

ALMA MATER STUDIORUM · UNIVERSITÀ DI BOLOGNA

---

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING  
ARTIFICIAL INTELLIGENCE

# An effective strategy for reducing the size of LLaMA-based Language Models

Supervisor:  
Prof. Francesco Conti

Presented by:  
Angelo Galavotti

Co-supervisor:  
Luca Bompani

Sessione I  
Anno Accademico 2024/2025



*Alla migliore madre del mondo,  
al miglior padre del mondo.*



## **Abstract**

Pending



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	The social impact of LLMs . . . . .	2
1.2	The environmental impact of LLMs . . . . .	3
1.3	Scope of this project . . . . .	3
1.4	Document structure . . . . .	3
<b>2</b>	<b>Background and Related Work</b>	<b>5</b>
2.1	Early language models . . . . .	5
2.2	The architecture of Transformers . . . . .	6
2.3	The structure of Large Language Models . . . . .	7
2.4	A look at the target hardware . . . . .	7
2.5	Relevant compression and optimization techniques . . . . .	7
<b>3</b>	<b>Methodology</b>	<b>9</b>
3.1	Preliminary research: Franken-LLaMA . . . . .	9
3.2	Pruning methods . . . . .	9
3.2.1	Depth-wise pruning . . . . .	9
3.2.2	Width-wise pruning (WANDA-based) . . . . .	9
3.2.3	LoRA . . . . .	9
3.3	Quantization and EoRA . . . . .	9
3.4	Evaluation criteria . . . . .	9
<b>4</b>	<b>Implementation</b>	<b>11</b>
4.1	Implementation of the Pipeline . . . . .	11

---

4.1.1	Depth-wise Pruning . . . . .	11
4.1.2	Width-wise Pruning (WANDA-based) . . . . .	11
4.1.3	LoRA . . . . .	11
4.1.4	Quantization and EoRA . . . . .	11
4.2	Evaluation . . . . .	11
<b>5</b>	<b>Experimental Results and Analysis</b>	<b>13</b>
5.1	Preliminary observation of the results . . . . .	13
5.2	Results on TriviaQA . . . . .	13
5.3	Results on WikiText . . . . .	13
<b>6</b>	<b>Conclusion and Future Work</b>	<b>15</b>
6.1	Future work . . . . .	15
6.2	Final remarks . . . . .	15
	<b>References</b>	<b>17</b>



# List of Figures

2.1	A visual representation of the LSTM architecture. The three gates (input, forget, and output) control the flow of information in and out of the cell state, allowing it to retain relevant information over long sequences. . . . .	6
2.2	On the left, a standard Transformer architecture, which consists of an encoder and a decoder. On the right, we can observe the same architecture, while highlighting the only the layers used by a GPT-style model, which is a decoder-only transformer. The main difference is that the decoder does not attend to the encoder's output, allowing for auto-regressive generation. . . . .	8



# Chapter 1

## Introduction

It is no secret that, in the last three years, *Large Language Models* (LLMs) have fundamentally transformed our relationship with technology. Their impact rivals the most significant innovations of the past century, such as the internet and the smartphone. When people contemplate Artificial Intelligence today, they immediately think of ChatGPT or Claude, which have seamlessly integrated into our daily routines. Yet these powerful tools come with significant environmental concerns. Their development and operation consume vast amounts of energy and water resource: modern data centers supporting these models require extensive cooling systems and electricity consumption that can rival small cities.

The computational complexity of these systems necessitates cloud-based deployment, which not only amplifies their environmental footprint but also fundamentally restricts user autonomy. This cloud dependency creates a concerning power dynamic where users have limited control over their tools, while simultaneously enabling extensive data collection practices and potential surveillance mechanisms that would be impossible with local, user-controlled alternatives.

In this more conversational opening chapter we will briefly examine the environmental and social impact of LLMs while highlighting the growing imperative for efficient, locally-deployable models that democratize access

without depleting our planet's resources. Afterwards, we will outline the scope of this project, which aims to explore the potential of compression techniques to make LLMs more efficient.

