new Python generation dataset consisting of 60 programming problems. This dataset extends the existing benchmark datasets, HumanEval and EvalPlus, allowing for deeper understanding of these models; (3) propose engineered prompts for generating and evaluating Python code; (4) propose a semi-manual evaluation of state-of-the-art methods for Python code generation; and (5) make all of our work, experiments, and dataset publicly available to advance research in code generation and provide access to a thorough analysis of CPU-friendly models, which can serve as strong alternatives to very large or closed models. We define a set of disadvantages in Section 10.

Our findings reveal two key points: 1) CPU-friendly models achieve competitive scores compared to advanced chatbot models like ChatGPT-3.5, ChatGPT-4, and Gemini, which require significant GPU resources for processing; 2) some models, such as dolphin-2.6-mistral-7b, excel in generating Python code but fall short in meeting the required output formats. Often, even if the generated code is correct, these models are penalized for not returning the specified options from the input prompt. A similar issue affects the openhermes-2.5-mistral-7b model, albeit with slightly lower scores than dolphin-2.6-mistral-7b. In contrast, models like Mistral demonstrate consistent performance in understanding prompts and generating code, leading to superior results on our dataset, though they perform worse on HumanEval and EvalPlus datasets.

Remarkably, llama.cpp project has made these competitive outcomes using standard computers possible, a development that seemed highly unlikely

just a few months or years ago. This progress highlights the feasibility of running sophisticated language models on CPUs today. However, it is important to mention that we did not utilize more powerful models such as Mixtral, an enhanced version of Mistral. Despite existing quantized versions, Mixtral still requires considerable computational resources which our machines could not support. For instance, the 2-bit quantized version of Mixtral requires 15.6 GB of storage and 18.14 GB of RAM, while the 8-bit version needs 49.62 GB of storage and 52.12 GB of RAM. Thus, these models do not align with our objective of demonstrating the performance of CPU-friendly models on standard machines.

In our future work, we plan to explore we plan to explore the expanded capabilities of CPU-friendly models across various code-related tasks, such as defect detection, cloze tests, code refinement, and code translation. Our key focus will be on creating a comprehensive framework to evaluate the performance of diverse LLMs in their quantized forms at different bit levels. This framework will integrate the Automatic Prompt Engineer (APE) (Zhou et al., 2023) framework, which generates various candidate prompts, executes them to generate responses, and evaluates the responses to select the most effective prompt. This integration can effectively enhance the quality of the prompts used in our experiments.

Moreover, the proposed framework will allow for a thorough and systematic analysis of how quantization affects both precision and task-specific outcomes in code-related contexts. Through this approach, we aim to gain deeper insights into the trade-offs between computational efficiency and model accuracy, ultimately guiding the optimization of LLMs for code-based appli-

cations.

10. Challenges and Limitations

Our work faces three main limitations. The first is the rapid emergence of new LLMs. At the time of this publication, several new models have been introduced, and we cannot yet assess the performance of their quantized versions. The second limitation concerns the dataset used. Our dataset presents a significant challenge for LLMs, as it requires them to accurately relate and comprehend three key elements: the problem, the context, and the possible solutions. This complexity tests the models’ ability to understand natural language and correctly generate Python code.

Another limitation of our work is that LLMs are usually trained on high-level imperative programming languages like Python, resulting in the underrepresentation of low-resource, declarative, and low-level languages. This bias limits the model’s versatility across different programming paradigms. A comprehensive review of cross-language generation and providing a specialized dataset for these underrepresented languages would be essential to address this gap and study the model’s applicability in more languages.

Similarly, existing datasets like HumanEval and EvalPlus are designed to evaluate LLMs’ ability to generate Python code based on natural language prompts. However, these datasets lack a variety of scenarios that would more thoroughly test the capabilities of quantized models in code generation, such as varying the structure of input text. This gap highlights the need for more comprehensive evaluation methods to fully understand the impact of quantization on code generation tasks.

Acknowledgments

We would like to acknowledge Novelis for their support in publishing this article. We are especially grateful for the assistance and contributions of their research team.

