

LLaMAntino 3 8B: Further Fine-Tuning on Italian Language Using Consumer Hardware

Bruno Guzzo

Fall 2024

Abstract

This project delves into the fine-tuning of a large language model (LLMs) for enhanced Italian language proficiency and its application in a chatbot. By employing 4-bit quantization, gradient checkpointing, and low-rank adaptation, we efficiently fine-tune the model on consumer hardware using well established Italian dataset. This method enables us to optimize the model's performance while minimizing computational overhead. To assess the impact of our fine-tuning strategy, we conduct a comprehensive evaluation using a series of Italian benchmarks, including those for question answering, commonsense reasoning, and natural language understanding. The project further explores the potential of incorporating a Retrieval-Augmented Generation (RAG) approach, based on a curated collection of Italian literature and novels, to enrich the chatbot's knowledge base and enhance its conversational capabilities.

1 Introduction

Large Language Models (LLMs) have rapidly emerged as a transformative force in the field of artificial intelligence, demonstrating remarkable capabilities in understanding and generating human-like text. These models, trained on massive datasets of text and code, leverage deep learning techniques to learn complex patterns and relationships within natural language. At the core of their architecture lie several fundamental components:

- **Tokenization:** The process of breaking down text into smaller units, called tokens, which can be words, subwords, or even characters. This allows the model to process and represent the input text in a structured manner.
- **Embedding:** Each token is then converted into a numerical vector, or embedding, which captures its semantic meaning in a high-dimensional space. This representation allows the model to understand relationships between words and their contextual significance.

- **Attention Mechanism:** This innovative mechanism, introduced in the seminal paper "Attention is All You Need" (Vaswani et al., 2017), enables the model to weigh the importance of different parts of the input sequence when processing information. This allows LLMs to effectively capture long-range dependencies and contextual nuances in text, a significant advancement over previous sequential models.
- **Transformers:** The dominant architecture for LLMs, transformers consist of multiple layers of self-attention and feed-forward neural networks. This structure allows for parallel processing of information, making them highly efficient and scalable for handling large datasets.

The evolution of LLMs over the past decade has been marked by significant advancements in model size, training data, and architectural innovations. Early models, such as recurrent neural networks (RNNs) and long short-term memory (LSTM) networks, faced limitations in capturing long-range dependencies and processing large amounts of data. The introduction of the transformer architecture revolutionized the field, enabling the development of models with billions of parameters. These models have demonstrated impressive capabilities in various tasks, including text generation, translation, question answering, and code generation, sparking widespread interest and research in both academia and industry.

2 The Base Model: Meta LLaMA

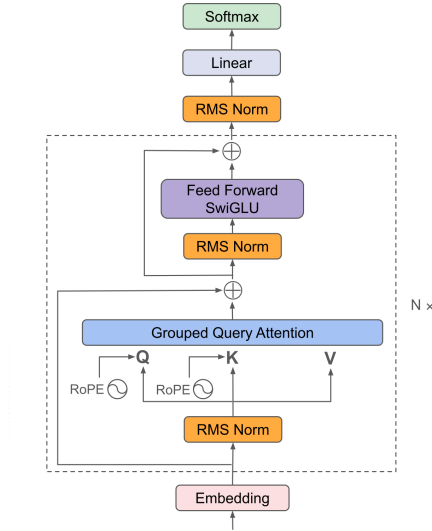


Figure 1: LLaMa model architecture

The LLaMA model is a series of large language models (LLMs) ranging from 7B to 65B parameters, trained on trillions of tokens. The goal of LLaMA is to achieve the best possible performance at various inference budgets. The model architecture is based on the transformer architecture, with the following modifications:

- **Pre-normalization:** The input of each transformer sub-layer is normalized, rather than normalizing the output. The **RMSNorm** normalizing function is used. RMSNorm simplifies the LayerNorm (used by transformer) by removing the re-centering operation, relying only on re-scaling. It normalizes the summed inputs to a neuron using the Root Mean Square (RMS) statistic.

$$\bar{a}_i = \frac{a_i}{\text{RMS}(\mathbf{a})} g_i, \quad \text{where } \text{RMS}(\mathbf{a}) = \sqrt{\frac{1}{n} \sum_{i=1}^n a_i^2}. \quad (1)$$