The future of AI depends not just on what these models can do, but how sustainably they can do it.

## 1.1 The social impact of LLMs

The widespread adoption of LLMs has created ripple effects across virtually every sector of society, fundamentally altering how we work, learn, and create. In education, these tools have sparked heated debates about academic integrity while simultaneously offering new possibilities for personalized learning and accessibility for students with disabilities [5]. The workplace has experienced perhaps the most dramatic shifts, with entire professions grappling with automation anxiety while others discover unprecedented productivity gains. Creative industries find themselves in particularly complex territory: writers, artists, and content creators must navigate between leveraging AI as a collaborative tool and protecting their intellectual property from being absorbed into training datasets without consent or compensation [6].

What strikes me most profoundly is how these models have democratized access to sophisticated capabilities that were once the exclusive domain of experts. A small business owner can now generate marketing copy that rivals professional agencies, students can receive tutoring in subjects where human expertise might be scarce, and non-programmers can write functional code with natural language instructions. Yet this democratization comes with a troubling caveat: it's entirely dependent on maintaining access to centralized, corporate-controlled systems. When OpenAI experiences an outage, millions of users worldwide suddenly lose access to tools they've integrated into their daily workflows. When pricing models change, entire business models built around AI assistance can become unsustainable overnight.

This dependency becomes even more concerning when we consider the data these systems collect. Every interaction, every query, every creative prompt may potentially become part of these companies' datasets, raising questions about privacy and intellectual property. As such, the need for transparency and user control over these systems has never been more urgent, and I personally believe that the future of AI must prioritize local, user-deployable models that empower individuals rather than centralizing power in the hands of a few corporations.

## **1.2 The environmental impact of LLMs**

### **1.3 Scope of this project**

### **1.4 Document structure**

Pending for later.



# Chapter 2

## Background and Related Work

Before explaining the details and implementation of the methodology used in this project, it is essential to provide an overview of the evolution of the inner workings of the Transformer architecture as well as Language Models in general. In addition, we will also discuss relevant compression techniques that have been developed in this context, and how they influenced this work. Finally, we will also shed some light on the target hardware, whose limitations have been a driving force behind the design choices made in this project.

### 2.1 Early language models

Before the Transformer architecture, language models were typically based on recurrent neural networks (RNNs). They consist in an extension of the Multi-Layer Perceptron (MLP) in which by integrating cycles, allowing for the network to take into account the previous inputs when making a prediction. For these reasons, they can learn long-range dependencies in sequential data, such as text.

However, RNNs can be very much prone to the vanishing gradient problem ([9]), which makes it difficult for them to train in the long run, as during backpropagation the contribution of the gradients from earlier layers diminishes exponentially.

This problem motivated the introduction of the Long Short-Term Memory (LSTM) units, which uses gated cells in order to store more information about the input. In particular, there are three types gate cells:

- *Input gate*: controls how much of the new input should be added to the cell state.
- *Forget gate*: determines how much of the previous cell state should be retained.
- *Output gate*: decides how much of the cell state should be outputted to the next layer.

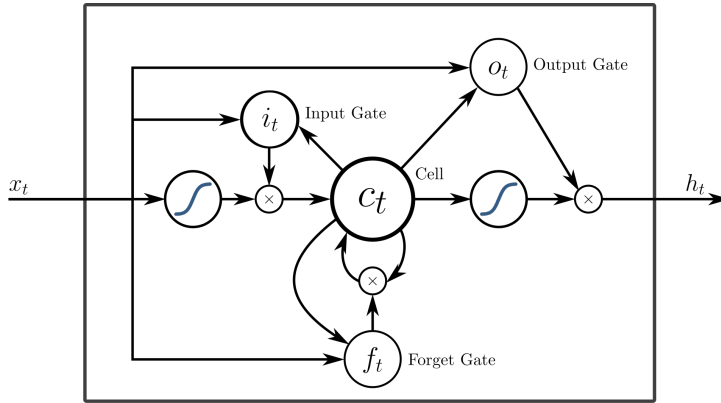


Figure 2.1: A visual representation of the LSTM architecture. The three gates (input, forget, and output) control the flow of information in and out of the cell state, allowing it to retain relevant information over long sequences.