References

- Agarap, A.F., 2019. Deep learning using rectified linear units (relu). [arXiv:1803.08375](https://arxiv.org/abs/1803.08375).
- Aghajanyan, A., Huang, B., Ross, C., Karpukhin, V., Xu, H., Goyal, N., Okhonko, D., Joshi, M., Ghosh, G., Lewis, M., Zettlemoyer, L., 2022. Cm3: A causal masked multimodal model of the internet. [arXiv:2201.07520](https://arxiv.org/abs/2201.07520).
- Ahmad, W., Chakraborty, S., Ray, B., Chang, K.W., 2021. Unified pre-training for program understanding and generation, in: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, Online. pp. 2655–2668. URL: <https://aclanthology.org/2021.naacl-main.211>, doi:10.18653/v1/2021.naacl-main.211.
- Ahmed, I., Roy, A., Kajol, M., Hasan, U., Datta, P.P., Reza, M.R., 2023. Chatgpt vs. bard: a comparative study. Authorea Preprints .
- AI@Meta, 2024. Llama 3 model card. URL: https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md.
- Ainslie, J., Lee-Thorp, J., de Jong, M., Zemlyanskiy, Y., Lebron, F., Sanghai, S., 2023. GQA: Training generalized multi-query transformer models from multi-head checkpoints, in: Bouamor, H., Pino, J., Bali, K. (Eds.), Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Singapore.

- pp. 4895–4901. URL: <https://aclanthology.org/2023.emnlp-main.298>, doi:10.18653/v1/2023.emnlp-main.298.
- Anand, Y., Nussbaum, Z., Duderstadt, B., Schmidt, B., Mulyar, A., 2023. Gpt4all: Training an assistant-style chatbot with large scale data distillation from gpt-3.5-turbo. <https://github.com/nomic-ai/gpt4all>.
- Anthropic, 2024. Claude 3.5 sonnet model card addendum. URL: https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf.
- Austin, J., Odena, A., Nye, M., Bosma, M., Michalewski, H., Dohan, D., Jiang, E., Cai, C., Terry, M., Le, Q., Sutton, C., 2021. Program synthesis with large language models. [arXiv:2108.07732](https://arxiv.org/abs/2108.07732).
- Ba, J.L., Kiros, J.R., Hinton, G.E., 2016. Layer normalization. [arXiv:1607.06450](https://arxiv.org/abs/1607.06450).
- Beltagy, I., Peters, M.E., Cohan, A., 2020. Longformer: The long-document transformer. URL: <https://arxiv.org/abs/2004.05150>, [arXiv:2004.05150](https://arxiv.org/abs/2004.05150).
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al., 2020. Language models are few-shot learners, in: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (Eds.), Advances in Neural Information Processing Systems, Curran Associates, Inc.. pp. 1877–

1901. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf.
- Chee, J., Cai, Y., Kuleshov, V., Sa, C.D., 2023. QuIP: 2-bit quantization of large language models with guarantees. Thirty-seventh Conference on Neural Information Processing Systems URL: <https://openreview.net/forum?id=xrk9g5vcXR>.
- Chen, M., Shao, W., Xu, P., Wang, J., Gao, P., Zhang, K., Qiao, Y., Luo, P., 2024. Efficientqat: Efficient quantization-aware training for large language models. arXiv preprint arXiv:2407.11062 .
- Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H.P.d.O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G.e.a., 2021. Evaluating large language models trained on code. ArXiv abs/2107.03374. URL: <https://api.semanticscholar.org/CorpusID:235755472>.
- Chen, W., 2023. Large language models are few(1)-shot table reasoners. arXiv:2210.06710.
- Chiang, W.L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J.E., Stoica, I., Xing, E.P., 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. URL: <https://lmsys.org/blog/2023-03-30-vicuna/>.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H.W., Sutton, C., Gehrmann, S., et al., 2022. Palm: Scaling language modeling with pathways. arXiv:2204.02311.