- **SwiGLU activation function:** The ReLU non-linearity is replaced with the SwiGLU (Swish-Gated Linear Unit) activation function. It is a combination of the **Swish activation function** and the **Gated Linear Unit (GLU)**.

$$\text{Swish}(x) = x \cdot \sigma(x), \quad \text{GLU}(x) = x \cdot \sigma(Wx + b) \quad (2)$$

$$\text{SwiGLU}(x) = x \cdot \sigma(\beta \cdot x) + (1 - \sigma(\beta \cdot x)) \cdot (Wx + b) \quad (3)$$

Where \mathbf{W} , \mathbf{b} and β are trainable parameters.

SwiGLU is **non-monotonic**, this allows it to capture complex non-linear relationships between input and output.

- **Rotary Position Embeddings:** Absolute positional embeddings are removed, and instead rotary positional embeddings (RoPE) are used at each layer of the network to compute query and key. **RoPE** is based on the idea of using a rotation matrix to control the inner product between query and value vectors. Let $x_m \in \mathbb{R}^d$ and $x_n \in \mathbb{R}^d$ be two word embeddings at positions m and n , respectively. Let q_m and k_n be the corresponding query and key vectors. RoPE encodes the relative position information by rotating q_m and k_n by an angle that is proportional to $m - n$. The rotated vectors are then used to compute the inner product. Then, **the**

rotation matrix is defined as follows:

$$R_{\Theta, m}^d = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{d/2} & -\sin m\theta_{d/2} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{d/2} & \cos m\theta_{d/2} \end{pmatrix} \quad (4)$$

Where $\Theta = \{\theta_i = 10000^{-2(i-1)/d}, i \in [1, 2, \dots, d/2]\}$ is a set of pre-defined parameters. The inner product of the rotated vectors is then computed as follows:

$$q_m^T k_n = (R_{\Theta, m}^d W_q x_m)^T (R_{\Theta, n}^d W_k x_n) = x_m^T W_q R_{\Theta, n-m}^d W_k x_n \quad (5)$$

Where $R_{\Theta, n-m}^d = (R_{\Theta, m}^d)^T R_{\Theta, n}^d$.

This inner product decays with increasing relative distance $|m - n|$. This is due to the rotation matrix $R_{\Theta, n-m}^d$ becoming more diagonal and thus decreasing the overall cosine value of the angle between q and k vectors.

params	dimension	n heads	n layers	learning rate	batch size	n tokens
6.7B	4096	32	32	$3.0e^{-4}$	4M	1.0T
13.0B	5120	40	40	$3.0e^{-4}$	4M	1.0T
32.5B	6656	52	60	$1.5e^{-4}$	4M	1.4T
65.2B	8192	64	80	$1.5e^{-4}$	4M	1.4T

Table 1: **LLaMA Model sizes, architectures, and optimization hyper-parameters**

3 LLaMA 3

LLaMA 3 is a series of large language models (LLMs) developed by Meta AI. It is the successor to the LLaMA model, featuring significant enhancements in terms of scale, performance, and capabilities. A key improvement from LLaMA to LLaMA 3 is the adoption of the **Grouped Query Attention** (GQA). It is a type of attention mechanism that groups the queries into a smaller number of groups where each group shares a single key and value. This strategy reduces the number of attention weights that need to be computed leading to faster training and computation.

Feature	LLaMA	LLaMA 3
<i>Normalization</i>	RMSNorm (before self-attention)	RMSNorm (before self-attention)
<i>Activation Function</i>	SwiGLU	SwiGLU
<i>Positional Embedding</i>	RoPE	RoPE
<i>Attention</i>	Standard Multi-Head Attention	Grouped Query Attention (GQA)
<i>Attention Mask</i>	No attention mask	Attention mask to prevent self-attention between different documents in the same sequence
<i>Tokenizer</i>	SentencePiece Byte Pair Encoding (BPE)	SentencePiece BPE with Tiktoken vocabulary
<i>Vocabulary Size</i>	32,000 tokens	128,000 tokens
<i>Context Window</i>	2048 tokens	128K tokens

Table 2: **LLaMA 3 and LLaMA key differences**

LLaMA 3 undergoes a post-training process that aligns it with human feedback. This process involves **Supervised Fine-Tuning** (SFT) and **Direct Preference Optimization** (DPO). SFT is conducted on a mixture of human-annotated and synthetically-generated data, while DPO utilizes human-annotated preference data.