## 2.2 The architecture of Transformers

The Transformer architecture, introduced by Vaswani et al. (2017) [TODO ADD REF HERE], represents a paradigm shift in sequence modeling, moving away from recurrent and convolutional approaches toward a purely attention-based mechanism. This innovation addressed fundamental limitations of previous architectures, particularly the sequential processing bottleneck that



hindered parallelization during training. **Attention Mechanism Evolution**  
The development of attention mechanisms began with Bahdanau et al. (2014), who introduced additive attention to improve neural machine translation by allowing the decoder to focus on relevant parts of the input sequence. This was refined by Luong et al. (2015) with multiplicative attention, which offered computational advantages. However, the breakthrough came with Vaswani et al. (2017), who demonstrated that attention mechanisms alone, without recurrence or convolution, could achieve state-of-the-art results across multiple tasks. The core innovation lies in the scaled dot-product attention mechanism: where Q, K, and V represent queries, keys, and values respectively, and  $d_k$  is the dimension of the key vectors. This formulation enables efficient computation while capturing long-range dependencies without the sequential constraints of RNNs.

## 2.3 The structure of Large Language Models

## 2.4 A look at the target hardware

## 2.5 Relevant compression and optimization techniques

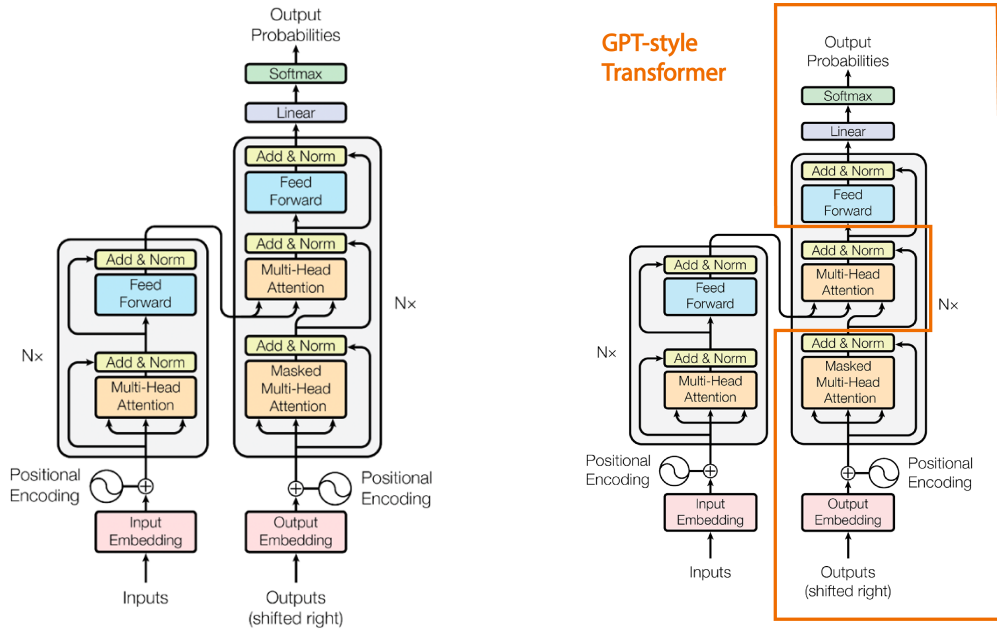


Figure 2.2: On the left, a standard Transformer architecture, which consists of an encoder and a decoder. On the right, we can observe the same architecture, while highlighting the only the layers used by a GPT-style model, which is a decoder-only transformer. The main difference is that the decoder does not attend to the encoder's output, allowing for auto-regressive generation.