- Christiano, P.F., Leike, J., Brown, T., Martic, M., Legg, S., Amodei, D., 2017. Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems* 30. URL: https://proceedings.neurips.cc/paper_files/paper/2017/file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf.
- Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., Schulman, J., 2021. Training verifiers to solve math word problems. [arXiv:2110.14168](https://arxiv.org/abs/2110.14168).
- Coursera, 2023. Most popular programming languages in 2024. URL: <https://www.coursera.org/articles/popular-programming-languages>.
- Cui, G., Yuan, L., Ding, N., Yao, G., Zhu, W., Ni, Y., Xie, G., Liu, Z., Sun, M., 2023. Ultrafeedback: Boosting language models with high-quality feedback. [arXiv:2310.01377](https://arxiv.org/abs/2310.01377).
- Dai, Z., Yang, Z., Yang, Y., Carbonell, J., Le, Q.V., Salakhutdinov, R., 2019. Transformer-xl: Attentive language models beyond a fixed-length context. [arXiv:1901.02860](https://arxiv.org/abs/1901.02860).
- Daniele, L., Suphavadeeprasit, 2023. Amplify-instruct: Synthetically generated diverse multi-turn conversations for efficient llm training. [arXiv preprint arXiv:\(coming soon\)](https://arxiv.org/abs/coming-soon) URL: <https://huggingface.co/datasets/LDJnr/Capybara>.
- Dao, T., Fu, D., Ermon, S., Rudra, A., Ré, C., 2022. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems* 35, 16344–16359.

- Dauphin, Y.N., Fan, A., Auli, M., Grangier, D., 2017. Language modeling with gated convolutional networks, in: Proceedings of the 34th International Conference on Machine Learning - Volume 70, JMLR.org. p. 933–941.
- Dery, L., Kolawole, S., Kagey, J.F., Smith, V., Neubig, G., Talwalkar, A., 2024. Everybody prune now: Structured pruning of llms with only forward passes. arXiv preprint arXiv:2402.05406 .
- Dettmers, T., Svirschevski, R., Egiazarian, V., Kuznedelev, D., Frantar, E., Ashkboos, S., Borzunov, A., Hoefler, T., Alistarh, D., 2023. Spqr: A sparse-quantized representation for near-lossless llm weight compression. URL: <https://arxiv.org/abs/2306.03078>, arXiv:2306.03078.
- Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2019. BERT: Pre-training of deep bidirectional transformers for language understanding, in: Burstein, J., Doran, C., Solorio, T. (Eds.), Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), Association for Computational Linguistics, Minneapolis, Minnesota. pp. 4171–4186. URL: <https://aclanthology.org/N19-1423>, doi:10.18653/v1/N19-1423.
- Ding, N., Chen, Y., Xu, B., Qin, Y., Zheng, Z., Hu, S., Liu, Z., Sun, M., Zhou, B., 2023. Enhancing chat language models by scaling high-quality instructional conversations. arXiv:2305.14233.

- Dong, Y., Jiang, X., Jin, Z., Li, G., 2023. Self-collaboration code generation via chatgpt. [arXiv:2304.07590](https://arxiv.org/abs/2304.07590).
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al., 2024. The llama 3 herd of models. URL: <https://arxiv.org/abs/2407.21783>, [arXiv:2407.21783](https://arxiv.org/abs/2407.21783).
- Frantar, E., Alistarh, D., 2023. Sparsegpt: Massive language models can be accurately pruned in one-shot. URL: <https://arxiv.org/abs/2301.00774>, [arXiv:2301.00774](https://arxiv.org/abs/2301.00774).
- Frantar, E., Ashkboos, S., Hoefler, T., Alistarh, D., 2023. Gptq: Accurate post-training quantization for generative pre-trained transformers. URL: <https://arxiv.org/abs/2210.17323>, [arXiv:2210.17323](https://arxiv.org/abs/2210.17323).
- Fried, D., Aghajanyan, A., Lin, J., Wang, S., Wallace, E., Shi, F., Zhong, R., Yih, S., Zettlemoyer, L., Lewis, M., 2023. Incoder: A generative model for code infilling and synthesis. The Eleventh International Conference on Learning Representations URL: <https://openreview.net/forum?id=hQwb-1bM6EL>.
- Gerganov, G., 2023. llama.cpp. URL: <https://github.com/ggerganov/llama.cpp>.
- Gerganov, G., 2024. Gguf. <https://huggingface.co/docs/hub/en/gguf>.
- Gunasekar, S., Zhang, Y., Aneja, J., Mendes, C.C.T., Giorno, A.D., Gopi, S., Javaheripi, M., Kauffmann, P., et al., 2023. Textbooks are all you need. [arXiv:2306.11644](https://arxiv.org/abs/2306.11644).