4 LLaMAntino 3 8B

LLaMAntino 3 8B 'ANITA' is Italian language LLM based on the Meta’s LLaMA 3 8 billion parameters model. It is developed and maintained by the **SWAP Team** of the Bari University, Italy. It is fine-tuned using a combination of **Supervised Fine-Tuning** (SFT) and **Dynamic Preference Optimization** (DPO) to achieve superior performance in the Italian language.

The fine-tuning process involved two main steps:

- **Model Supervised Fine-Tuning:** The starting point is the **meta-llama/Meta-Llama-3-8B-Instruct** model, an instruction tuned version of the base LLaMa 3 8B model. Then, it has been fine-tuned using a dataset named **Chat-Error/wizard_alpaca_dolly_orca**, this dataset is a combination of three well-known datasets: WizardLM.Orca, Dolly-v2.Orca, and Alpaca.Orca. It consists of 100K prompts organized into system, instruction, input, and output fields in English language. The

model was fine-tuned using the **Unsloth framework** on an NVIDIA H100 64GB GPU card.

- **Model Direct Preferences Optimization:** After SFT, DPO was applied to align the model’s outputs with human preferences. For this purpose, the **mlabonne/orpo-dpo-mix-40k** dataset has been used, which contains approximately 40k examples from various sources like **Capybara-Preferences**, **distilabel-intel-orca-dpo-pairs**, and **ultrafeedback-binarized-preferences-cleaned**. The DPO process was also conducted using Unsloth on an NVIDIA H100 64GB GPU card.

To adapt the model to the Italian language, we used the **gsarti/clean_mc4_it** dataset, which is a cleaned version Italian split of the multilingual colossal Common Crawl’s **Web Crawl Corpus** (mc4). This dataset has been used with the same SFT strategy as before, formatting the prompts according to the standard LLaMA 3 template.

Then, the model has been evaluated using several benchmarks, including **MMLU**, **HellaSwag**, **A12 Reasoning Challenge** (arc_c), TruthfulQA, Winogrande, and GSM8K. The results showed outstanding performance compared to similar and larger models, demonstrating the effectiveness of our fine-tuning strategy.

Metric	Value
Avg.	0.6160
Arc.IT	0.5714
Hellaswag.IT	0.7093
MMLU.IT	0.5672

Table 3: **LLaMAntino 3 ‘ANITA’ Main Benchmark Result**

5 Optimizations

Prior to running fine-tuning experiments on LLaMAntino 3 model various optimization strategy is needed. This strategies ensure that the model can run and train on a consumer based hardware without causing out of memory errors or being too slow. The optimization strategy used is the following.

5.1 Quantization

Large neural networks, are computationally expensive and memory-intensive. This makes them difficult to deploy on resource-constrained devices. Quantization addresses this by reducing the numerical precision of the network’s weights

and activations, typically from 32-bit floating-point to lower bit-depth representations like 8-bit or 4-bit integers. The core idea behind quantization is to map the range of real-valued weights and activations to a smaller set of quantized values. A common approach is to use an affine mapping:

$$r = S(q - Z) \quad (6)$$

Where \mathbf{r} is the real value, \mathbf{q} is the quantized value, \mathbf{S} is the scale factor and \mathbf{Z} is the zero-point.

To minimize accuracy loss, quantized training simulates the effects of quantization during the training process. This allows the network to adapt to the lower precision representation. During training, weights and activations are quantized in the forward pass, but gradients are still computed and updated in floating-point.

5.2 Low Rank Adaptation (LoRA)

The core idea of Low-Rank Adaptation (LoRA) is rooted in the observation that the changes in a model’s weights during adaptation to a new task often exhibit a low rank. This means that the updates to the weight matrices can effectively be represented by a much smaller set of basis vectors. LoRA leverages this property to significantly reduce the number of trainable parameters during adaptation.