# Chapter 3

## Methodology

Intro here

### 3.1 Preliminary research: Franken-LLaMA

### 3.2 Pruning methods

#### 3.2.1 Depth-wise pruning

#### 3.2.2 Width-wise pruning (WANDA-based)

#### 3.2.3 LoRA

### 3.3 Quantization and EoRA

### 3.4 Evaluation criteria



# Chapter 4

## Implementation

nada

### 4.1 Implementation of the Pipeline

#### 4.1.1 Depth-wise Pruning

#### 4.1.2 Width-wise Pruning (WANDA-based)

ù

#### 4.1.3 LoRA

#### 4.1.4 Quantization and EoRA

### 4.2 Evaluation

How were the models evaluated?



# Chapter 5

## Experimental Results and Analysis

intro here

### 5.1 Preliminary observation of the results

Before examining the results based on In this section, we will present and comment on some examples of text generated by the

### 5.2 Results on TriviaQA

### 5.3 Results on WikiText





# Chapter 6

## Conclusion and Future Work

### 6.1 Future work

- speak about Distillation
- speak about experimenting with other quantization methods
- speak about adding support for other llms
- kV cache compression

### 6.2 Final remarks

The objective of this project was to research and implement a methodology for compressing LLaMA based LLMs,



# Bibliography

- [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin, *Attention is All You Need*, <https://arxiv.org/abs/1706.03762>.
- [2] Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova, *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*, <https://arxiv.org/abs/1810.04805>.
- [3] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, *Improving Language Understanding by Generative Pre-Training*, [https://cdn.openai.com/research-covers/language-unsupervised/language\\_understanding\\_paper.pdf](https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf).
- [4] OpenAI, *GPT-4 is OpenAI's most advanced system, producing safer and more useful responses*, <https://openai.com/index/gpt-4/>.
- [5] Mike Perkins, *Academic Integrity considerations of AI Large Language Models in the post-pandemic era: ChatGPT and beyond*, [https://www.researchgate.net/publication/368775737\\_Academic\\_integrity\\_considerations\\_of\\_AI\\_Large\\_Language\\_Models\\_in\\_the\\_post-pandemic\\_era\\_ChatGPT\\_and\\_beyond](https://www.researchgate.net/publication/368775737_Academic_integrity_considerations_of_AI_Large_Language_Models_in_the_post-pandemic_era_ChatGPT_and_beyond).
- [6] Daniel Mügge, *AI Is Threatening More Than Just Creative Jobs—It's Undermining Our Humanity*, <https://www.socialeurope.eu/ai-is-threatening-more-than-just-creative-jobs-its-undermining-our-humanity>.

- 
- [7] Arpan Suravi Prasad, Moritz Scherer, Francesco Conti, Davide Rossi, Alfio Di Mauro, Manuel Eggimann, Jorge Tomás Gómez, Ziyun Li, Syed Shakib Sarwar, Zhao Wang, Barbara De Salvo, Luca Benini *Siracusa: A 16 nm Heterogenous RISC-V SoC for Extended Reality with At-MRAM Neural Engine*, <https://arxiv.org/abs/2312.14750>.
  - [8] Robin M. Schmidt, *Recurrent Neural Networks (RNNs): A gentle Introduction and Overview*, <https://arxiv.org/pdf/1912.05911>
  - [9] Chris Nicholson, *A Beginner's Guide to LSTMs and Recurrent Neural Networks*. <https://skymind.ai/wiki/lstm>.
  - [10] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, et al., *PyTorch: An Imperative Style, High-Performance Deep Learning Library*, <https://arxiv.org/abs/1912.01703>.
  - [11] Angelo Galavotti, *FRANKEN-LLAMA code repository*, <https://github.com/AngeloGalav/franken-llama>
  - [12] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, Yejin Choi *HellaSwag: Can a Machine Really Finish Your Sentence?* <https://arxiv.org/abs/1905.07830>.
  - [13] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, et al., *LLaMA: Open and Efficient Foundation Language Models*, <https://arxiv.org/abs/2302.13971>.
  - [14] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, et al., *Llama 2: Open Foundation and Fine-Tuned Chat Models*, <https://arxiv.org/abs/2307.09288>.
  - [15] Llama Team, AI @ Meta, *The Llama 3 Herd of Models*, <https://arxiv.org/abs/2407.21783>.