- Hartford, E., 2023. Dolphin dataset. URL: <https://huggingface.co/datasets/cognitivecomputations/dolphin>.
- Heidari, A., Navimipour, N.J., Zeadally, S., Chamola, V., 2024. Everything you wanted to know about chatgpt: Components, capabilities, applications, and opportunities. *Internet Technology Letters* , e530.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., Steinhardt, J., 2021. Measuring massive multitask language understanding. *arXiv:2009.03300*.
- Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., de Las Casas, D., Hendricks, L.A., Welbl, J., Clark, A., et al., 2022. Training Compute-Optimal Large Language Models. *arXiv e-prints* .
- Huang, W., Liu, Y., Qin, H., Li, Y., Zhang, S., Liu, X., Magno, M., Qi, X., 2024. Billm: Pushing the limit of post-training quantization for llms. *arXiv preprint arXiv:2402.04291* .
- Inc., M., TsinghuaNLP, 2024. Minicpm: Unveiling the potential of end-side large language models.
- Jiang, A.Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D.S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L.R., Lachaux, M.A., Stock, P., Scao, T.L., Lavril, T., Wang, T., Lacroix, T., Sayed, W.E., 2023. Mistral 7b. *arXiv:2310.06825*.
- Jiang, A.Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford,

- C., Chaplot, D.S., de las Casas, D., Hanna, E.B., Bressand, F., et al., 2024. Mixtral of experts. [arXiv:2401.04088](https://arxiv.org/abs/2401.04088).
- Kim, S., Hooper, C., Gholami, A., Dong, Z., Li, X., Shen, S., Mahoney, M.W., Keutzer, K., 2024. Squeezellm: Dense-and-sparse quantization. URL: <https://arxiv.org/abs/2306.07629>, [arXiv:2306.07629](https://arxiv.org/abs/2306.07629).
- Kulal, S., Pasupat, P., Chandra, K., Lee, M., Padon, O., Aiken, A., Liang, P., 2019. Spoc: Search-based pseudocode to code. URL: <https://arxiv.org/abs/1906.04908>, [arXiv:1906.04908](https://arxiv.org/abs/1906.04908).
- Lancaster, T., 2023. Artificial intelligence, text generation tools and chatgpt—does digital watermarking offer a solution? *International Journal for Educational Integrity* 19, 10.
- Leswing, K., 2024. Nvidia ceo jensen huang announces new ai chips: ‘we need bigger gpus’. URL: <https://www.cnbc.com/2024/03/18/nvidia-announces-gb200-blackwell-ai-chip-launching-later-this-year.html>.
- Li, R., allal, L.B., Zi, Y., Muennighoff, N., Kocetkov, D., Mou, C., Marone, M., Akiki, C., LI, J., Chim, J., et al., 2023a. Starcoder: may the source be with you! *Transactions on Machine Learning Research* URL: <https://openreview.net/forum?id=KoF0g41haE>. reproducibility Certification.
- Li, S., Ning, X., Hong, K., Liu, T., Wang, L., Li, X., Zhong, K., Dai, G., Yang, H., Wang, Y., 2023b. Llm-mq: Mixed-precision quantization for efficient llm deployment. *The Efficient Natural Language and Speech Processing Workshop with NeurIPS* 9.

- Li, Y., Bubeck, S., Eldan, R., Giorno, A.D., Gunasekar, S., Lee, Y.T., 2023c. Textbooks are all you need ii: phi-1.5 technical report. [arXiv:2309.05463](https://arxiv.org/abs/2309.05463).
- Li, Y., Choi, D., Chung, J., Kushman, N., Schrittwieser, J., Leblond, R., Eccles, T., Keeling, J., Gimeno, F., Dal Lago, A., et al., 2022. Competition-level code generation with alphacode. *Science* 378, 1092–1097. URL: <http://dx.doi.org/10.1126/science.abq1158>, doi:10.1126/science.abq1158.
- Lin, J., Tang, J., Tang, H., Yang, S., Dang, X., Gan, C., Han, S., 2023. Awq: Activation-aware weight quantization for llm compression and acceleration. [arXiv:2306.00978](https://arxiv.org/abs/2306.00978).
- Lin, S., Hilton, J., Evans, O., 2021. Truthfulqa: Measuring how models mimic human falsehoods. [arXiv:2109.07958](https://arxiv.org/abs/2109.07958).
- Liu, C., Bao, X., Zhang, H., Zhang, N., Hu, H., Zhang, X., Yan, M., 2023a. Improving chatgpt prompt for code generation. [arXiv:2305.08360](https://arxiv.org/abs/2305.08360).
- Liu, J., Xia, C.S., Wang, Y., ZHANG, L., 2023b. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation. Thirty-seventh Conference on Neural Information Processing Systems URL: <https://openreview.net/forum?id=1qvx610Cu7>.
- Liu, J., Xia, C.S., Wang, Y., ZHANG, L., 2023c. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation, in: Oh, A., Neumann, T., Globerson, A., Saenko, K., Hardt, M., Levine, S. (Eds.), *Advances in Neural Information Processing Systems*, Curran Associates, Inc.. pp. 21558–21572. URL:

https://proceedings.neurips.cc/paper_files/paper/2023/file/43e9d647ccd3e4b7b5baab53f0368686-Paper-Conference.pdf.

Liu, Z., Oguz, B., Zhao, C., Chang, E., Stock, P., Mehdad, Y., Shi, Y., Krishnamoorthi, R., Chandra, V., 2023d. Llm-qat: Data-free quantization aware training for large language models. arXiv preprint arXiv:2305.17888 .

López Espejel, J., Ettifouri, E.H., Yahaya Alassan, M.S., Chouham, E.M., Dahhane, W., 2023a. Gpt-3.5, gpt-4, or bard? evaluating llms reasoning ability in zero-shot setting and performance boosting through prompts. *Natural Language Processing Journal* 5, 100032. URL: <https://www.sciencedirect.com/science/article/pii/S2949719123000298>, doi:<https://doi.org/10.1016/j.nlp.2023.100032>.

López Espejel, J., Yahaya Alassan, M.S., Chouham, E.M., Dahhane, W., Ettifouri, E.H., 2023b. A comprehensive review of state-of-the-art methods for java code generation from natural language text. *Natural Language Processing Journal* 3, 100013. URL: <https://www.sciencedirect.com/science/article/pii/S2949719123000109>, doi:<https://doi.org/10.1016/j.nlp.2023.100013>.

Ma, X., Fang, G., Wang, X., 2023a. Llm-pruner: On the structural pruning of large language models. *Advances in neural information processing systems* 36, 21702–21720.

- Ma, Y., Cao, Y., Sun, J., Pavone, M., Xiao, C., 2023b. Dolphins: Multimodal language model for driving. [arXiv:2312.00438](https://arxiv.org/abs/2312.00438).
- Manyika, J., 2023. An overview of bard: an early experiment with generative ai.
- Mukherjee, S., Mitra, A., Jawahar, G., Agarwal, S., Palangi, H., Awadallah, A., 2023. Orca: Progressive learning from complex explanation traces of gpt-4. [arXiv:2306.02707](https://arxiv.org/abs/2306.02707).
- Mulia, A.P., Piri, P.R., Tho, C., 2023. Usability analysis of text generation by chatgpt openai using system usability scale method. *Procedia Computer Science* 227, 381–388. URL: <https://www.sciencedirect.com/science/article/pii/S1877050923017040>, doi:<https://doi.org/10.1016/j.procs.2023.10.537>. 8th International Conference on Computer Science and Computational Intelligence (ICCSCI 2023).
- Omar, R., Mangukiya, O., Kalnis, P., Mansour, E., 2023. Chatgpt versus traditional question answering for knowledge graphs: Current status and future directions towards knowledge graph chatbots. *arXiv preprint arXiv:2302.06466* .
- OpenAI, 2022. Introducing chatgpt. <https://openai.com/blog/chatgpt>.
- OpenAI, 2023. Gpt-4 technical report. *arxiv* URL: <https://arxiv.org/pdf/2303.08774.pdf>.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C.L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J.,