Consider a **pre-trained weight matrix** $W \in \mathbb{R}^{d \times k}$ within a neural network layer, such as a fully connected layer or an attention head in a Transformer model. During adaptation, this weight matrix is typically updated through gradient descent, resulting in a change ΔW . LoRA posits that ΔW can be approximated by a **low-rank decomposition**:

$$\Delta W \approx BA \quad (7)$$

Where $A \in \mathbb{R}^{r \times k}$ and $B \in \mathbb{R}^{d \times r}$ are matrices with $r \ll \min(d, k)$. Here, r represents the **rank of the decomposition**, which is assumed to be much smaller than the original dimensions of the weight matrix. The forward pass through the layer is then modified to incorporate the low-rank update:

$$h = Wx + BAx = (W + BA)x \quad (8)$$

Where x is the input to the layer and h is the output and ΔWx is also scaled by $\frac{\alpha}{r}$, where α is a constant in r .

During training, **the original weight matrix W is frozen**, and only the matrices A and B are updated. Both A and B are initialized using a random Gaussian distribution. By training only the low-rank matrices, the number of trainable parameters is significantly reduced, leading to **substantial memory savings**.

5.3 Gradient Checkpointing

Gradient checkpointing is a technique used to reduce the memory consumption of deep neural network training. It is based on the idea that instead of storing all intermediate activations of a neural network during the forward pass, only a subset of these activations need to be stored. The remaining activations can be recomputed during the backward pass, thereby trading off computation time for memory savings. Let's consider a neural network with n layers. We can denote the activations of the network at layer i as a_i . The forward pass of the network can be written as:

$$a_i = f_i(a_{i-1}, w_i) \quad (9)$$

Where f_i is the function computed by layer i and w_i are the weights of layer i . The backward pass of the network can be written as:

$$\frac{\partial L}{\partial a_{i-1}} = \frac{\partial L}{\partial a_i} \frac{\partial a_i}{\partial a_{i-1}} \quad (10)$$

Where L is the loss function. To compute the gradients of the loss function with respect to the activations, we need to store all the activations a_i during the forward pass. This can be **memory-intensive**, especially for deep neural networks. Gradient checkpointing addresses this issue by only storing a subset of the activations, called checkpoints. Let's denote the checkpoints as c_k , where k is the index of the checkpoint. The forward pass of the network with gradient checkpointing can be written as:

$$a_i = f_i(a_{i-1}, w_i) \text{ if } i \text{ is not a checkpoint} \quad (11)$$

$$c_k = a_i \text{ if } i \text{ is a checkpoint} \quad (12)$$

During the backward pass, the activations that are not checkpoints are recomputed from the nearest checkpoint ($a_i = f_i(a_{i-1}, w_i)$). The backward pass with gradient checkpointing can be written as:

$$\frac{\partial L}{\partial a_{i-1}} = \frac{\partial L}{\partial a_i} \frac{\partial a_i}{\partial a_{i-1}} \text{ if } i \text{ is not a checkpoint} \quad (13)$$

$$\frac{\partial L}{\partial c_k} = \frac{\partial L}{\partial a_i} \frac{\partial a_i}{\partial c_k} \text{ if } i \text{ is a checkpoint} \quad (14)$$

6 LLaMAntino 3 8B Further Fine-Tuning

The objective of this project is to test whether is possible to extend further the capabilities of LLaMAntino 3 model to generalize and handle Italian language tasks, such as chatting. Since the adoption of already discussed optimization strategy has proved to be needed run this kind of model locally, we run some of the same benchmark proposed by the authors to check how these optimizations impact the overall model performance. A description of those benchmarks is present in the appendix section. To save time, we choose to adopt a subset 0.20% of the proposed datasets used by the benchmarks chosen by the authors.

Metric	Original Value	Our Value
Arc.IT	0.5714	0.5256
Hellaswag.IT	0.7093	0.6775

Table 4: **LLaMAntino 3, 4-bit Quantization, 20% Dataset**

6.1 Italian Wikipedia Dataset

The first way we choose to try extend the model generalization capability and enhance its knowledge regarding the Italian language it to fine tune the model to the whole Italian Wikipedia dataset, available online. This first idea has been justified since part of this dataset - and its variants - has been commonly used in the training of successful LLMs like the Open AI’s GPT family. The dataset we choose has roughly 1.83M rows with field such as title, text and url, more details can be found here.