- 
- [16] Hanjuan Huang, Hao-Jia Song, Hsing-Kuo Pao, *Large Language Model Pruning*, <https://arxiv.org/abs/2406.00030>.
- [17] Bo-Kyeong Kim, Geonmin Kim, Tae-Ho Kim, Thibault Castells, Shinkook Choi, Junho Shin, Hyoungh-Kyu Song, *Shortened LLaMA: Depth Pruning for Large Language Models with Comparison of Retraining Methods*, <https://arxiv.org/abs/2402.02834>.
- [18] Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, Danqi Chen *Sheared LLaMA: Accelerating Language Model Pre-training via Structured Pruning*, <https://arxiv.org/abs/2310.06694>.
- [19] ModelCloud, *GPTQModel*, <https://github.com/ModelCloud/GPTQModel>.
- [20] Elias Frantar, Saleh Ashkboos, Torsten Hoefer, Dan Alistarh *GPTQ: Accurate Post-Training Quantization for Generative Pre-trained Transformers*, <https://arxiv.org/abs/2210.17323>
- [21] Mingjie Sun, Zhuang Liu, Anna Bair, J. Zico Kolter, *A Simple and Effective Pruning Approach for Large Language Models*, <https://arxiv.org/abs/2306.11695>.
- [22] Shih-Yang Liu, Maksim Khadkevich, Nai Chit Fung, Charbel Sakr, Chao-Han Huck Yang, Chien-Yi Wang, Saurav Muralidharan, Hongxu Yin, Kwang-Ting Cheng, et al., *EoRA: Fine-tuning-free Compensation for Compressed LLM with Eigenspace Low-Rank Approximation*, <https://arxiv.org/abs/2410.21271>.
- [23] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, *LoRA: Low-Rank Adaptation of Large Language Models*, <https://arxiv.org/abs/2106.09685>.
- [24] Mandar Joshi, Eunsol Choi, Daniel S. Weld, Luke Zettlemoyer, *TriviaQA: A Large Scale Distantly Supervised Challenge Dataset for Reading Comprehension*, <https://arxiv.org/abs/1705.03551>.

- 
- [25] Stephen Merity, Caiming Xiong, James Bradbury, Richard Socher, *Pointer Sentinel Mixture Models*, <https://arxiv.org/abs/1609.07843>.
  - [26] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J. Liu, *Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer*, <https://arxiv.org/pdf/1910.10683>.
  - [27] Chen Liang, Haoming Jiang, Zheng Li, Xianfeng Tang, Bin Yin, Tuo Zhao, *HomoDistil: Homotopic Task-Agnostic Distillation of Pre-trained Transformers*, <https://arxiv.org/abs/2302.09632>.

# Acknowledgements

Ringrazio il professor Ozalp Babaoglu per la sua disponibilità come relatore, nonostante fosse prossimo alla pensione. Sono molto grato di essere uno dei suoi ultimi tesisti della sua carriera. I suoi contributi nell'ambito dei sistemi operativi sono inestimabili.

Ringrazio il professor Francesco Giacomini, per avermi dato l'accesso a una delle esperienze formative più importanti della mia vita. Lo ringrazio inoltre per avermi assistito attentamente nello sviluppo del progetto e nella scrittura di questa tesi. Grazie anche per la tua pazienza nei miei confronti.

Ringrazio tutto il personale dell'INFN CNAF per aver reso l'esperienza ancora più gradevole e per avermi fatto visitare il centro di calcolo, trattandomi sempre come se fossi un loro collega.

Un ringraziamento speciale va alla mia famiglia, per avermi sempre dato la spinta di andare avanti e per credere in me, senza negarmi mai nulla.

Ringrazio tutti gli amici che ho conosciuto nel corso durante questi tre anni: Leon, Drif, Giaco, Baldo, Adriano, Donnoh, Vir, Matteo, Pino, Denis, Samuele, Alice e tanti altri. Non avrei potuto chiedere dei compagni migliori.

Infine, ringrazio Leti, per essermi stata accanto nei momenti più bui e avermi dato tutto l'affetto che esiste in questo mondo, anche quando non me lo meritavo.