- Kelton, F., Miller, L.E., Simens, M., Aspell, A., Welinder, P., Christiano, P.F., Leike, J., Lowe, R.J., 2022. Training language models to follow instructions with human feedback. ArXiv abs/2203.02155.
- Park, Y., Hyun, J., Cho, S., Sim, B., Lee, J.W., 2024. Any-precision llm: Low-cost deployment of multiple, different-sized llms. URL: <https://arxiv.org/abs/2402.10517>, arXiv:2402.10517.
- Plevris, V., Papazafeiropoulos, G., Jiménez Rios, A., 2023. Chatbots put to the test in math and logic problems: A comparison and assessment of chatgpt-3.5, chatgpt-4, and google bard. AI 4, 949–969.
- Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., 2018a. Improving language understanding by generative pre-training. arxiv .
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., 2018b. Language models are unsupervised multitask learners. arxiv URL: <https://d4mucfpksywv.cloudfront.net/better-language-models/language-models.pdf>.
- Rae, J.W., Borgeaud, S., Cai, T., Millican, K., Hoffmann, J., Song, F., Aslanides, J., Henderson, S., Ring, R., Young, S., et al., 2021. Scaling Language Models: Methods, Analysis & Insights from Training Gopher. arXiv e-prints , arXiv:2112.11446doi:10.48550/arXiv.2112.11446, arXiv:2112.11446.
- Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., Manning, C.D., Finn, C., 2023. Direct preference optimization: Your language model is secretly a reward model. arXiv:2305.18290.

- Ramachandran, P., Zoph, B., Le, Q.V., 2017. Swish: a self-gated activation function. arXiv: Neural and Evolutionary Computing URL: <https://api.semanticscholar.org/CorpusID:196158220>.
- Ramírez, B.G., Espejel, J.L., del Carmen Santiago Díaz, M., Linares, G.T.R., 2024. Sólo escúchame: Spanish emotional accompaniment chatbot. URL: <https://arxiv.org/abs/2408.01852>, arXiv:2408.01852.
- Rozière, B., Gehring, J., Gloeckle, F., Sootla, S., Gat, I., Tan, X.E., Adi, Y., Liu, J., et al., 2024. Code llama: Open foundation models for code. arXiv:2308.12950.
- Scholak, T., Schucher, N., Bahdanau, D., 2021. PICARD: Parsing incrementally for constrained auto-regressive decoding from language models, in: Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics. pp. 9895–9901. URL: <https://aclanthology.org/2021.emnlp-main.779>.
- Shazeer, N., 2020. Glu variants improve transformer. arXiv:2002.05202.
- Smith, S., Patwary, M., Norick, B., LeGresley, P., Rajbhandari, S., Casper, J., Liu, Z., Prabhunoye, S., Zerveas, G., Korthikanti, V., et al., 2022. Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model. arXiv:2201.11990.
- Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., Hashimoto, T.B., 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.

Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A.M., Hauth, A., et al., K.M., 2023. Gemini: A family of highly capable multimodal models. [arXiv:2312.11805](https://arxiv.org/abs/2312.11805).

Team, G., Mesnard, T., Hardin, C., Dadashi, R., Bhupatiraju, S., Pathak, S., Sifre, L., Rivière, M., Kale, M.S., Love, J., et al., 2024. Gemma: Open models based on gemini research and technology. [arXiv:2403.08295](https://arxiv.org/abs/2403.08295).

Team, H., 2024a. Evalplus leaderboard. <https://evalplus.github.io/leaderboard.html>.

Team, H., 2024b. Magicoder-oss-instruct-75k. <https://huggingface.co/datasets/ise-uiuc/Magicoder-OSS-Instruct-75K>.

Teknium, 2023. Openhermes 2.5 mistral 7b - gguf. URL: <https://huggingface.co/teknium/OpenHermes-2.5-Mistral-7B>.

Thoppilan, R., Freitas, D., Hall, J., Shazeer, N., Kulshreshtha, A., Cheng, H.T., Jin, A., Bos, T., Baker, L., Du, Y., Li, Y., Lee, H., et al., 2022. Lamda: Language models for dialog applications. [arXiv](https://arxiv.org/abs/2201.08230) .

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., Lample, G., 2023a. Llama: Open and efficient foundation language models. [arXiv:2302.13971](https://arxiv.org/abs/2302.13971).

Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al., 2023b. Llama 2: Open foundation and fine-tuned chat models. [arXiv:2307.09288](https://arxiv.org/abs/2307.09288).