6.1.1 Fine-Tuning process

The fine-tuning has been made possible by the adoption of the optimizations described above and the overall fine-tuning for one epoch has took roughly 50 hours on a 8GB GPU. We choose a LoRA rank (r) equal to 32, this allowed us to train roughly 83M parameters i.e. 8.3% of the total model parameters.

Metric	Original Value	Our Value
Arc.IT	0.5256	0.4188
Hellaswag.IT	0.6775	0.5182

Table 5: **LLaMAntino 3, Full Wikipedia IT, $r = 32$, $\alpha = 8$**

Despite the resource employed there has been a sensible decrement of the chosen benchmark metrics. Our hypothesis is that due to the LoRA adaptation the number of trained parameter has not been sufficient. We have expected not to degrade the model performance at least for the ARC benchmark which is based on a grade-school level multiple-choice science questions, thus in our view such a fine-tune strategy should have improved the model on this field.

6.2 Empirical results

To gain a glimpse of how this fine-tune has changed the model we propose a set of prompts given to both non-tuned and tuned models. We have then preserved the same value of **top-p**=0.95 and **temperature**=0.6 along with a maximum

output size of 256 tokens.

Wikipedia IT Fine-Tuned Model:

1. **Prompt:** Scrivi una breve storia per bambini in rima sulla vita di Leonardo da Vinci, usando un linguaggio semplice e fantasioso.

Response: *Leonardo da Vinci nacque a Vinci, in Toscana, in una famiglia numerosa di artisti. Il padre fu un pittore e la madre una modella.*

Quando Leonardo era bambino, suo padre lo portò a Firenze, dove il piccolo si appassionò alla pittura, ai disegni e alle sculture. Il padre lo fece studiare pittura con un pittore locale, ma Leonardo preferì andare a Milano per studiare scultura con il famoso pittore Andrea del Verrocchio. Il maestro lo aiutò a diventare un ottimo scultore.

Ma Leonardo aveva anche un grande interesse per le macchine, le invenzioni e le scoperte scientifiche. In seguito si dedicò a molti altri interessi: architettura, musica, matematica, ingegneria, letteratura e medicina. Inoltre, fu un pioniere dell'aerostatica e del volo. Studiò l'aerostato e ne realizzò alcuni modelli. Nel 1495, fu il primo a volare, ma non riuscì a compiere un vol...

2. **Prompt:** Scrivi un sonetto in stile petrarchesco che esprima il dolore per la perdita di una persona cara.

Response: *In un giorno di tenebra, di lutto e di pianto
Nel quale l'Amore e la Memoria si fanno notare
La tua scomparsa, oh cara, è un dolore
Inestinguibile, che non lascia spazio a consolazione.*

*La tua bellezza, oh cara, è un fuoco
Che brucia nel mio cuore, senza spegnersi
La tua bocca, oh cara, è una fonte
Da cui sgorga un fiume di parole.*

*La tua voce, oh cara, è un suono
Che risuona in ogni parte del mio cuore
La tua presenza, oh cara, è un'ombra
Che mi segue ovunque vada.*

*In un giorno di lutto, di pianto e di dolore
Nel quale l'Amore e la Memoria si fanno notare
La tua scomparsa, oh cara, è un dolore
Inestinguibile, che non lascia spazio a consolazione.*

3. **Prompt:** Descrivi le principali differenze tra la fisica classica e la fisica quantistica, evidenziando le implicazioni filosofiche di queste differenze.

Response: *La fisica classica è la teoria scientifica che descrive il mondo*

fisico in termini di oggetti e processi macroscopici, come l'acqua, l'aria, la terra, i corpi celesti, la materia, le forze e le loro interazioni, in termini di strutture e processi discreti e deterministici. Il mondo fisico è considerato come un insieme di entità separate e indipendenti, con proprietà definite e determinate, e con leggi e processi che ne determinano l'evoluzione. Il mondo fisico è quindi considerato come un insieme di oggetti e processi, e le loro proprietà e processi sono determinati da leggi e da equazioni che descrivono le leggi della natura. La fisica classica è basata sulla rappresentazione di una realtà oggettiva che esiste indipendentemente dallo studio di essa.

La fisica quantistica, al contrario, è basata sulla rappresentazione di una realtà non oggettiva, che non può essere considerata come un insieme di entità separate e indipendenti. La realtà...

Original Model:

1. **Prompt:** Scrivi una breve storia per bambini in rima sulla vita di Leonardo da Vinci, usando un linguaggio semplice e fantasioso.

Response: *In Italia, un paese bello e nuovo,
Viveva un genio, di cui parleremo a due.
Leonardo, un nome famoso e grande,
Disegnava, dipingeva, e inventava in pale.*

*Volava con gli uccelli, in cielo alto e chiaro,
Studiava l'acqua, e come muove il vento in aria.
Fece un cavallo, di legno, in movimento reale,
E un'arma, per difendersi, senza ferire il reale.*

*Disegnò un'uomo, nudo, con precisione alta,
E dipinse una Madonna, con un sorriso soave e alta.
Inventò un'elica, e un aereo, in sogno li vide,
E scrisse, disegnò, e sognò, finché il destino lo diede.*

2. **Prompt:** Scrivi un sonetto in stile petrarchesco che esprima il dolore per la perdita di una persona cara.

Response: *In mortal vita, dove si va a no,
Troppo presto te seppi perdere, amica,
E con te, la luce che mi faceva no,
E la speranza, che in te riposava.*

*La tua memoria, in me ferma stays,
E con essa, il lamento e la pietra,
Che sul cuore, pesa, inestrimabile ways,
E la notte, che non ha più serenata.*

Tanto, in te, avevo la mia pace,

*E in te, la mia vita, era tutta inerte,
Ora, in te non c'è più, né spazio,
E la mia, è in tormento, inestirpata.*

3. **Prompt:** Descrivi le principali differenze tra la fisica classica e la fisica quantistica, evidenziando le implicazioni filosofiche di queste differenze.

Response: *Le differenze fondamentali:*

1. ***Determinismo vs. Indeterminismo**:* Fisica classica: eventi determinati da condizioni iniziali e leggi di natura; Fisica quantistica: eventi casuali, descrivibili solo da probabilità di esito.
2. ***Osservazione e misura**:* Fisica classica: proprietà oggettive, indipendenti dall'osservatore; Fisica quantistica: proprietà in relazione all'osservatore, collegate alla misura.
3. ***Località vs. Non-località**:* Fisica classica: azioni isolate, spazio e tempo separati; Fisica quantistica: azioni non-locali, entanglement e quantità di azione non spaziotemporalmente separate.

Implicazioni filosofiche:

La fisica quantistica sconvolge l'idea di oggettività e determinismo, suggerendo un ruolo attivo dell'osservatore.

La non-località quantistica sfida la nostra comprensione di spazio e tempo, sollevando questioni sulla natura della realtà.

6.3 Conclusion

The evidence suggests that the fine-tuned model has a **biased language style** toward a more formal and scientific one, as the one used in the Wikipedia pages. It is of particular note that even for the tasks where creativity and informality are required - e.g. prompt 1 and 2 - the fine-tuned model tend to preserve a more **comprehensible** and **semi-formal language**.

We can further infer that the original training set must have included a collection of creative and informal Italian novels. It is noteworthy as well that for the scientific tasks - i.e. prompt 3 - the fine-tuned model tend to follow a more **discursive pattern** where the original model tend to **break the information** in small chunks. Thus, the fine-tuned model is capable of generating text like the one we can found in scientific and academic books.

7 Testing Italian LLMs using INVALSI

Testing the language comprehension of an LLM is no easy task, particularly when we are dealing with numerous tasks possible. One way to test the general capability of a language model to comprehend a given language is to use the national reference test for the education system. For the Italian language, such a test are developed by INVALSI each year to test the ability of middle and

high school Italian students. In our case we choose to develop a small benchmark of **80 Italian language comprehension questions** extracted from the INVALSI tests of the years 2008, 2009 and 2010 for the middle school. The final dataset include three JSON file - each per year - containing a list of object with information regarding the context and a list of questions with the real answer. The process of extracting the text, questions and answers has been done using **Google Gemini Flash** by providing a custom prompt and the original INVALSI file.

7.1 Prompt And Metrics

Each question is feed into the model by defining a custom prompt which include some basics instruction, the context and the question.