- Tunstall, L., Beeching, E., Lambert, N., Rajani, N., Rasul, K., Belkada, Y., Huang, S., von Werra, L., Fourier, C., Habib, N., et al., 2023. Zephyr: Direct distillation of lm alignment. [arXiv:2310.16944](https://arxiv.org/abs/2310.16944).
- Vailshery, L.S., 2024. Most used programming languages among developers worldwide as of 2023. URL: <https://www.statista.com/statistics/793628/worldwide-developer-survey-most-used-languages/>.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.u., Polosukhin, I., 2017. Attention is all you need. *Advances in Neural Information Processing Systems* 30.
- Wang, B., Komatsuzaki, A., 2021. Gpt-j-6b: A 6 billion parameter autoregressive language model. <https://github.com/kingoflolz/mesh-transformer-jax>. URL: <https://github.com/kingoflolz/mesh-transformer-jax>.
- Wang, Y., Wang, W., Joty, S., Hoi, S.C., 2021. CodeT5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing* .
- Wei, J., Tay, Y., Bommasani, R., Raffel, C., Zoph, B., Borgeaud, S., Yogatama, D., Bosma, M., Zhou, D., Metzler, D., et al., 2022a. Emergent abilities of large language models. [arXiv:2206.07682](https://arxiv.org/abs/2206.07682).
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., hsin Chi, E.H., Xia, F., Le, Q., Zhou, D., 2022b. Chain of thought prompting elicits reasoning

- in large language models. ArXiv abs/2201.11903. URL: <https://api.semanticscholar.org/CorpusID:246411621>.
- Wei, Y., Wang, Z., Liu, J., Ding, Y., Zhang, L., 2023. Magicoder: Source code is all you need. arXiv preprint arXiv:2312.02120 .
- Wong, M.F., Guo, S., Hang, C.N., Ho, S.W., Tan, C.W., 2023. Natural language generation and understanding of big code for ai-assisted programming: A review. *Entropy* 25, 888.
- Wu, X., Duan, R., Ni, J., 2023. Unveiling security, privacy, and ethical concerns of chatgpt. *Journal of Information and Intelligence* URL: <https://www.sciencedirect.com/science/article/pii/S2949715923000707>, doi:<https://doi.org/10.1016/j.jiixd.2023.10.007>.
- Xiao, G., Lin, J., Seznec, M., Wu, H., Demouth, J., Han, S., 2023. Smoothquant: Accurate and efficient post-training quantization for large language models, in: *International Conference on Machine Learning*, PMLR. pp. 38087–38099.
- Yu, Z., He, L., Wu, Z., Dai, X., Chen, J., 2023. Towards better chain-of-thought prompting strategies: A survey. arXiv:2310.04959.
- Zhang, B., Sennrich, R., 2019. Root mean square layer normalization. arXiv:1910.07467.
- Zheng, L., Chiang, W.L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., et al., E.P.X., 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. arXiv:2306.05685.

Zhong, W., Cui, R., Guo, Y., Liang, Y., Lu, S., Wang, Y., Saied, A., Chen, W., Duan, N., 2023. Agieval: A human-centric benchmark for evaluating foundation models. [arXiv:2304.06364](https://arxiv.org/abs/2304.06364).

Zhou, Y., Muresanu, A.I., Han, Z., Paster, K., Pitis, S., Chan, H., Ba, J., 2023. Large language models are human-level prompt engineers. The Eleventh International Conference on Learning Representations URL: <https://openreview.net/forum?id=92gvk82DE->.

Zhuang, Y., Yu, Y., Wang, K., Sun, H., Zhang, C., 2024. Toolqa: A dataset for llm question answering with external tools. Advances in Neural Information Processing Systems 36.

Appendix A. Prompt Engineering

As discussed in Section 5, each evaluated model requires a customized template that aligns with its specific tokens for beginning and ending a turn, as well as its role management. For example, the Mistral-7B-Instruct-v0.2 model uses two roles, “user” and “assistant”. In this model, the conversation starts with the user’s input, followed by the assistant’s response. Table A.15 shows an example of the final template used as input for Mistral-7B-Instruct-v0.2, where the prompt begins with the special tokens “`[INST]`”, with the user prompt text highlighted in teal, and ends with “`[/INST]`”.