Informazioni: *Rispondi alla seguente domanda di italiano.
Per domande a scelta multipla rispondi indicando esclusivamente la lettera.
Ad esempio: A, B o C.
Per le domande aperte rispondo i modo breve e conciso.*

Contesto: $\{test_context\}$

Domanda: $\{test_question\}$

Doing so each questions carries the right context and it is isolated from the others esuring no bias between the questions. The model generated response is evaluated by comparing it to the real one using the **ROUGE-1** metric (precision), then we extract the mean of this value across questions. We also used a **multi-run strategy**, this means that the final metric is the mean of the metrics obtained with multiple benchmark run with the same parameters. This allow us to stabilize the overall metrics by reducing the impact of model randomness. Typically, we use two to three run epochs for this benchmark.

7.2 Models comparison

To compare the original model and our fine-tuned version we choose to use the same value of **top-p**=0.95 and **temperature**=0.6 along with a maximum output size of 8. Such a small number of allowed output token is justified by the small answers size in the dataset. Then, we also choose to use **two epoch** for the benchmark due to time reason, though this is enough to reduce the randomness of the metrics. As seen before, to execute this benchmark we have used once again the 4-bit quantization to prevent out of memory errors. The results is listed in the table below.

LLaMA 3 8B	LLaMAntino 3 8B	Our Model
0.12875	0.183020	0.081875

Table 6: **INVALSI Language Comprehension**, 4-bit quantization

7.3 Conclusion

Despite a very **bad performance** of all three evaluated model, it is evident how our fine tuning strategy has not improved the model user this aspect. It is worth noting that none of the considered models perform well in this benchmark. In our view this behavior stems from how the model has been designed and trained, as shown before both the original LLaMAntino 3 model and our fine-tuned version has proven to be efficient in generating response with **robust context** but struggle when a very brief text is needed as response.

This behavior can certainly be moderated by an accurate fine-tuning process, like **reward training** to force the model in a certain type of output.

Extending further the idea of using INVALSI test as benchmark could lead to solid benchmark stragey to test how well the newer LLMs performs in understanding deeply the Italian language. In fact, this idea has been recently exploded in the Italian academic community leading to a public classification of the most common LLMs regarding the Italian language. More denials could be found here.

8 Model Implementation

We finally propose two implementation of our model, one as a chat-bot and the second in a RAG System in which it can extract information from a knowledge base of plain-text books and files.

8.1 RAG Implementation

Retrieval-Augmented Generation (RAG) is a powerful technique that enhances the capabilities of large language models (LLMs) by incorporating external knowledge sources. This approach addresses a fundamental limitation of LLMs: their reliance on the knowledge encoded within their parameters during training. By integrating external information, RAG systems can generate more accurate, informative, and contextually relevant responses, significantly improving the overall user experience.

Our implementation enrich the model knowledge using a library of plain-text files. Despite the possibility of using several documents and books, our implementation use only one book due to the limited VRAM availabe. However, it is capable of generating helpful information regarding the chosen topic. An

example can be found below.

Prompt: Speiga in breve il primo canto della divina commedia.

Response: *Il primo canto della Divina Commedia di Dante Alighieri è il canto in cui l'autore introduce il lettore nella sua opera, esortandolo a seguirlo nel viaggio verso il Paradiso terrestre e l'eternità. Il canto è anche l'unico nella Commedia in cui Dante si rivolge direttamente al lettore, esortandolo a partecipare alla sua esperienza di poeta, che ha come obiettivo la salvezza dell'anima.*

Retrieval-Augmented Generation is a promising approach that significantly enhances the capabilities of large language models. As research in this area continues to advance, we can expect RAG systems to play an increasingly important role in the development and deployment of AI-powered solutions.

8.2 Chat-Bot Implementation

LLMs can also be used as a chat-bot to help and entertain the user while providing quick and comprehensible information. This is especially important in context such as voice assistant and customer assistance. We propose a chat-bot based on our fine-tuned model capable of entertain a semi-real-time conversation about different topics.

Such speed is reached by limiting the maximum number of generated token to 64. Our implementation is also capable of storing the chat history using a chat template that involve using a system prompt to guide the model behavior.

The complete chat template include the system prompt, the user questions and the model responses listed in order. This enables the model to follow the conversation topic providing off-topic responses. We avoid to show it here due to complex form of the special token used.

9 Conclusion

Present your findings and data. Use tables and figures to illustrate your results.

Appendix

A ARC Dataset

B HellaSwag Dataset

References