In contrast, models like Llama-3.1, Zephyr, and Dolphin use three roles: “system”, “user” and “assistant”. Unlike the template in Table Table A.15, these models require the prompt to be split, with part of it (highlighted in olive) assigned to the system role, and the rest (in teal) used as the user prompt. Additionally, each of these models uses different special tokens to mark the beginning and end of roles, highlighted in pink. Refer to Table A.16, Table A.18 and Table A.17 for the complete prompts used with Llama-3.1, Zephyr, and Dolphin, respectively.

```
<s> [INST] You are an expert Python programmer. You will receive
three keywords: problem, variables, and options.
Your task consists of resolving the problem statement using Python
code. There are two important keys:

1. The Python code should use the variable names given in the
'Variables'.
2. The Python code should return the string of the option items.

Example 1:
Problem: If the sum of the price and VAT exceeds 10, click on the
Cancel button. Otherwise, click on the Accept button.
Variables: price, VAT, total
Options: a) click on Cancel button b) click on Accept button

Code:
def compute_sum(price, VAT):
    if price + VAT > 10:
        return "click on Cancel button"
    else:
        return "click on Accept button"

Write Python code given:
Problem: $problem,
Variables: $variables,
Options: $options

Provide only Python code; do not include any additional comments.
[/INST]
```

Figure A.15: Prompt template for Mistral-7B-Instruct-v0.2 model. **Note:** Text highlighted in pink represents special tokens, while teal indicates the user's input.

```

<|begin_of_text|><|start_header_id|>system<|end_header_id|>
You are an expert Python programmer. You will receive three
keywords: problem, variables, and options.
Your task consists of resolving the problem statement using Python
code. There are two important keys:

1. The Python code should use the variable names given in the
'Variables'.
2. The Python code should return the string of the option items.

Example 1:
Problem: If the sum of the price and VAT exceeds 10, click on the
Cancel button. Otherwise, click on the Accept button.
Variables: price, VAT, total
Options: a) click on Cancel button    b) click on Accept button

Code:
def compute_sum(price, VAT):
    if price + VAT > 10:
        return "click on Cancel button"
    else:
        return "click on Accept button"

<|eot_id|><|start_header_id|>user<|end_header_id|>
Write Python code given:
Problem: $problem,
Variables: $variables,
Options: $options

Provide only Python code; do not include any additional comments.
<|eot_id|><|start_header_id|>assistant<|end_header_id|>

```

Figure A.16: Prompt template for Meta-Llama-3.1-8B-Instruct model. **Note:** Text highlighted in pink represents special tokens, olive indicates the system's input, and teal represents the user's input.

```
<|system|>
You are an expert Python programmer. You will receive three
keywords: problem, variables, and options.
Your task consists of resolving the problem statement using Python
code. There are two important keys:

1. The Python code should use the variable names given in the
'Variables'.
2. The Python code should return the string of the option items.

Example 1:
Problem: If the sum of the price and VAT exceeds 10, click on the
Cancel button. Otherwise, click on the Accept button.
Variables: price, VAT, total
Options: a) click on Cancel button b) click on Accept button

Code:
def compute_sum(price, VAT):
    if price + VAT > 10:
        return "click on Cancel button"
    else:
        return "click on Accept button"

</s>
<|user|>
Write Python code given:
Problem: $problem,
Variables: $variables,
Options: $options

Provide only Python code; do not include any additional comments.
</s>
<|assistant|>
```

Figure A.17: Prompt template for zephyr-7b-beta model. **Note:** Text highlighted in pink represents special tokens, olive indicates the system's input, and teal represents the user's input.

```

<s>[INST] <<SYS>>
You are an expert Python programmer. You will receive three
keywords: problem, variables, and options.
Your task consists of resolving the problem statement using Python
code. There are two important keys:

1. The Python code should use the variable names given in the
'Variables'.
2. The Python code should return the string of the option items.

Example 1:
Problem: If the sum of the price and VAT exceeds 10, click on the
Cancel button. Otherwise, click on the Accept button.
Variables: price, VAT, total
Options: a) click on Cancel button b) click on Accept button

Code:
def compute_sum(price, VAT):
    if price + VAT > 10:
        return "click on Cancel button"
    else:
        return "click on Accept button"

<</SYS>>
Write Python code given:
Problem: $problem,
Variables: $variables,
Options: $options

Provide only Python code; do not include any additional comments.
[/INST]

```

Figure A.18: Prompt template for dolphin-2.6-mistral-7b model. **Note:** Text highlighted in pink represents special tokens, olive indicates the system's input, and teal represents the user's input.